

Ph129
Lecture Notes

Statistics for Physicists - Eddis, Drijard, Ross

Advanced Theory of Statistics - Kendall & Stewart

Handbook of the Poisson Distribution, Haight

Selected Papers on Noise and Stochastic Processes - various authors -
including famous S.O. Rice papers on noise
and several RMP papers on statistical physics

Ailey - On the theory of stochastic processes and their application to the theory
of cosmic radiation

Fuller in "Comment - Models and Essays" [510 + 589] p. 105
On probability problems in the theory of counters.

Stochastic Processes - Glazier

Decision Theory
 Student F test
 Goodness of Fit
 Distribution independent tests

Generation of Monte (Pseudo) Random Numbers
 Monte Carlo Integrator

Things I don't understand very well)

Minimal Spanning Tree: Investigates
 Clustering

Generation of Optimal Decision Algorithms

(binary trees, Wind's magnetic field

Noise
 Nonparametric representations of data

Topics requested Sam.

Very little on how to run at
 Las Vegas e.g. combinatoric analysis.

Course

Either Homework ($> 50\%$) + oral
 or Prepare discussion of some
 advanced topic at end of course
 either 1 or student will present.

Books.

For the first part of course (the decision theory), choose between

Statistical and Computational Methods in Data Analysis by S. Brandt - North-Holland.

Statistical methods in Experimental Physics by W.T. Eadie et al., North-Holland.

Brandt's book has more examples (programs) and is simpler. (Bevington (old notes) is also good at a low level).

and for the mathematical background, see
Mathematical methods of Statistics
by H. Cramer (Princeton University Press)

Subjects:

Laws of Statistics

Basic Properties following from laws:

Bayes theorem

Continuous/Discrete random variables

Characteristic Functions

Sums of Random Variables

Probability Distributions

Central Limit Theorem

Probability Distributions.

Binomial

Poisson

Gaussian

Estimation of Parameters

maximum likelihood

χ^2

Method of moments

Bias / Error

Bayes v. Non Bayes

Robust (distribution free methods).

Minimization of Functions.

Phizics Logistics

Teacher G. Fox

Room 207 Booth

Please call x3765, x6673 if want to
ask questions.

Typically around 8.15am → 5.30pm

Graders : 1) J. Miller

x 2918

Room 108 Sloan Annex

2) M. Bucher

x 4861

60 w. Bridge.

They will alternate problem sets

D. Statistics

D1: Introduction

There are four main topics in statistics or probability theory.

- 1) "Combinatoric Brilliance": i.e. how to win at Las Vegas.
- 2) "Design of Experiments" - Histograms, variance, skewness etc. e.g. how to take reliable public opinion polls.
- 3) "Formal Theory": Lebesgue measure, central limit theorem etc.
- 4) "Parameter Estimation": χ^2 , maximum likelihood, method of moments.
- * Further there are some topics which are really techniques in numerical analysis
- 5) Find minimum/maximum of a function

of several variables.

6) Generation of random numbers and pseudo random variables : use of these in monte carlo evaluation of integrals.

most mathematical books cover 1) to 3). Typical is

Feller: An Introduction to Probability Theory and its applications : vols 1,2. This only covers 1-3 and is very good on 1.

Hogg: Introduction to Mathematical Statistics : is quite readable and covers 1) to 4).

Borngarten: Data Reduction and Error Analysis for the Physical Sciences: This covers 2, 4 and 5 . It is perhaps the best simple handbook of the

formulae. However it is naive on 4), 5) and 5) and if you start doing complicated things on χ^2 /maximum likelihood/minimization, Bevington will lead you astray. Another and even shorter general reference is Mathews and Walker: chapter 14: this mentions all the important points.

For 4), there are several notes by physicists -

- P.T. Solmitz, Ann. Rev. Nucl. Sci. 14, 375 (1964)
- J. Orear, UCRL-8417 (1958). "Notes on Statistics for Physicists"
- S. Yellin, DESY 72/12 : Preprint (1972)

K. McDonald CTSI Internal Report #57

Arndt and MacGregor Methods in Computational Physics - Vol 6.

Annis, Cheston and Primakoff, Rev. Mod. Phys. 25, 818 (1953)

white very good on (3) and formal aspects of (4) is Cramer: Mathematical Methods of Statistics.

There is also a recent book on more practical aspects of 4) i.e. a translation of Cramer into physics applications:
W.T. Eadie, D. Drijard, F.E. James, M. Roos, B. Sadoulet
 Statistical methods in Experimental Physics.

6) I will postpone till we do numerical analysis while there is no good reference on 5) - even though if one existed (and was read) it would save millions of dollars of computer time. Tolerable is:

Fox: (not me): Optimization methods for Engineering Design.

Bad (i.e. irrelevant mathematics) is
J.W. Daniel: The approximate minimization of functionals.

So the plan of course is:

D2 : Formal Points : definitions, Bayes Theorem, Gaussians, Central Limit Theorem.

D3 : Binomial and Poisson Distributions:

Las Vegas ... > Biology (Birth and Death Processes).

D4 Estimation of ~~some~~ Parameters

Maximum Likelihood, χ^2 , and method of moments, Goodness of fit, minimization techniques.

STATISTICS FOR PHYSICISTS

Geoffrey Fox

I: BASIC TECHNIQUES AND THEORY

- Formal Theory
- Special distributions including those needed for combinatorics
- Statistical Inference
 - Parameter Estimation
 - Decision Theory

H2

D2	Basic Definitions	S 88-3
	Bayes	S 88-7
	Continuous Distributions	S 88-13
	Random Variables	S 88-15
	Joint Distributions	S 88-16
	means, moments, Correlations	S 88-19
	Central Limit Theorem	S 88-32
	Monte Carlo Integration	S 88-38
	Extensions to Correlated Functions	S 88-41
D3	Binomial and Poisson Distributions	S 88-49
	Compound Poisson Distribution	S 88-55
	Additivity Property and Gaussian	S 88-59
	Limits of Binomial and Poisson Distributions	
	Birth and Death Processes	S 88-65

D4	Estimation of Parameters	S88-74
	Maximum Likelihood Method	
	Method and Heuristic Justification	S88-75
	Example	S88-86
	Proof	S88-91
	χ^2 Method	S88-98
	Relation to Maximum Likelihood	
	Example	S88-104
	Counting (Binned) Experiment	S88-108
	Event Experiment	S88-113
	Method of Moments	S88-119
D5	Minimization	S88-128
	Basic Derivative Method	
	Error Calculation	S88-130
	Real World	S88-132
	Eigenvalue and Marquardt	

Goodness of fit

χ^2	S88-148
Confidence level	S88-151
Bartlett's S function	S88-154
Hypothesis Testing	S88-157
Student's t distribution	S88-160
Robust estimation	S88-163

Generation of Random Numbers	R1
in $[0, 1]$ uniform	
General Distributions	R9

Useful Books for Statistics Part of Course

- G.P. Yost "Lectures on Probability and Statistics" - Copies handed out in class (contact G. Fox if you need further information)
- W.T. Eadie et al. "Statistical Methods in Experimental Physics" - A little Sophisticated for this course but technically excellent
(North Holland)
- W.H. Press et al. "Numerical Recipes - The Art of Scientific Computing" Chapter 7, 13, 14 - Computer Related Issues This has FORTRAN and PASCAL programs available
- J. Mathews and R.L. Walker (Benjamin) "Mathematical Methods of Physics" Chapter 14

D2 Formal Points

D2/1 Formal Definitions

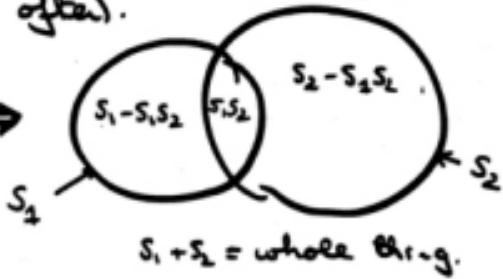
(i) A major difficulty in formulating the theory of probability is that one must relate it to the real world. Cramer (mathematical methods of statistics - Princeton) amusingly notes that some of the axiomatic approaches to probability are comparable to definitions of geometry in terms of the fundamental idea of a point as a limit of a spot of chalk dust of smaller and smaller size.

I will first dispose of my 24 lectures on formal probability by noting that the theory can be based on the following axioms.

Consider a variable point X in a space Y where Y is some subspace of \mathbb{R}^k .

(k dimensional real Euclidean Space). Further consider the family of all Borel sets S in \mathbb{Y} . You will of course recall that a Borel set is gotten by performing (finite or denumerably many times) on intervals $[a, b]$, (a, b) , $[a, b]$, $[a, b] \dots$ the operations of addition, subtraction or multiplication. (and continuing process denumerably often).

modern math \Rightarrow



Axiom 1: To every S corresponds a non-negative number $P(S)$ which is called the probability of the event X belonging to S .

Axiom 2: $P(\mathbb{Y}) = 1$

Axiom 3: $P(S)$ is a completely additive

Set function:

$$\text{i.e. } P(S_1 + S_2 + \dots) = P(S_1) + P(S_2) + \dots$$

If the S_i are disjoint ~~$S_i S_j = 0$~~ $S_i S_j = 0 \quad i \neq j$.

Note that it follows from the axioms that (a) $0 \leq P(S) \leq 1$: any S . Proof: $P(S) \geq 0$ was postulated. Take $S_1 = S$ $S_2 = Y - S = S^*$

(the complement of S) in axiom 3.

$$\text{but } P(S_1) + P(S_2) = P(Y) = 1 \Rightarrow P(S_2) = P(S) \leq 1.$$

(b) $P(S_1 + S_2) \leq P(S_1) + P(S_2)$, any S_1, S_2 with equality if S_1 and S_2 are disjoint.

$$\text{Proof } S_1 + S_2 = (S_1 - S_2 S_2) + S_2 \quad \text{as}$$

$P(S_1 + S_2) = P(S_1 - S_2 S_2) + P(S_2 - S_2)$ as sets disjoint. But $P(S_1 - S_2 S_2) \leq P(S_1)$ - See (a). This proves desired result.

Anyhow the variable point X is now called a random variable and $P(S)$ is the probability function of X .

w) To relate this definition to the outside world, we must resort to the vanishing chalk dust approach as espoused by Mathews and Walker.

If we flip an (unbiased) coin N times we expect $\lim_{N \rightarrow \infty} \frac{N_h}{N} = \frac{1}{2}$ where N_h is the number of times it turns up heads. we define:

$$P(\text{heads}) = \lim_{N \rightarrow \infty} \frac{N_h}{N} = \frac{1}{2}$$

This is a random variable with only two possible values - ~~say~~ heads or tails. Alternatively we could say $x=0$ and 1 . Then in our previous language any Boxd set containing 0 and not 1 has $P(S)=\frac{1}{2}$, any containing 0 and 1, $P(S)=1$ etc.

27/2 Bayes Theorem: Conditional Probabilities

Take two possible (Borel) subsets of $Y \in \mathbb{R}^n$. Call them A and B.

A could be given $x=0$
 B $x=1$ only.



$(A+B)^c$

Generally Y is divided into four subspaces by A and B and following chalk dust method - imagine N experiments, give the following frequencies.

$$N_1 = n(A-AB)$$

$$N_2 = n(B-AB)$$

$$N_3 = n(AB)$$

$$N_4 = n((A+B)^c) \text{ - complement of } (A+B).$$

of events $\sim (\star)$ falling in subspace \star .

$$\text{then } N_1 + N_2 + N_3 + N_4 = N. \quad \text{we}$$

previously defines $P(\star)$ and now we expect

$$\begin{aligned} P(A) &= \lim_{N \rightarrow \infty} \frac{N_1 + N_2}{N} \\ P(B) &= \lim_{N \rightarrow \infty} \frac{N_2 + N_3}{N} \\ P(A+B) &= \lim_{N \rightarrow \infty} \frac{N_1 + N_2 + N_3}{N} \end{aligned} \quad (D2.1)$$

NOW we define conditional probability $p(S/A)$ { read as probability of S given A has occurred } as $p(SA)/p(A)$. Note this is a probability according to axioms if we restrict Y to n.

Then in terms of frequencies.

$$p(A/B) = \frac{N_3}{N_2 + N_3} \quad \text{as } N \rightarrow \infty \quad (D2.2)$$

$$p(B/A) = \frac{N_3}{N_1 + N_3} \quad \text{as } N \rightarrow \infty$$

whence:

$$P(B) P(A/B) = P(A) P(B/A) \quad (D2.3)$$

which is Bayes law. of course (D2.3) is obvious from defⁿ, we didn't need frequencies to prove it.

n_{tot} pictures with some property of interest

$n(A)$ Number found by an undergraduate during finals week

$n(B)$ Number found by a graduate student working on a theory in 26 dimensions

example (page 11 Eadie et al.)

A: Scan Number One

B: Scan Number Two

We can tag members of $C = AB$ (A and B).

$p(C) = p(A)p(B)$ if Scans A, B independent
 (In practice not true as say short tracks systematically missed!).

$$\therefore \frac{n(C)}{n_{\text{tot}}} = \frac{n(A)}{n_{\text{tot}}} \cdot \frac{n(B)}{n_{\text{tot}}}$$

$$\text{Or } n_{\text{tot}} = \frac{n(A)n(B)}{n(C)} \text{ is desired}$$

estimate of total number of events

In practice second scan will only be on a subset of data sample.

Note if A, B independent

$$p(A/B) = \frac{p(AB)}{p(B)} = p(A) \text{ independent of } B.$$

Example: page 18 of Yost

$$B \left\{ \begin{array}{ll} B_1 = 0 & \text{ill} \\ B_2 = 1 & \text{well} \end{array} \right.$$

A = "Test Positive"

$$P(B/A) = \frac{P(B)P(A/B)}{P(A)}$$

Apply Separately $B = B_1 \text{ or } B_2$

$$P(B/A) = \frac{P(B_2)P(A/B_2) + P(B_1)P(A/B_1)}{P(A) = .001x.99 + .02x.999}$$

$$P(B_1/A) = \text{"ill"} = .047$$

$$P(B_2/A) = \text{"well"} = .953$$

The Bayes Controversy

Now we have derived Bayes theorem where A and B are (Borel) sets in a space populated by a random variable. Now a useful "applicator" of the theorem is take one of them

A : to be Borel set given a set of experimental results - this is correct.

B : to be a set of hypotheses e.g. the set in space of values of a theoretical parameter

e.g. we measure Scattering angle θ and wish to find α when probability of observing θ is $\frac{1}{2}(1 + \alpha \cos \theta)$.

A is a set in θ space.

B is a set in α space.

Then experimenters measure $P(A/B)$ i.e. values of θ subject to known (by the almighty and those on fourth floor) single value of α .

Bayes says necessarily

$$\text{i.e. } P(B/A) = P(A/B) \frac{P(B)}{P(A)}$$

and so this gives you ~~you~~ what you really want - the "probability distribution of α ". Consequent from measurement of θ . Problem with this is that:

a) α is not a random variable - it only takes one value. answer: Our knowledge of α can be considered as a random variable.

- b) what is $p(B)$? We will discuss this so called a prior knowledge of B later on - it essentially summarizes all previous (other) knowledge of B gained from previous (other) experiments.
- c) what is $p(A)$? The probability of measurement of Θ given any possible value of α ? Nobody knows what this is! You essentially determine it by normalizing $p(B_1/A)$ or equivalently note that relative probabilities

$$\frac{P(B_1/A)}{P(B_2/A)} = \frac{p(A/B_1)}{p(A/B_2)} \frac{p(B_1)}{p(B_2)}$$

$p(A)$ are indeed independent of the noxines

D2/3 Continuous Distributions

To be formal and reveal relics of a primeval mathematical training: take random variables in \mathbb{R}_1 and consider special Borel set $\{x : \text{it has probability } F(x)\}$ with $F(\infty) = 1$.

Now in some cases we can differentiate $F(x)$ and so write:

$$F(x) = \int_{-\infty}^x p(x) dx$$

where $p(x)$ is probability density. One can clearly interpret $p(x) dx$ as the probability that random variable X takes values in range x to $x+dx$.

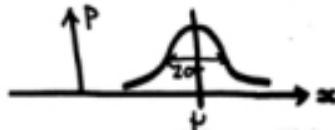
(a) In the case of tossed coin, let heads be $x=1$ and tails $x=0$.

Then (for an unbiased) coin

$$p(x) = \frac{1}{2} [\delta(x) + \delta(x-1)]$$

$$(a) p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left[-\frac{1}{2}\frac{(x-\mu)^2}{\sigma^2}\right] \quad (B2.4)$$

This is the probability density of the normal or Gaussian distribution - so called because it was first discussed by de Moivre in 1733.



Note that we have ensured that

$$\int_{-\infty}^{+\infty} p(x) dx = 1$$

Also the associated distribution function is:

$$F(x) = \bar{F}(x) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^x dy \exp\left[-\frac{1}{2}\frac{(y-\mu)^2}{\sigma^2}\right]$$

which cannot be done analytically except by relating it to the error integral:

$$\text{erf}(x) = \frac{2}{\sqrt{\pi}} \int_0^x \exp(-y^2) dy$$

D2/4: Functions of a Random Variable

If x is a random variable, then so is $y = f(x)$. Let x have probability distribution $p(x)$ and $y = q(y)$. Then probability of x lying in $[x, x+dx]$ is $p(x)dx$: let $[y, y+dy]$ be associated interval in y .

$$\text{Clearly } p(x)dx = q(y)dy$$

$$\text{and } \frac{dy}{dx} = |f'(x)|$$

$$\text{So } p(x) = q(f(x)) |f'(x)| \quad (\text{D2.5})$$

There are some obvious modifications of this if the $x \leftrightarrow y$ relation is not single-valued.

22/5: Joint Probability Distributions

We can now consider the important extension of (22/5.4) to more than one variable. Suppose we have two random variables x and y each taking values in \mathbb{R}_2 . Then we can conceive of a random variable z which is the ordered pair (x, y) . Even if we know the probability distributions of x and y ,

$$\begin{aligned} x &\rightarrow p(x) dx \\ y &\rightarrow q(y) dy \end{aligned}$$

The distribution of z is of course not defined. Suppose that z has density $r(x, y)$ so that probability that z lies in volume element $dx dy$ is $r(x, y) dx dy$. A priori, all we know is that r satisfies

$$\left. \begin{aligned} \int r(x, y) dy &= p(x) \\ \int r(x, y) dx &= q(y) \end{aligned} \right\} \text{These are called marginal distributions (EDJRS p. 18)}$$

(a) Example: Bayes law for densities.

Let event X be x lies in $[x, x+dx]$
 " " Y " y " " $[y, y+dy]$

Let $T(x|y) = t(x|y) dx$ be probability
 of x given y . By definition

$$T(x|y) = \frac{\Pr(XY)}{\Pr(Y)} = \frac{T(x,y) dx dy}{q(y) dy}$$

or $t(x|y) = r(x,y)/q(y)$ a conditional
 distribution (EDJRS
 page. 18)

Similarly let $S(Y|X) = s(y|x) dy$ be
 probability of Y given X . By definition:

$$s(y|x) = r(x,y)/p(x)$$

or eliminating $r(x,y)$

$$t(x|y) q(y) = s(y|x) p(x) \quad (D2.6)$$

This is a density version of Bayes law.

(b) Example: $y = f(x)$.

When y is a known function of x ,

$$\text{clearly } s(y|x) = \delta(y - f(x))$$

Then using above equ. for $S(y|x)$, we can in this case find

$$r(x, y) = \delta[y - f(x)] p(x).$$

(i) Independence:

If the distribution of z satisfies

$$r(x, y) = p(x) q(y) \quad (D2.7)$$

then we say that x and y are statistically independent.

(ii) Function of 2 variables

Suppose rather than z in \mathbb{R}_2 , we consider $h(x, y)$ in \mathbb{R}_2 : a ^{conditional} function of x and y . Then $w = h(x, y)$ is also a random variable with distribution $d(w)$ where

$$d(w) = \iint dx dy \delta[w - h(x, y)] r(x, y) \quad (D2.8)$$

This may perhaps be proved very
grettly: it is obvious if you just do
one-say x -integral in (D2.8) ^{by} and
use same argument in D2/4.

D2/6 Means - moments etc

(a) Given a one dimensional distribution:

$$\text{mean: } \langle x \rangle = \int_{-\infty}^{+\infty} x p(x) dx$$

$$\text{n'th moment: } \langle x^n \rangle = \int_{-\infty}^{+\infty} x^n p(x) dx$$

Expectation: $\langle f(x) \rangle = \int_{-\infty}^{+\infty} f(x) p(x) dx$ which "follows" from n'th moment by Taylor expansion of $f(x)$.

This can be considered as mean of the random variable $Z = f(x)$.

n'th central moment:

$$\langle (x - \langle x \rangle)^n \rangle = \int_{-\infty}^{+\infty} (x - \langle x \rangle)^n p(x) dx$$

Variance: $V_x = \langle (x - \langle x \rangle)^2 \rangle = 2^{\text{nd}} \text{ central moment}$

Standard deviation σ_x : $V_x = \sigma_x^2$ where

σ_x is used much more than V_x .

$$\begin{aligned} \text{Note } V_x &= \langle x^2 \rangle - 2 \langle x \rangle \langle x \rangle + \langle x \rangle^2 \\ &= \underline{\langle x^2 \rangle} - \underline{\langle x \rangle^2} \end{aligned} \quad (B2.9)$$

a result which is used continuously.

(b) n-dimensional distribution:

The mean now becomes an n-dimensional vector

$$\langle x_R \rangle = \int \otimes dx_i \quad x_R p(x_1 \dots x_n)$$

The variance becomes an $n \times n$ moment matrix M

$$M_{R\ell} = \int \otimes dx_i \quad [x_R - \langle x_R \rangle] [x_\ell - \langle x_\ell \rangle] \\ \cdot p(x_1 \dots x_n)$$

Note for we can show - analogously to (02.9) - that

$$M_{R\ell} = \langle x_R x_\ell \rangle - \langle x_R \rangle \langle x_\ell \rangle \quad (02.10)$$

Also note that M is symmetric and positive semi definite. [This follows because $\forall R, \ell, y_R^T M^{-1} y_\ell - \text{any } y - \text{ clearly has a } g(y) \geq 0 \text{ integrand}].$

Now the diagonal terms are called variances which are ≥ 0 : their

Square roots are standard deviations again. Further

$$k \neq l : \rho_{kl} = M_{kl} / \sqrt{M_{kk} M_{ll}} \text{ is } \leq 1$$

and is called correlation coefficient.

Note that if $\rho_{kl} = 0$ - all $k \neq l$ - then we say that $x_1 \dots x_n$ are uncorrelated. Further if $x_1 \dots x_n$ are independent, then (D2.10) shows at once that they are uncorrelated: the converse does not follow of course. (except under special circumstances e.g. x_i Gaussianly distributed).

One can define higher order terms $M_{klmn\dots}$ but these are not too important as in usual applications $x_1 \dots x_n$ are \sim functions of $N \gg n$ observations.

The central limit theorem - two which

we will come - implies that only the mean and $M_{R\bar{R}}$ are important asymptotically (in N).

(c) use of moment matrix:

Often the results of an experiment on $x_1 \dots x_n$ are quoted in terms of means \pm errors : the errors are standard deviations = $\sqrt{\text{diagonal elements of moment matrix}}$. Now this is throwing away information - the off diagonal terms in M - which can be important. Thus suppose some years after the experiment, a theorist decides that $y = T \bar{x}$ (T an $n \times n$ matrix) is useful. What are errors on y ?

$$\langle y_R y_L \rangle = T_{R\bar{R}} \langle x_R x_L \rangle T_{L\bar{L}}^T$$

$$\text{or } M_y = T M_x T^T \quad (02.11)$$

To find diagonal terms in $\frac{M^2}{2} - v$.
 the errors² on the matrix y — requires
 off diagonal terms in M_x . An example
 from of high energy physics is in the
 decay of a Spin 1 particle, one can
 define angular distribution by
 density matrix elements.

$$x_1 = g_{00}$$

$$x_2 = g_{1-1}$$

$$x_3 = g_{10}$$

in some Lorentz frame? Some years ago
 these were always given in one frame
 — the so called t channel frame. Recently
 it was found that the values in another
 frame (the "s channel") were more
 meaningful physically. Unfortunately
 latter values y were related — by
 precisely such a linear transformation
 $y = T x$ and one can no longer
 estimate errors.

The x_i and y_i in this example correspond to expressing angular decay

$$W(\theta, \phi) = 1 + \sum_{i=1}^3 x_i X_i(\theta, \phi)$$

$$\text{or } W(\theta, \phi) = 1 + \sum_{i=1}^3 y_i Y_i(\theta, \phi)$$

where y_i are linear combinations of the x_i corresponding to the

$X_i(\theta, \phi)$ being linear combinations of the $Y_i(\theta, \phi)$.

The details are different, but you can see the essential points from

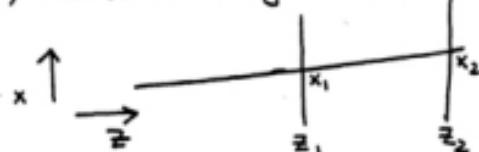
$$X_1 = \cos \theta \quad Y_1 = P_1(\cos \theta)$$

$$X_2 = \cos^2 \theta \quad Y_2 = P_2(\cos \theta)$$

$$X_3 = \cos^3 \theta \quad Y_3 = P_3(\cos \theta)$$

More Examples of Correlations

i) measurement of Tracks with Chambers



$$\text{Track is } x = a + b z$$

↑ ↑
 intercept slope

$$\text{where } b = \frac{(x_2 - x_1)}{(z_2 - z_1)}$$

$$a = \frac{z_2 x_1 - z_1 x_2}{z_2 - z_1}$$

x_i are independent random variables
but a and b are clearly correlated

$$\langle (b - \langle b \rangle)(a - \langle a \rangle) \rangle = - \frac{z_2 \sigma_1^2 - z_1 \sigma_2^2}{(z_2 - z_1)^2}$$

where x_i has standard deviation σ_i

(Note if $\sigma_1 = \sigma_2$ and $z_1 = -z_2$ i.e.

origin z midpoint between chambers) then
 a and b are uncorrelated) However

the origin $z=0$ (in our experiment) was
a convenient surveyors mark and so

was not chosen to minimize correlation!
 In the case of many chambers the slope and intercept will still be correlated but there will be no simple way to get rid ~~of~~ it. In practice it is important to keep full correlation matrix e.g. best vertex



needs it.

- 2) measurement of Regge pole intercept
 in E350,

$$\frac{dx}{dt dx} = g(t) (1-x)^{1-2\alpha(t)}$$

where $x \leq 1$ ($-5 \leq x \leq 1$ in experiment)

$g(t)$ is unknown fixed (at fixed t) constant of little interest.

$\alpha(t)$ is interesting.

In a (non linear) fit

$g(t)$ and $\alpha(t)$ are strongly correlated



so that $\int_{.5}^1 g(t) (1-x)^{1-2\alpha(t)}$ is

approximately preserved. Unfortunately we only want $\alpha(t)$ (predicted by theory) and cannot use fact that a particular linear combination of $\alpha(t)$ and $g(t)$ is predicted much better than either one.

In practice - when off diagonal terms are not given - one assumes they are zero - i.e. variables are uncorrelated. This gives:

$$\sigma_{y_i}^2 = \sum_{j=1}^n (T_{ij})^2 \sigma_{x_j}^2 \quad (D2.12)$$

which is the usual ~~rule~~ rule of addition of errors in quadrature. One has to emphasize that this rule is only valid where x_j uncorrelated.

(d) Error in lack of Correlation Assumption

$$\text{If } y = a_1 x_1 + a_2 x_2$$

Then $\sigma_y^2 = a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2 + 2a_1 a_2 \sigma_1 \sigma_2 \rho_{12}$
and all one can say in general is that:

$$\sigma_y^2 \leq (|a_1 \sigma_1| + |a_2 \sigma_2|)^2$$

The uncorrelated estimate is

$$\tilde{\sigma}_y = \sqrt{a_1^2 \sigma_1^2 + a_2^2 \sigma_2^2}.$$

and the true σ_y satisfies

$$0 \leq \sigma_y \leq \sqrt{2} \tilde{\sigma}_y$$

where extremes are only attained when $|a_1\sigma_1| = |a_2\sigma_2|$ which is most dangerous case. Other ratios $|a_1\sigma_1|/|a_2\sigma_2| \neq 1$ are better e.g. $a_1\sigma_1 = 0 \Rightarrow \sigma_y = \tilde{\sigma}_y$!

For n variables, all one can say is

$$0 \leq \sigma_y \leq \sqrt{n} \tilde{\sigma}_y.$$

which is clearly pretty bad for large n .

(2) Parameters of a Gaussian Distribution

Suppose x is Gaussian i.e.

$$p(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{(x-\mu)^2}{2\sigma^2}\right\}$$

Then

$$\langle x \rangle = \int_{-\infty}^{+\infty} x p(x) dx$$

$$= \mu + \int_{-\infty}^{+\infty} (x - \mu) p(x) dx = \mu$$

as last integral clearly vanishes from antisymmetry. So μ is mean.

$$\text{Variance} = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} (x-\mu)^2 \rho(x) dx$$

$y = x - \mu$
multiply
and divide
by σ^2 .

$$= \sigma^3 / \sqrt{2\pi}\sigma \int_{-\infty}^{+\infty} dy \frac{y^2}{\sigma^2} e^{-y^2/2\sigma^2}$$

$$= \sigma^3 / \sqrt{2\pi}\sigma \frac{\partial}{\partial \sigma} \left\{ \int_{-\infty}^{+\infty} e^{-y^2/2\sigma^2} dy \right\}$$

$$= \sigma^2 / \sqrt{2\pi} \frac{\partial}{\partial \sigma} [\sqrt{2\pi}\sigma]$$

$$= \sigma^2$$

so σ is Standard deviation - as its nomenclature suggested.

(f) Additive of Independent Random Variables

Suppose $y = \frac{1}{n} \sum_{i=1}^n x_i$ where the x_i are independent and all have the same distribution with mean μ and standard deviation σ . Then clearly $\langle y \rangle$ is also μ .

$$\begin{aligned}\langle (y-\mu)^2 \rangle &= \frac{1}{n^2} \sum_{i,j} \langle (x_i - \mu)(x_j - \mu) \rangle \\ &= \sigma^2/n \quad \text{as all off} \\ &\quad \text{diagonal terms vanish.}\end{aligned}$$

So exactly y has mean μ and std. deviation σ/\sqrt{n} . Later we will show that as $n \rightarrow \infty$ (i.e. not exactly) all other moments of y vanish faster than $1/n$ and so y becomes Gaussian.

D2/7 Central Limit Theorem

(i) Moment Generating Function

This which is also called characteristic function is defined by

$$\varphi(k) = \int_{-\infty}^{+\infty} dx e^{ikx} p(x)$$

$$= 1 + ik\langle x \rangle - \frac{k^2}{2!} \langle x^2 \rangle + \dots$$

(a) Note the Fourier inversion theorem

$$gives \quad p(x) = \frac{1}{2\pi} \int_{-\infty}^{+\infty} dk e^{-ikx} \varphi(k)$$

(b) most importantly put $h=x+y$
where x and y are independent. From
D2/5 it follows that distribution
density of h is

$$w(h) = \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} \delta(h-x-y) p(x) q(y) dx dy.$$

so

$$\varphi_h(k) = \int_{-\infty}^{+\infty} dh w(h) e^{ikh}$$

$$\begin{aligned}
 &= \int_{-\infty}^{+\infty} \int_{-\infty}^{+\infty} e^{ik(x+y)} p(x) q(y) dx dy \\
 &= \varphi_x(k) \varphi_y(k)
 \end{aligned}$$

So characteristic functions multiply for addition of two independent random variables. This result can be obviously be extended to sum of n independent random variables.

$$(2) \quad \varphi(k) = \sum_n k^n E((x)^n) / n!$$

\uparrow
expectation value

$$\varphi(k) = 1 + L \langle x \rangle k - \frac{\langle x^2 \rangle k^2}{2} \text{ etc.}$$

For a Gaussian distribution

$$\begin{aligned}
 \varphi(k) &= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{+\infty} \exp(ixk) \exp \left[-\frac{1}{2} \frac{(x-\mu)^2}{\sigma^2} \right] \\
 &= \exp \left\{ i\mu k - \frac{1}{2} k^2 \sigma^2 \right\}
 \end{aligned}$$

(d) Instead, define the Cumulant as 588-34

$$K(k) = \ln \varphi(k)$$

I have never used this!

Note that

$$K(k) = \frac{1}{2} \langle x \rangle ik - \underbrace{(\langle x^2 \rangle - \langle x \rangle^2) k^2}_{\sigma^2}$$

i.e. as $x \rightarrow x + p$ only shifts K by additive constant ikp , we only see "central" moments for coefficients of k^2 and above.

(ii) Central Limit Theorem

Let $y = \frac{1}{n} \sum_{i=1}^n x_i : x_i$ independent

(a) If x_i are Gaussian then so is y : in fact if y Gaussian exactly, then so must all the x_i be Gaussian.

(b) De Moivre in 1733 showed that if x_i binomially distributed, then y was always Gaussian for large n . (even though Gauss didn't exist in 1733).

(c) Laplace, in 1812, showed that de Moivre's result was generally true.

Cramer p.215. Central Limit Theorem of Laplace

If $x_1 \dots x_n$ are independent random variables with the same distribution, then y is asymptotically Gaussian with mean μ and standard deviation σ/\sqrt{n} (where x_i has mean μ and s.d. σ , and both μ and σ exist a.i.e. are finite)

Central Limit Theorem - Lindeberg

If $x_1 \dots x_n$ are independent random variables with $\mu_i = \langle x_i \rangle$

$$\sigma_i^2 = \langle (x_i - \mu_i)^2 \rangle$$

$$g_i^3 = \langle |x_i - \mu_i|^3 \rangle$$

where g_i must exist. (not necessary in Laplace's version where all x_i had same distribution). Let

$$g^3 = \sum_{i=1} g_i^3$$

$$\text{and } = \sum_{i=1}^n \sigma_i^2$$

Then if $\lim_{n \rightarrow \infty} S/\Sigma \rightarrow 0$, then \mathbf{y} is Gaussianly distributed with mean $\mu = \bar{x} = \sum_{i=1}^n x_i/n$ and s.d. Σ/n .

Note the $\lim_{n \rightarrow \infty} S/\Sigma$ does $\rightarrow 0$ when S ends in Laplace's case as $S \propto n^{1/3}$
and $\Sigma \propto n^{1/2}$

Last week we stated the Central Limit Theorem
 Now we must prove one version of it

Proof: (Laplace)

From our discussion of characteristic functions (c.f.)

$$\varphi_y = [\varphi_x(k/n)]^n$$

where φ_y / φ_x are c.f. of y and x respectively.

$$\begin{aligned} \text{Now } \varphi_x(k) &= \exp(ik\mu) \int dx e^{ik(x-\mu)} p(x) \\ &= \exp(ik\mu) \left[1 - \frac{k^2 \sigma^2}{2} + O(k^3) \right] \end{aligned}$$

$$\begin{aligned}
 \text{So } \Psi_y &= \exp(iky) \left[1 - \frac{k^2 \sigma^2}{2n^2} + O\left(\frac{(k/n)^3}{n}\right) \right]^n \\
 &= e^{iky} \left\{ \exp \left[-\frac{k^2 \sigma^2}{2n^2} + O\left(\frac{(k/n)^3}{n}\right) \right] \right\}^n \\
 &= \exp(iky) \exp \left\{ -\frac{k^2 \sigma^2}{2n} + O\left(\frac{k^3}{n^2}\right) \right\} \\
 &\quad \downarrow \\
 \text{and Fourier transform of a Gaussian,} \\
 \text{is well known to be a Gaussian, so} \\
 p(y) &= \frac{1}{\sqrt{2\pi}} \frac{1}{\sigma} \exp \left\{ -\frac{(y-\mu)^2}{2\sigma^2} \right\}
 \end{aligned}$$

and y is Gaussianly distributed as claimed.

(iii) Monte-Carlo Integration

lets consider this important numerical technique from a formal point of view. we want to evaluate

$$I = \int_a^b f(x) dx$$

we can surely write:

$$I = \int_{a-\infty}^{a+\infty} f(x) p(x) dx$$

where $p(x) = \frac{1}{(b-a)} \quad a \leq x \leq b$
 $= 0 \quad \text{outside } [a, b].$

$$\text{let } \bar{I} = \frac{1}{n} \sum_{i=1}^n f(\tilde{x}_i)$$

where \tilde{x}_i are distributed according
 to $p(x)$. Now apply central limit theorem
 to $\tilde{x}_i = f(x_i)$.

$$\tilde{x}_i \text{ has mean } \mu = \int_a^b f(x) p(x) dx$$

$$\text{and } \sigma^2 = \int_a^b f^2(x) p(x) dx - \mu^2$$

Then according to theorem

$$\bar{I} \rightarrow I \pm \sigma/\sqrt{n}$$

or I can be estimated as \bar{I}
 with error σ/\sqrt{n} . Clearly a good
 monte carlo evaluation requires small σ .
 (i.e. Small "variance" of f). We will
 return to this in numerical methods.

Note that in any Monte Carlo calculation one should calculate error as well as value of integral.

$$\sigma \text{ involves } \int_a^b f^2(x) p(x) dx \\ = \langle f^2(x) \rangle$$

which is estimated as

$$\Rightarrow \frac{1}{n} \sum_{i=1}^n f^2(x_i)$$

So one should accumulate $\sum f^2$ as well as $\sum f$.

(iv) Extension of Central Limit Theorem

An important extension of the central limit theorem (necessary when discussing moments) is

Theorem

Let x_i be independent random variables with the same distribution and let $f(x)$, $g(x)$ be any functions over this distribution. Suppose

$$\mu_f = \langle f \rangle$$

$$\mu_g = \langle g \rangle$$

$$\sigma_f^2 = \langle (f - \mu_f)^2 \rangle$$

$$\sigma_{fg} = \langle (f - \mu_f)(g - \mu_g) \rangle$$

$$\sigma_g^2 = \langle (g - \mu_g)^2 \rangle$$

$$M = \begin{bmatrix} \sigma_f^2 & \sigma_{fg} \\ \sigma_{fg} & \sigma_g^2 \end{bmatrix} \text{: moment matrix}$$

Finally set $F = \frac{1}{n} \sum_{i=1}^n f(x_i)$

$$G = \frac{1}{n} \sum_{i=1}^n g(x_i)$$

Then (F, g) is asymptotically a (2-dimensional) Gaussian distribution with

$$\hat{p}(F, g) \propto \exp \left\{ -\frac{n}{2} (y-\mu)^T B(y-\mu) \right\}$$

where we have lapsed into vector notation : $y = \begin{bmatrix} F \\ g \end{bmatrix}$, $\mu = \begin{bmatrix} \mu_F \\ \mu_g \end{bmatrix}$ and $B^{-1} = n$.

Proof: just as in ordinary central limit theorem except you consider 2nd dimensional characteristic function $\varphi(k_1, k_2)$.

Corollary : if F, g are defined as above, then F/g is asymptotically Gaussian with mean μ_F/μ_g and Standard deviation $\sqrt{\frac{1}{n} \langle (f_{F,g} - g_{F,g})^2 \rangle}$

Examples (a) $g=1$, mean = μ_F , $\sigma = \sqrt{\frac{1}{n} \sigma_F^2}$
 i.e. ordinary central limit theorem
 (b) $F=g$, mean = 1, $\sigma = 0$.

This extension is important as it allows you to use correlations between functions calculated from the same data (observations).

Example

1) Suppose we select a sample of events from some set e.g. the sample maybe those falling into some bin of a histogram. We wish to estimate fraction of events with ^{given} property.

Define random variable Θ by $\Theta = 1$ event passes selection criterion
 $= 0$ event fails criterion.

Then $f = \Theta$ $\langle f \rangle = p$, say.
 $g = 1$. $\langle g \rangle = 1$

The estimate of fraction

is just
$$\frac{\sum_{i=1}^N \Theta_i}{\sum_{i=1}^N 1}$$

Then error in estimate is

$$\sigma = \frac{P}{\sqrt{N}} \sqrt{\langle (f/p - 1)^2 \rangle}$$

$$= \frac{P}{\sqrt{N}} \sqrt{\langle (f^2/p^2 - 2f/p + 1) \rangle}$$

$$= \sqrt{\frac{P(1-p)}{N}} \quad \text{as } f^2 = f.$$

This is (familiar) formula for Standard deviation of a binomial distribution..

The naive estimate is

$$\langle f \rangle = N \sum_{i=1}^n f(x_i)$$

$$\bullet \langle f^2 \rangle = N \sum_{i=1}^n f^2(x_i) = \langle f \rangle$$

$$\sigma = \langle f^2 \rangle - \langle f \rangle^2 = p - p^2$$

Error is σ/\sqrt{N} in agreement with above.

i) Another example:

Suppose a scattering experiment observes θ with probability $p(\theta)$

Then typically, one would expand

$$p(\theta) = \sum_{n=0}^{\infty} a_n P_n(\cos \theta) \text{ as a series}$$

of Legendre Polynomials

$$\begin{aligned} \text{Then } \langle P_{n_1}(\cos \theta) \rangle &= \int P_{n_1}(\cos \theta) p(\theta) d\theta \\ &= a_{n_1} \frac{2}{2n_1 + 1} \end{aligned}$$

$\therefore \langle P_{n_1}(\cos \theta) \rangle$ need to be calculated
as averages we observed θ_k

$$\langle P_{n_1}(\cos \theta) \rangle = \frac{1}{N} \sum_{k=1}^N P_{n_1}(\cos \theta_k)$$

Even though $P_{n_1}(\cos \theta)$ and $P_{n_2}(\cos \theta)$
are orthogonal polynomials, the
means $\langle P_{n_1}(\cos \theta) \rangle$, $\langle P_{n_2}(\cos \theta) \rangle$
are not uncorrelated. Both

are calculated from same set
 $\{\theta_k\}$ and fluctuate in a correlated
 fashion as observations θ_k fluctuate.

$$f = P_{n_1}(\cos \theta)$$

$$g = P_{n_2}(\cos \theta)$$

$$F = \frac{1}{N} \sum_k P_{n_1}(\cos \theta_k)$$

$$G = \frac{1}{N} \sum_k P_{n_2}(\cos \theta_k)$$

F and G will be independent
 if $\int p(\theta) P_{n_1}(\cos \theta) P_{n_2}(\cos \theta) d\theta = 0$

This is not in general true.

A third example

- iii) Consider a histogram in which N events with weight w_i ($i=1 \dots N$) and x -value x_i fall.

$$\langle x \rangle = \frac{\sum_{i=1}^N w_i x_i = S_3}{\sum_{i=1}^N w_i = S_1} \quad \text{c.f. Monte Carlo calculation of } \frac{\int x w(x) dx}{\int w(x) dx}$$

with $f = w x$ $g = w$

the error in this calculation is:

$$\sigma = \frac{S_3}{S_1} \sqrt{\frac{1}{N} \left\langle \left(\frac{N w x}{S_3} - \frac{N w}{S_1} \right)^2 \right\rangle}$$

$$\sigma^2 = \frac{S_3^2}{S_1^2} N \left\langle \left[\frac{w^2 x^2}{S_3^2} - \frac{2 w^2 x}{S_3 S_1} + \frac{w^2}{S_1^2} \right] \right\rangle$$

$$= \frac{1}{S_1^2} [S_6 - 2 \langle x \rangle S_4 + \langle x^2 \rangle S_2]$$

$$S_2 = \sum_{i=1}^N w_i^2$$

$$S_4 = \sum_{i=1}^N w_i^2 x_i^2$$

$$S_6 = \sum_{i=1}^N w_i^2 x_i^2$$

D3 Jee förfun in Las Vegas

D3/1: Binomial and Poisson Distributions

We now come to change of space and consider more combinatorial aspects of probability. Much the nicest reference for this is Feller who has a host of beautiful examples.

(i) Bernoulli Trials

(a) A series of Bernoulli Trials is the generalized coin tossing problem. We have a sequence of independent "trials" - each of which has only two outcomes: Success or "heads": probability p
 Failure or "tails": probability q
 where of course $p+q=1$.

Stiffocrats of affirmative action and women's lib, will note that

generations of math books have downgraded tails compared to heads. Anyhow we follow in their footsteps and set $b(k, n, p)$ as the probability of k successes in n trials. This is easily calculated by writing out all possible outcomes as, for instance,

ht h tt ...
n

There are $\binom{n}{k}$ configurations with k "h"s and each has probability $p^k q^{n-k}$. So :

$$b(k, n, p) = \binom{n}{k} p^k q^{n-k} \quad (D3.1)$$

This is the binomial distribution because $b(k, n, p)$ are just terms in binomial expansion

$$1 = (p+q)^n = \sum_{k=0}^n b(k, n, p).$$

(b) Example:

A course consists of 50 lectures given to 20 people - all of whom get up before 9 a.m. - and attend all classes. A grand total of 1000 Smiles is observed. What is the chance that Frank Nagy (a physics graduate student in Kingsway) Smiled exactly once during whole course.

Sol^t: Use previous formalism, setting each trial as a smile; and Success as Frank Nagy Smiled.

A priori each smile has equal probability of belonging to each person

$$\text{so } p = \frac{1}{20} \quad q = \frac{19}{20}$$

$$\text{and } n = 1000.$$

Thus:

$$\Pr \{ \text{Frank Nagy Smiles once and once only} \} =$$

$$\left(\frac{1000}{20}\right) \cdot \frac{1}{20} \cdot \left(\frac{19}{20}\right)^{999}$$

$$\approx 50 e^{999 \log_{10} 19/20}$$

which isn't very
big

$$\approx \frac{50 e^{-51}}{10^{-19}}$$

(ii) Poisson limit:

(a) This is a very important special case of (D3.3) - the binomial distribution (D3.1) which describes anything from telephone calls to events observed in particle physics. In (D3.3), take limit of $n \rightarrow \infty$ and $p \rightarrow 0$ such that $np = \lambda$ remains finite. k is also fixed.

$$b(k, n, \lambda/n) = \frac{n!}{k!(n-k)!} \left(\frac{\lambda}{n}\right)^k \left(1-\frac{\lambda}{n}\right)^{n-k}$$

$$\text{now } \frac{n!}{(n-k)!} = n(n-1) \dots (n-k+1) \rightarrow n^k \text{ as}$$

$n \rightarrow \infty$ for fixed k . (You can also get this using Stirling's formula for the

Gamma function)

$$\text{Also } (1 - \lambda_{1/n})^{n-k} \approx (1 - \lambda_n)^k \rightarrow e^{-\lambda}$$

$$\begin{aligned} \text{so: } P(k, n, \lambda_n) &\rightarrow \frac{\lambda_n^k}{k!} \frac{\lambda_n^{n-k}}{(n-k)!} e^{-\lambda} \\ &= P(k, \lambda) = e^{-\lambda} \frac{\lambda^k}{k!} \quad (03.2) \end{aligned}$$

\Rightarrow The Poisson distribution.

$$\text{Note } \sum_{k=0}^{\infty} P(k, \lambda) = 1 \text{ as it should}$$

(A) Example: Given circumstances of previous example, what is chance that Frank Nagy smiles exactly once in a given lecture:

Solt: Again we set each trial = a smile and we have a 1000 of them. Now Success = "Frank Nagy Smiled this smile in our given lecture" and it is not very probable $P = \frac{1}{20.50} = \frac{1}{1000}$.

So n is large, p is small but
 $np = \lambda = 1$ is reasonable.

So the Poisson approximation to the exact (in this example) Binomial distribution is valid and chance is e^{-1} i.e. quite probable that F.N. Smiles once and once only in a given lecture.

The binomial estimate is

$$\left(\frac{1}{2}\right)^{\frac{1}{1000}} \left(\frac{999}{1000}\right)^{999} = \left[\frac{999}{1000}\right]^{999}$$

which is indeed $\approx e^{-1}$.

Compound Poisson Distribution

Let $y = \sum_{i=1}^N y_i$

y_i is Poisson with mean μ

N is Poisson with mean λ

$\Pr(y = \text{integer } r)$ is

$$\sum_{N=0}^{\infty} \Pr(N \text{ takes given value})$$

$$\Pr\left(\sum_{i=1}^N y_i, N \text{ fixed has } r \text{ events}\right)$$

as for fixed N , y is Poisson with mean $N\mu$, we find

$\Pr(y = \text{integer } r)$

$$= \sum_{N=0}^{\infty} \frac{\lambda^N e^{-\lambda}}{N!} \cdot \frac{(N\mu)^r e^{-N\mu}}{r!}$$

More generally we could give y_i a non Poisson distribution.

Examples are

i) ionizer - p 53 EDJRS

Particularly interesting is sensitivity to assumption that Poisson note that N variables with each with Note the Poisson distribution of 4 objects with $\langle \text{number of objects} \rangle = \lambda$ is much broader than Ordinary Poisson of near 4λ . This is of importance in particle physics where

ii) we believe that $q\bar{q}$ "clusters" are produced independently in a high energy collision. These clusters are Poissonly distributed and decay ~~into~~ $4 \rightarrow 6$ final particles.

Show $f_2 = \sigma^2 - \text{mean} = 0$ for Poisson
 > 0 for cluster model.

This is a compound process with y_i probably distributed Gaussianly

For example ; consider Poisson distribution of clusters -
each of which always contains 4 particles

For a Poisson distribution x

$$\text{Mean} = \lambda \quad \sigma^2 = \lambda$$

$$f_2(x) = \sigma_x^2 - \langle x \rangle = 0$$

Let random variable $y = 4x$

$$\langle y \rangle = 4\lambda$$

$$\langle y^2 \rangle = \langle (4x)^2 \rangle = 16 \langle x^2 \rangle$$

$$= 16 (\langle x^2 \rangle - \langle x \rangle^2 + \langle x \rangle)$$

$$= 16 (\sigma_x^2 + \langle x \rangle)$$

$$= 16 (\lambda + \lambda^2)$$

$$\begin{aligned} f_2(y) &= \langle y^2 \rangle - \langle y \rangle^2 - \langle y \rangle \\ &= 16(\lambda + \lambda^2) - 16\lambda^2 - 4\lambda \\ &= 12\lambda > 0 \end{aligned}$$

i.e., y is always greater than x

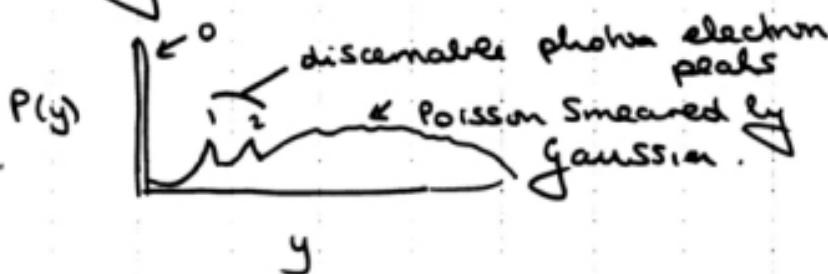
iii) Cherenkov light

Total pulse height observed

$$y = \sum_{i=1}^{N_\gamma} y_i$$

y_i = response in phototube of
a single photon - probably Gaussian
(i.e. limit of large Poisson)

N_γ = number of photons given off
in Cherenkov medium which is
certainly Poisson.



(ii) Additivity Property of Binomial or Poisson Distributions

These two distributions share an important additivity property which is a) obvious and b) leads to confusion when used in conjunction with certain

limit theorem.

(a) Suppose the independent random variables x_1 and x_2 have binomial distributions with parameters (n_1, p) and (n_2, p) respectively. Here $x_{ni} = k$ if $\exists k$ successes. Then it is obvious from source of binomial distribution that $x_1 + x_2$ has a binomial distribution with parameter $(n_1 + n_2, p)$.

(b) Suppose independent random variables y_1 and y_2 have Poisson distributions λ_1 and λ_2 . Then $y_1 + y_2$ has Poisson distribution with parameter $\lambda_1 + \lambda_2$. This is obvious either algebraically or by using (a) with $n_1 = \lambda_1/p$ and $n_2 = \lambda_2/p$. and let $n_i \rightarrow \infty$, $p \rightarrow 0$, fixed λ_1, λ_2 .

(iv) Gaussian limit of Binomial distribution

(a) M. and W. confused me by discussing two limits of binomial distribution.

We have already discussed one of these - the Poisson limit - which is really distinctive limit.

However we can also see a Gaussian limit which is just a special case of central limit theorem. Although this limit has a particular interpretation because of the additivity property we just discussed. Note ~~that~~ also that central limit theorem was discovered in general in 1812 whilst de Moivre discovered its binomial special case in 1733.

(b) Let x_i be independent random

variables - each of which has the same $b(k, 1, p)$ distribution. i.e. for each i : $x_i = 0$ is failure and $x_i = 1$ is success with relative probabilities $1-p$ and p respectively. Then $y = \sum_{i=1}^n x_i$ is in fact, (k), the number of successes in n Bernoulli trials. Then

(a) from additivity (iii): y is Binomial with distribution $b(k, n, p)$

(b) From Standard central limit theorem, y is Gaussian with mean np and standard deviation $\sqrt{n}\sigma$. Here μ and σ are mean and S.D. of $b(k, 1, p)$ distribution. i.e.

$$\mu = 0 \cdot (1-p) + 1 \cdot p = p$$

$$\sigma^2 = \langle x^2 \rangle - \langle x \rangle^2$$

$$= 0 \cdot (1-p) + 1 \cdot p - p^2 = p(1-p)$$

We conclude that for fixed p ,
 and $n \rightarrow \infty$, $b(k, n, p)$ tends to Gaussian distribution with mean np and standard deviation $\sqrt{p(1-p)n}$.

(c) we can do the same trick for the Poisson distribution. Fix λ finite and let x_i be independent with each a Poisson with parameter λ . Then from additivity $y = x_1 + \dots + x_n$ has Poisson distribution with parameter $n\lambda = \lambda$. Meanwhile from central limit theorem, y is Gaussian with mean np and standard deviation $\sqrt{n}\sigma$. Here again μ and σ are mean and s.d. of individual distribution.

$$\begin{aligned} \mu &= \sum k \exp(-\lambda) \frac{\lambda^k}{k!} \\ &= \lambda e^{-\lambda} \frac{d}{d\lambda} \left(\sum \frac{\lambda^k}{k!} \right) = \lambda e^{-\lambda} \frac{d}{d\lambda} (e^\lambda) \\ &= \underline{\underline{\lambda}}. \end{aligned}$$

$$\begin{aligned}
 \sigma^2 &= \langle x^2 \rangle - \lambda^2 \\
 &= e^{-\lambda} (\lambda \frac{\partial}{\partial \lambda})^2 (e^\lambda) - \lambda^2 \\
 &= e^{-\lambda} \left(\lambda \frac{\partial}{\partial \lambda} \right) (\lambda e^\lambda) - \lambda^2 \\
 &= \lambda + \lambda^2 - \lambda^2 = \lambda = \mu
 \end{aligned}$$

So, for large λ , any Poissonly distributed random variable y has mean λ and standard deviation $\sqrt{\lambda}$.

We will soon show that ^{counts in} any counting experiment (e.g. number of raindrops in a bowl or particles in a high energy physics experiment) is described by a Poisson random variable with (large) parameter $\lambda = \text{number of expected counts}$. This above result then translates into well known result that counting experiment with N events, has error \sqrt{N} .

D3/2 Birth and Death Processes

(a) We will briefly sketch some techniques which at their most elementary, justify Poisson for a counting experiment. The more advanced applications are important in many practical problems. (Telephones (random calls) and biology (random birth/deaths)) involving happenings distributed more or less randomly in time.

(b) Telephone calls

We start with a discrete example:

Suppose your favorite secretary receives N telephone calls in an hour. we wish to find probability that she receives R ^{calls} in a given time t .

Consider every second: } probability
 $p = N/3600$ of receiving one call : }
 essentially no probability of receiving
 > 1 call. Thus each second can be
 described by variable y_i ($i=0, 1, 2, \dots$)
 one call which has Poisson distribution
 with parameter $\lambda = N/3600$. The
 number of calls in time t seconds
 is $\sum_{i=0}^t y_i$ which from additivity is
 also Poisson with parameter $\sum_{i=0}^t \lambda_i$
 $= pt$. So

$$P_k(t) = e^{-pt} \frac{(pt)^k}{k!} \text{ is probability}$$

of receiving k calls in t seconds.

Note mean number received is pt
 or for $t = 3600$, just N as input
 data claimed.

One could also derive this

by taking each second as a Bernoulli trial with probability p of success: then we get distribution $b(k, t, p)$ which has $t \rightarrow \infty, p \rightarrow 0$: pt finite \therefore we can employ usual Poisson limit.

(c) Basic Birth Process

it is natural to take limit of time interval (second is above example) to zero, and let t be a gross continuous and not a discrete valued variable. Suppose we are studying events (e.g. telephone calls as above) that

- a). Occur independently of previous history of world.
- b). In time $[t, t+dt]$, the probability of an event happening is $\lambda dt + o(dt)$ and

the probability of >1 event happening is $\circ(dt)$.

Define $P_k(t)$ as the probability of k events happening in time $0 \leq t$. Then

$$P_k(t+dt) = P_k(t) [1 - \lambda dt]$$

↑
i.e. k events + $P_{k-1}(t) \lambda dt + \circ(dt)$
at $t+dt$

$$\text{or } \frac{dP_k(t)}{dt} = -\lambda P_k(t) + \lambda P_{k-1}(t) \quad : k \neq 0$$

For $k=0$, the above is modified as
 P_{0-1} term doesn't exist i.e.

$$\frac{dP_0(t)}{dt} = -\lambda P_0(t) \quad , \quad k=0$$

Solving the last gives

$$P_0(t) = e^{-\lambda t} \quad \text{and then induction shows } P_k(t) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}$$

and we have achieved Poisson again
 It is now obvious why any counting experiments for random events

Should have a Poisson Distribution.

(d) Birth and Death Process (Feller p.40)

Consider a system with a denumerable number of states E_n whose probability of occupation $p_{n(t)}$ varies with t . For Poisson E_n is "we have observed n events" (note only one E_n occupied at a time) and \exists chance λdt of a transition $E_n \rightarrow E_{n+1}$. This is a pure birth process.

We can generalize λ to $\lambda_n(t)$: a function of n and t . Also we can have death processes which also allow relapses $E_{n+1} \rightarrow E_n$ with chance $\mu_n(t)dt$. So:

chance $E_n \rightarrow E_{n+1}$	chance $\lambda_n dt + o(dt)$
$E_n \rightarrow E_{n-1}$	$" \quad \mu_n dt + o(dt)$
Any other change	$" \quad o(dt)$

We can at once, derive a similar

differential equation for P_n : probability of system being in state E_n .

$$\frac{dP_n(t)}{dt} = -(\lambda_n + \mu_n) P_n(t)$$

$$+ \lambda_{n-1} P_{n-1}(t) + \mu_{n+1} P_{n+1}(t)$$

with $\mu_0 = 0$ and $\lambda_{-1} = 0$ so that above is still correct with $n=0$.

(e) Example of Birth/Death Process

Consider a telephone exchange with an infinite number of lines and let E_n be state "n lines are busy". Then we apply above formalism with

$$\mu_n = n \mu \quad \mu \text{ fixed. } (\frac{1}{\mu} \text{ is average length of call})$$

$$\text{but } \lambda_n = \lambda \text{ independent of } n.$$

One can find the steady state solution - independent of both t and initial state:

$$P_n(t) = e^{-\lambda t / \mu} \left(\frac{\lambda}{\mu} \right)^n / n! \quad \text{as}$$

probability of n lines being busy.

It is also easy to do - the ~~same~~ seemingly harder problem - of a finite number of telephone lines in exchange. There is an example of this sort in problem set about parking lots.

Similar techniques can be used to decide on number of servicemen needed to repair a set of machines of prescribed reliability. Here

~~Death~~ = machine Starts working
~~Birth~~ = machine Breaks down.

and E_n = n machines broken. There is a beautiful discussion of this on pages 416-420 of Feller.

(5) Markov Chains

One can generalize either the discrete (case (b)) or continuous ((c) to (e)) cases by considering a system

with again states ϵ_1 and now not only transitions $\epsilon_{1 \rightarrow 2}$ but a general probability P_{1k} or P_{jk} of transition from ϵ_j to ϵ_k at each trial.

This theory includes general random walk problem and is discussed in Vol 1. Feller. The mathematics is not very deep - Solving recurrence relations - and the techniques are not useful for the basic problem in physics. (Statistics of data reduction). So we will not consider further. People with biological interest should however read further in Feller.

This Markov theory is relevant in discussions of the so called Metropolis method in Monte Carlo integration

These notes by refer to some additional references:

Cramer: Mathematical methods of statistics
published by Princeton.

plus

F.T. Solmitz Ann. Rev. Nucl. Sci. 14, 375 (1964)

and some old notes I have copies of

J. Olear, "Notes on Statistics for Physicists"
UCRL-8417 (1958)

S. Yellin, DESY 72/12 : Preprint (1972)

K. McDonald; CTSCL Internal Report #57.

D4 Estimation of Parameters

So far we have been concerned with rather straightforward problems on random variables with known probability distributions. Now we turn to a fundamental scientific problem: given some observations $x_1 \dots x_n$ (= values of a random variable) of some physical quantity, how do we extract good theoretical parameters describing system we have observed and how do we estimate error in our determination. This will require rather more careful thought as to relation of probability theory and real world. We will first describe the principle of maximum likelihood - everything else can be derived from this.

D4/1 Maximum Likelihood Principle

(i) Example:

(a) Rod Quivering in an Earthquake

x_1, x_2, \dots, x_n are n measurements of the length of a rod. Owing to a Gaussian earthquake (or a Gaussian f.e. normal) darkness) occurring during measurement, the x_i are scattered about the true value: you visit your favourite theorist (or do yourself after sobering up) who believes that the probability distribution of x_i is ^{Gaussian} named with mean $\alpha_0 =$ ~~true length~~ true length and standard deviation σ_0 .

We may write down the (relative) probability of each measurement x_i as:

$$p_i(x_i) = c \exp \left\{ -\frac{1}{2\sigma_0^2} (x_i - \alpha_0)^2 \right\}$$

The likelihood function for whole experiment of n events is defined as :

$$L(\underline{x_1 \dots x_n}) = \prod_{i=1}^n c \exp \left\{ -\frac{1}{2\sigma_0^2} (x_i - \alpha_0)^2 \right\}$$

$$= C \exp \left\{ -\frac{1}{2\sigma_0^2} \sum_i (x_i - \alpha_0)^2 \right\}$$

This is - upto a constant - the probability of the x_i given the true length α_0 and certain theoretical assumptions we represent as T .

$$L(x_1 \dots x_n) = \Pr \{ x_1 \dots x_n | \alpha_0, T \}$$

in conditional probability notation. T is for instance the notion that distribution Gaussian and the value of σ_0 taken (the latter being calculable from the strength of earthquake or amount of beer drunk).

(b) Enter Mr. Bayes

Now what we really want is the probability distribution of an estimate α of α_0 . We can formally get this by use of Bayes law: first extend $p(x_i) = \exp\left\{-\frac{1}{2\sigma_0^2}(x_i - \bar{x})^2\right\}$ to be defined for any value of α : not necessarily $\alpha = \text{true value } \alpha_0$. Introduce:

$P(\alpha | x_i, T)$ = Probability true length is α given measurements x_i and pre-judice T . Then by Bayes + or rather a slight extension of it.

$$P(\alpha | x_i, T) = \frac{P(x_i | \alpha, T) P(\alpha | T)}{P(x_i | T)} \quad (D4.1)$$

which only differs from ordinary Bayes by extra $|T$ in each probability.

This is the essence of the maximum likelihood principle. (D4.2) converts the likelihood $L(x_1 \dots x_n) = P(x_i | \alpha, T)$ which is a function of x_i for fixed probability } of α for fixed

α, T into $P(\alpha | z_i; T)$ which is the desired probability distribution of α for given measurements z_i and prejudget T .

(c) Interpretation of (b), (c)

(a) we can ignore $P(z_i | T)$ because it is independent of α and fixed for our given z_i and T . It has the practical point that likelihood is never normalized sensibly.

(b) $P(\alpha | T)$ philosophically represents a prior knowledge of α before experiment (measurement of z_i) ~~is performed~~. This is shaky - not because we don't really know how to define what T ("given T ") means - but, for instance, suppose we wish to parameterize "complete ignorance" of α before

experimental performed.

Natural is $P(\alpha|T) = \text{constant}$

But Suppose we put $\beta = \alpha^2$ and
then β has a priori probability $P(\beta|T)$

$$P(\beta|T) d\beta = P(\alpha|T) d\alpha$$

$$\text{i.e. } P(\beta|T) = \frac{1}{2\sqrt{\beta}} \neq \text{constant}$$

So the concept of "complete
ignorance" is not well defined. Jeffreys
(a big bad English-oxford-book cited in
mathews and
waller) Spends 200's of pages trying
to avoid this difficulty. His efforts are
so much garbage.

Fortunately there is no trouble for
large n because of both an obvious
numerical reason (see next section (i)(a))
and a not so obvious mathematical
reason (See section (ii)).

(a) The Likelihood Maximizes

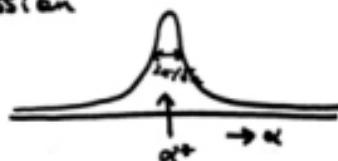
Consider, again, the explicit form of the likelihood of a red quivering in the earthquake.

$$\mathcal{L}(x_1 \dots x_n) = c \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \alpha)^2 \right\}$$

According to (34.1), we now regard \mathcal{L} as a function of α and so rewrite it:

$$\mathcal{L}(\alpha) = c \exp \left\{ -\frac{n}{2\sigma^2} \left(\alpha - \frac{1}{n} \sum_{i=1}^n x_i \right)^2 - \underbrace{\frac{1}{2\sigma^2} \left[\sum_{i=1}^n x_i^2 - \frac{1}{n} (\sum_{i=1}^n x_i)^2 \right]}_{\text{constant}} \right\}$$

The last term is a constant independent of α , and so as a function of α , \mathcal{L} is just a Gaussian which maximizes at $\alpha = \frac{1}{n} \sum_{i=1}^n x_i$ and has width σ/\sqrt{n} .



So as $\ell = \Pr(\alpha | x_i, \tau)$ we say
that maximum likelihood estimate of
 α is $\hat{\alpha} = \frac{1}{n} \sum_{i=1}^n x_i$ and its standard
deviation is σ/\sqrt{n} .

Note that whatever $P(\alpha)$ you
put in the shape of $\ell \cdot P(\alpha | \tau)$ is
essentially independent of P as $n \rightarrow \infty$
and we get a Gaussian peaked at
 $\hat{\alpha}$ riding on top of constant value $P(\alpha | \tau)$.

(ii) General maximum likelihood method

(a) The above is glib in many ways -
especially in its treatment of errors. we will
repair this later with real life theorems.

First we state the obvious generalization.

The general maximum likelihood
method is to put $P(\alpha | \tau) = 1$ in (34.2) and
form $\ell(\alpha) = \prod_{i=1}^n p(x_i; \alpha)$ for any

measurement x_i whose probability $p(x_i, \alpha)$ depends on value x_i and theoretical parameter(s) α .

The maximum likelihood estimate of α is $\hat{\alpha} = \hat{\alpha}^*$ where $L(\hat{\alpha}^*)$ is a maximum. The error is the standard deviation of the distribution of $\hat{\alpha}$ as a function of α i.e.

$$L(\alpha) = \exp \left\{ \text{const} - \frac{1}{2} \frac{1}{\sigma^2} (\alpha - \hat{\alpha}^*)^2 \right\}$$

(b) One must issue some obvious warnings. For instance L could have several maxima



For instance, this would happen if in earthquake problem you were still drunk in morning and decided to

find $\hat{\alpha}$. Then $\hat{\alpha}$ would peak at $\pm \hat{\alpha}^*$.
 Less obviously one should remember that theorems (to come) only say that maximum likelihood method works for large n . ~~Then it is a~~
~~universal~~ ~~method~~ ~~universally~~ good method. However for small n it is much more likely to give nonsense than many other methods. Thus...

(c) Warning Example

Consider probability distribution appropriate for golanzahlen except:

$$p(\varphi, \alpha) = (1 + \alpha \cos \varphi) / 2\pi : -1 \leq \alpha \leq 1.$$

If we had but one event

$$\mathcal{L}(\alpha) = (1 + \alpha \cos \varphi_1) / 2\pi$$

if $\cos \varphi_1 > 0$: $\alpha = 1$ (i.e. end of range)
 maximizes \mathcal{L}

or $\cos \varphi_1 < 0$: $\alpha = -1$ (other end of range)
 maximizes \mathcal{L} .

If you have 2 events:

$$L(\alpha) = (1 + \alpha(\cos\varphi_1 + \cos\varphi_2) + \alpha^2 \cos\varphi_1 \cos\varphi_2)/k\pi^2$$

This has a maximum at end of range
 again,
 unless $\cos\varphi_1, \cos\varphi_2 < 0$.

(d) So usually people warn you by saying: "If distribution of α not Gaussian, then do not simply quote maximum and standard deviation, rather plot the likelihood". Unfortunately this is dubious, because for small n , nobody has proved any theorems and $P(\alpha|T)$ philosophy ~~important~~
 while for large n , theorems are proven on the relevance of the likelihood method but they also prove L is Gaussian!

So correct warning is: "If

distribution not Gaussian - everything is OK if either:

(a) just many sharp humps as in the example, data can be fitted with > 1 value of parameter.

(b) If several parameters, can be Gaussian in some and flat in others if depends weakly (or not at all!) on them.

(c) Gaussians maybe distorted by effects due to end of range.

If (a), (b) or (c) don't apply and χ^2 isn't Gaussian, then one can only be unhappy!

(iii) Lifetime Example

All books give the example of determining the lifetime τ_0 of a particle from the observation of decay times $t_1 \dots t_n$ of n decays.

The probability distribution is Poisson:

$$p(t) dt = \frac{1}{\tau} e^{-t/\tau} dt$$

Now dt is meaningless -

it is independent of τ and only affects normalization of likelihood. The unknown $P(x_i|\tau)$, $P(x_i|\tau)$ terms is (obvi) made this arbitrary anyway!

$$\text{likelihood } L(\tau) \propto \frac{1}{\tau} e^{-t_1/\tau} \dots \frac{1}{\tau} e^{-t_n/\tau}$$

$$\text{or } L(\tau) = \exp \left[-\frac{1}{\tau} \sum_{i=1}^n t_i - n \ln \tau \right]$$

where we lazily convert \propto into an = sign.

maximum is $\frac{dL}{dz} = 0$ or

$$\frac{1}{T^2} \sum_{i=1}^n t_i - \frac{n}{T} = 0$$

$$\text{or } \bar{t}^* = \frac{1}{n} \sum_{i=1}^n t_i \text{ or the maximum}$$

likelihood estimate is just the mean:
such a simple result is only true because
of particular Poisson example:

We can use this example, to
illustrate an interesting theoretical
point about the error.

(a) According to the central limit
theorem, the variable $y = \frac{1}{n} \sum_{i=1}^n t_i$ is a
sum over independent random variables
and has mean $\langle y \rangle = \tau_0$ and
standard deviation σ

$$\sigma^2 = \frac{1}{n} \left\{ \langle t_i^2 \rangle - \langle t_i \rangle^2 \right\}$$

$$= \frac{1}{n} \left\{ \frac{1}{\tau_0} \int t^2 e^{-t/\tau_0} dt - \tau_0^2 \right\}$$

$$\text{or } \sigma = \tau_0 / \sqrt{n}$$

(D4.2)

(v) Multiplication of Experiments

If we have two experiments that are independent but depend on same theoretical parameter α .

Then $L_1 = \prod_{i=1}^{n_1} p(x_i, \alpha)$ is likelihood of expt #1

$L_2 = \prod_{i=1}^{n_2} q(y_i, \alpha)$ is likelihood of expt #2

So and clearly the combined likelihood is

$$L = L_1 L_2$$

i.e. one should multiply likelihoods when combining experiments. We will now apply this to case when each L_i is Gaussian in α

24/2 χ^2 method

This is standard and fully correct (best) method for combining results from several experiments. It is in fact used (because it's easy, conceptually and for computer) even when it's inapplicable e.g. if there are correlations between "supposed "independent" random variables or if ~~other~~ basic experiments have too few events to be Gaussian.

The problem is : given N experimental measurements x_i and errors σ_i plus estimates \bar{s}_i of the x_i which are functions of n theoretical parameters. Then the probability of each observation is

$$p_i = \frac{1}{\sqrt{2\pi} \sigma_i} \exp \left\{ -\frac{(x_i - \bar{s}_i)^2}{2\sigma_i^2} \right\}$$

$$\text{The likelihood } L = \prod_{i=1}^N p_i$$

$$\propto \exp \left\{ - \sum_{i=1}^N (x_i - \xi_i)^2 / 2\sigma_i^2 \right\}$$

So, the principle of maximizing the likelihood is equivalent to minimizing χ^2 where

$$\chi^2 = \sum_{i=1}^N (x_i - \xi_i)^2 / \sigma_i^2 \quad (44/2.2)$$

wt n theoretical parameters in ξ_i .

- (i) Notice that χ^2 is the sum $\sum_{i=1}^N v_i^2$ where v_i are independent Gaussian random variables of unit standard deviation and zero mean.
- we can calculate χ^2 distribution exactly (M. and W. pages 393-395) but in practice this is rarely necessary.
- Thus we can easily show that - for fixed ξ_i - that

$$\begin{aligned} \langle \chi^2 \rangle &= N \\ \sigma_{\chi^2}^2 = \cancel{\dots} &= 2N \end{aligned} \quad (D4/2.2)$$

and this is sufficient for large N as central limit theorem says χ^2 will be Gaussian with above mean and standard deviation.

(iii) Now (D4/2.2) is not quite correct because if we determine the theoretical parameters to minimize χ^2 then implicitly ξ_i are also random variables that ~~are~~^{are} functions of x_j . It is not too hard to show that in this case one should just replace N by $N-n$ (number of degrees of freedom) in (D4/2.2). In any case, (D4/2.2) is very important because it allows an easy criterion for the goodness of fit. Thus if value

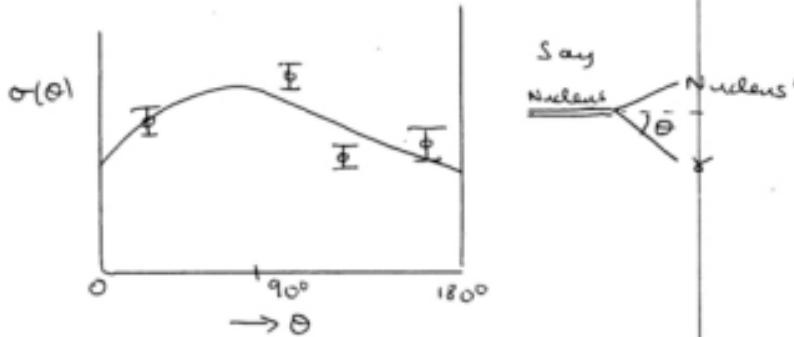
of χ^2 is many (says) standard deviations away from its mean, then the theoretical form is in doubt. The usual maximum likelihood method does not have this advantage : remember L was unnormalized and so its value had no significance for the goodness of fit. This is quite a difficult problem and an elementary discussion is given in Macdonald notes page 14 onwards.

(iii) The only way to get a feeling for χ^2 is to use it : So I want here to ~~just~~^{briefly} dwell on the important techniques used to minimize χ^2 to get good theoretical parameters and estimate errors in the same. m. and w. Solve

(v) Example of use of χ^2

Consider the example in Mathews and Walker page 340

$$\sigma(\theta) = a_0 + a_1 \cos\theta + a_2 \cos^2\theta + a_3 \cos^3\theta$$



The measurements x_i 's are cross sections and the theoretical predictions f_i are linear in unknown parameters a_j .

$$f_i = \sum_{m=1}^n c_{im} a_m$$

$$\text{minimize : } \sum_{L=1}^N \frac{(x_L - f_L)^2}{\sigma_L^2}$$

$$x_i = \text{const. } N_i$$

$$\sigma_i = \text{const. } \sqrt{N_i} \quad \sigma_i = x_i / \sqrt{N_i}$$

where i th bin has N_i events.

$$\sum_{i=1}^N \frac{(x_i - \hat{\xi}_i)}{\sigma_i^2} \frac{\partial \hat{\xi}_i}{\partial a_m} = 0$$

or $\sum_{\ell} m_{m\ell} \hat{a}_{\ell}^* = \theta_m$

* denotes
 χ^2 or maximum
 likelihood estimate

$$\begin{aligned} {}^*\chi^2 \text{ matrix} &= \sum_i C_{im} ({}_i \ell) && \text{Positive Symmetric} \\ &= \sum_i \frac{(\partial \hat{\xi}_i / \partial a_m)(\partial \hat{\xi}_i / \partial a_{\ell})}{\sigma_i^2} \end{aligned}$$

$$\ell_m = \sum_i \frac{x_i C_{im}}{\sigma_i^2} = \sum_i \frac{x_i}{\sigma_i^2} \frac{\partial \hat{\xi}_i}{\partial a_m}$$

$$\underline{\hat{a}^*} = m^{-1} \underline{\ell_m}$$

(vi) we can find errors in either of the two ways

- regard a^* as a function of x_i and calculate standard deviations from central limit theorem
- Shape of likelihood.

$$\mathcal{L} = \exp(-\frac{1}{2} \chi^2).$$

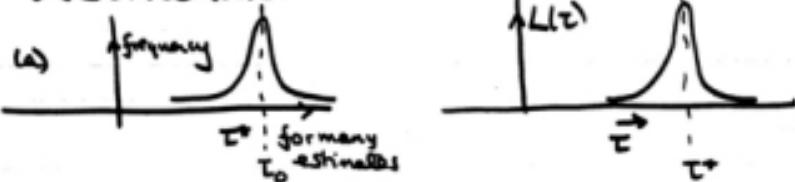
we have - as far as shape of \mathcal{L} concerned

$$\begin{aligned} \mathcal{L} &= \exp \left(\text{Independent of } a \right. \\ &\quad \left. + \text{no terms linear in } a-a^* \right. \\ &\quad \left. \text{due to "maximum" condition} \right. \\ &\quad - \frac{1}{2} (a-\underline{\alpha})^T M (a-\underline{\alpha}) \end{aligned}$$

i.e. M is error matrix for a .

The mean and standard deviation of \bar{y} refer to distribution of values of \bar{y} in many experiments - each of n observations.

(ii) However maximum likelihood method of finding error is, off hand, quite different. Rather than looking at distribution of \bar{y} over infinitely many experiments each of n observations, it considers τ dependence of L for one given experiment of n observations.



Then we found it became Gaussian and calling maximum likelihood standard deviation $\hat{\sigma}$

$$L(\bar{\tau}) \propto \exp \left\{ -\frac{1}{2\hat{\sigma}^2} (\bar{\tau} - \bar{\tau}^*)^2 \right\}$$

$$\text{or } \frac{1}{n} \bar{\sigma}^2 = - \left. \frac{d^2}{d\tau^2} \ln L(\tau) \right|_{\tau=\tau^*}$$

$$= \frac{2}{\tau^3} \sum_{i=1}^n t_i - \frac{n}{\tau^2} \Big|_{\tau=\tau^*}$$

$$\text{or } \bar{\sigma} = \tau^* / \sqrt{n} \quad (D4.3)$$

as $\tau^* = \tau_0 + O(\frac{1}{n})$ we see that (D4.2) and (D4.3) are in fact the same asymptotically (large n , that is).

Note that physics looks (Solvay is an exception) at motion method (a) of finding error in maximum likelihood method. On the other hand, mathematics looks never mention (b) as a serious method. They do however prove (a) and (b) are the same for large n (as we indicated in example above) for general distribution.

As far as I can see, method (a) is only correct definition of error. It follows that method (b) is wrong when it differs from method (a). (b) does differ from (a) for small n when, in particular, likelihood isn't Gaussian. So warnings for small n ("plot likelihood") are very dubious: the correct statement is that your results are simply unreliable and error is not calculable from likelihood function alone. The error is defined but must be calculated as in (a).

Note that when math books discuss maximum likelihood estimate for small n - it is only value w^t they discuss not

Shape of L and method (a) of finding emr.

(i) Asymptotic Validity of Maximum Likelihood method

Cramer chapters 22 and 23.

Theorem:

Let $P(x_i, \alpha)$ be a normalized probability for a result x_i given theoretical parameter α . we have n such observations x_i :

$L = \prod_{i=1}^n P(x_i, \alpha)$ is usual likelihood function:

Let α^* be maximum of L wrt α .

(a) If we fix n and conduct many such series of n observations

$$\langle \alpha^* \rangle = \alpha_0 + O(1/n)$$

Note it is not $= \alpha_0$ exactly. (α_0 is true theoretical parameter). If for any estimate $\langle \alpha_{\text{estimate}} \rangle = \alpha_0$, then we say

estimate is unbiased : if $\langle \text{estimate} \rangle = \alpha_0 + b(\alpha_0)$ then we have a biased estimate with bias $b(\alpha_0)$. So maximum likelihood estimate is for finite n biased and is only asymptotically unbiased. However bias is smaller than standard deviation and so it is of no real importance. Bias is discussed in detail by Cramer.

(b) We can regard $\hat{\alpha}$ as a random variable : then its standard deviation is $O(\sqrt{n})$ and to this order we can correctly calculate standard deviation in the "physics way" already discussed. (i.e. from a dependence of $L(\alpha)$).

(c) There is no other estimate of α_0 which has a smaller standard deviation asymptoticallyⁱⁿ, i.e. maximum

likelihood method is best ~~as~~ as long as n large

Proof: (c) can be found in Cramer p.479 or Hoel p.379. For (a) and (b), write

$$\ln L(\alpha) = \sum_{i=1}^n \ln P(x_i, \alpha)$$

Expand L about true theoretical values,

$$\ln L(\alpha) = L_0 + (\alpha - \alpha_0) L_1 + \frac{(\alpha - \alpha_0)^2}{2} L_2$$

$$L_0 = \sum_{i=1}^n \ln P(x_i, \alpha_0)$$

$$L_1 = \sum_{i=1}^n \frac{\partial}{\partial \alpha} \ln P(x_i, \alpha_0)$$

$$L_2 = \sum_{i=1}^n \frac{\partial^2}{\partial \alpha^2} \ln P(x_i, \alpha_0)$$

We get maximum likelihood estimate, by calculating

$$L_1 + (\alpha^* - \alpha_0) L_2 = 0$$

$$\text{i.e. } \alpha^* = \alpha_0 - \frac{L_1}{L_2} \quad (D4.4)$$

The error in this is $O\left\{(\alpha^* - \alpha_0)^2 \frac{L_3}{L_2}\right\}$

which is $O\left(\frac{1}{n}\right)$ as we will soon show that $(\alpha^* - \alpha_0)^2$ is $O\left(\frac{1}{n}\right)$.

Now L_1 and L_2 are random variables and (D4.4) defines α^+ as a random variable. We can calculate:

$$\begin{aligned}\langle \alpha^+ \rangle &= \alpha_0 - \langle L_1 / L_2 \rangle \\ &= \alpha_0 - \cancel{\langle L_1 \rangle / \cancel{\langle L_2 \rangle}} \\ &= \alpha_0 - \cancel{\langle \frac{\partial}{\partial \alpha} \ln P(x, \alpha_0) \rangle} / \cancel{\langle \frac{\partial}{\partial \alpha^2} \ln P(x, \alpha_0) \rangle}\end{aligned}$$

Use theorem on F/G, proved at end of D2: note this is only valid to $O(\gamma_n)$. Now the numerator is

$$\begin{aligned}\langle \frac{\partial}{\partial \alpha} \ln P(x, \alpha_0) \rangle &= \text{by definition} \\ &\int dx P(x, \alpha_0) \frac{\partial}{\partial \alpha} \ln P(x, \alpha_0) \\ &= \int dx \frac{\partial}{\partial \alpha} P(x, \alpha_0) = \frac{\partial}{\partial \alpha} \left[\int dx P(x, \alpha_0) \right] \\ &= \frac{\partial}{\partial \alpha} \cdot 1 \quad \text{as } P \text{ is truly normalized.} \\ &= 0. \quad \text{This proves part (a) that} \\ \langle \alpha^+ \rangle &= \alpha_0 \text{ with corrections from both} \\ L_1 / L_2 \text{ and } L_1 / L_2 &\neq \text{quite } 0 \text{ terms which are both } \gamma_n. \quad \text{These corrections}\end{aligned}$$

have been calculated by Yellin.

Now to prove part (b): first consider "physics" method. This takes Standard deviation of α^* as $\sqrt{-1/L_2}$ or rather this evaluated at $\alpha = \alpha^*$ not $\alpha = \alpha_0$ but these are same asymptotically.

$$\sigma_{\text{physics}}^2 = -1 / \sum_{i=1}^n \left[\frac{\partial^2}{\partial \alpha^2} \ln P(x_i, \alpha_0) \right]$$

$$\langle \sigma_{\text{physics}}^2 \rangle = -1/n \langle \frac{\partial^2}{\partial \alpha^2} \ln P(x_0, \alpha_0) \rangle$$

$$\text{where } \langle \frac{\partial^2}{\partial \alpha^2} \ln P(x, \alpha) \rangle = \int dx P \left\{ -\frac{(\partial P / \partial \alpha)^2}{P} + \frac{\partial^2 P / \partial \alpha^2}{P} \right\}$$

Now second term is just $\frac{\partial^2}{\partial \alpha^2} \int P$ and is zero again because of normalization condition. So

$$\langle \sigma_{\text{physics}}^2 \rangle = 1 / \left\{ n \int dx (\partial P / \partial \alpha)^2 / P \right\} \quad (D4.5)$$

we can calculate the correct stochastic standard deviation in α^* by taking

(D4.4) and using our F/g theorem
 with $f_i = \frac{\partial}{\partial \alpha} \ln P(x_i, \alpha)$, $g_i = \frac{\partial^2}{\partial \alpha^2} \ln P(x_i, \alpha)$
 generally

$$\sigma^2 = \frac{1}{n} \left\langle \left(f_i \right)_g - g \left\langle f_i \right\rangle_g \right\rangle^2$$

but $\left\langle f_i \right\rangle_g = 0$ - proved above, so:

$$\sigma^2 = \frac{1}{n} \left\langle f_i^2 \right\rangle_g - \left\langle f_i \right\rangle_g^2$$

It is easy to see that $\left\langle f_i^2 \right\rangle_g = \int \left(\frac{\partial P}{\partial \alpha} \right)^2 / p_2 \cdot P d\alpha$
 and this $= \left\langle g_i \right\rangle_g$. So $\sigma^2_{\text{statistical}}$ is just $1/n \left\langle g_i \right\rangle_g$ which is exactly equal to $\overrightarrow{\sigma^2_{\text{physics}}}$ from (D4.5). The expectation value is no sweat: from central limit theorem $\left\langle \sigma^2_{\text{physics}} \right\rangle$ and $\sigma^2_{\text{physics}}$ differ only by relative error $O(\sqrt{n}) \rightarrow 0$ as $n \rightarrow \infty$.

the simplest problem i.e. \mathbf{f}_i are linear in theoretical parameters.

Berntsen should be read before tackling a real problem.

(ii) Note that if x_i are calculated from the same population e.g.

$$x_1 = \frac{1}{n} \sum_{i=1}^n y_i$$

$$x_2 = \frac{1}{n} \sum_{i=1}^n y_i^2 \quad \text{etc.}$$

Then they are not independent but we can use the theorem proved at the end of Section 32 to show that x_i will (for large n) be Gaussian but there is a correlation matrix

$$\text{i.e. } L \propto \exp \left\{ -\frac{1}{2} (\mathbf{z} - \mathbf{f})^\top \mathbf{M} (\mathbf{z} - \mathbf{f}) \right\}$$

as discussed in Solndy.

By diagonalizing \mathbf{M} we can (in principle) reduce this to

the usual χ^2 method with x_i replaced by eigenvectors of $\underline{\underline{M}}$.

This idea is embodied in the method of moments when one calculates several moments from the same distribution of events. In fact one usually drops (forgets) off diagonal terms in $\underline{\underline{M}}$ and forms χ^2 from non-independent random variables. This is done in high energy physics when fitting the moments g_{00}, g_{11}, g_{20} of a 1⁻ particle decay. Of course, in this case, the off diagonal elements are never given by the experimentalist recording his data: so one is forced to set them zero.

vii) One should also mention that there are two versions of χ^2 depending on whether σ_i taken from "experiment" or estimated theoretical. In the latter case, it will be a function of the parameters to be varied.

For example, Suppose x_1 is number of events in some counting experiment i.e. σ_i (experimental) = $\sqrt{x_1}$ but σ_i (theory) = $\sqrt{S_1}$. Now σ_1 is only Gaussian for large x_1 but then $\sqrt{x_1}$ and $\sqrt{S_1}$ are essentially same. ~~use standard theorem~~ If $x_1 = 0$, then theoretical σ is ~~to be~~ preferred! - although χ^2 method itself is in difficulties as $\pi(x_1)$ not Gaussian. Best to go back to maximum likelihood and put in true Poisson distribution for x_1 .

D4/3 Counter Physics

For obvious reasons, I am forced to take my examples from high energy physics.



Take m happy smiling counters - all arranged in a row : Suppose they count random events that fall in their province. Then each of them has a parameter λ_i associated with them such that distribution of number n_i of events in counter i is Poisson with parameter $\lambda_i t$ (t is time). This follows from our theory of Birth processes. In HEP, $\lambda_i \propto \int_{BIN} dt \frac{d\sigma}{dt}$ where BIN is t -slice covered by counter.

This is a straightforward application of likelihood:

$$p_i = e^{-\lambda_i t} (\lambda_i t)^{n_i} / n_i! \quad (D4/3.1)$$

(n_i : observed, λ_i : theory : t maybe known or not).

$$\mathcal{L} = \prod_{i=1}^m p_i \quad (D4/3.2)$$

$$\text{or } \ln \mathcal{L} = - \sum_{i=1}^m \ln n_i! - t \sum_{i=1}^m \lambda_i + \sum_{i=1}^m n_i \ln(\lambda_i t)$$

(iii) Unnormalized Theory Experiment

Suppose t - which includes beam intensity etc - is unknown. Then we use (D4/3.2) to determine θ . (Regard t as a theoretical parameter)

$$\frac{\partial (\ln \mathcal{L})}{\partial t} = 0.$$

$$\text{or } t = \sum_{i=1}^m n_i / \sum_{i=1}^m \lambda_i = N / \Lambda$$

where N = total number of events
 Λ = Total Poisson parameter
 Summed over all counters.

Now if we ~~forget~~^{eliminate} t, we get

$$\ln L = - \sum_{i=1}^m n_i \ln n_i! - N + N \ln N$$

$$+ \sum_{i=1}^m n_i \ln (\lambda_i / \bar{\lambda})$$

or dropping terms independent of λ_i :

$$\ln L = \sum_{i=1}^m n_i \ln (\lambda_i / \bar{\lambda}) \quad (14/3.3)$$

$$= \sum_{j=1}^N \ln \hat{p}(g_j)$$

where g_j is counter where event added and $\hat{p}(k) = \lambda_k / \bar{\lambda}$ is probability of landing in k'th counter.

Note \hat{p} is normalized $\sum_{k=1}^m \hat{p}(k) = 1$.

This is discrete maximum likelihood method: we have N observations:

each has only m possible outcomes and the $L = \prod_{j=1}^N \hat{p}(g_j)$ is product of probabilities - normalized properly.

If λ_i are functions of some theoretical parameters, we find them by maximizing (4/3.3). In particular if all values of $\mu_i = \lambda_i/N$ are free parameters, we maximize (4/3.3) subject to constraint $\sum_{i=1}^n \mu_i = 1$: this gives us

$$\mu_i = n_i/N$$

(ii) Normalized Experiments

Now we suppose it is no longer an unknown theoretical parameter, but a known number.

$$\ln L = \text{const}$$

$$+ \underbrace{\{ N \ln(t\Lambda) - t\Lambda \}}_0 + \underbrace{\sum_{i=1}^n n_i \ln(\lambda_i/\Lambda)}_{(b)}$$

(b) is as before. (a) can be regarded as $\ln p_N$ where p_N is probability of getting N (which should now be

regarded as an experimental observable total counts. ($p_N = e^{-t\lambda} (t\lambda)^N$ which is up to constant, usual Poisson).

This requires some care to apply previous theorems as probabilities ~~must be~~
~~not~~ normalized in proofs in 34/2.
However theorems do hold.

Note one can call this extended maximum likelihood method.

(iv) Gaussian limit

of course if ~~all~~ each n_i is large

$$p_i \propto \exp \{ n_i - \lambda_i t \} / \lambda_i t$$

we get back χ^2 method. - note σ_i is theoretical value! (replace $\lambda_i t$ in denominator by n_i to get experimental value). So we only get something different from maximum likelihood method, if some n_i small and Poisson \neq Gaussian.

Day 4 Bubble chamber physics

(i) This is like D4/3 except that every event (not every counter) has its own probability p_i . (given events labelled by continuous parameter(s), p_i is a continuous probability density). So far, we haven't discussed how to apply maximum likelihood method to counting problem with continuous probability distribution. (Lifetime example was continuous but not a counting experiment). What we do is divide range of continuous parameters $\{0 \leq \theta \leq 2\pi\}$ in $(1 + p \cos \theta)/2\pi$ polarization distribution into m small bins. m should be so large that (a) probability of 2 events in a given experiment falling in same

bin is geo. (b) theoretical p_i constant over bin.

Then we have m observations - each with one of two results.

Success - 1 event in bin
Failure - 0 events in bin.

Let N be total number of events.
then likelihood is

$$\mathcal{L} = \prod_{\text{successes}} p(\text{success}) \prod_{\text{failures}} p(\text{failure})$$

$$= \prod_{i=1}^N \exp(-p_i t dx_i) (p_i dx_i t)$$

$$\textcircled{3} \prod_{j=1}^{m-N} \exp(-p_j dx_j t)$$

writing out \mathcal{L} in terms of Poisson probabilities for 0 or 1 events.

Now $m \gg N$ and m bins dx_i cover space. Thus

$$\mathcal{L} = \exp\left[-\int p(x) dx t\right] \prod_{i=1}^N p_i dx_i t ..$$

Now - fortunately - arbitrary dx_i ,
 $i=1..N$ are independent of theoretical
parameters and so can be dropped.

Now (at further cost to the significance
of normalization of L).

So our fundamental formula for
 L in a continuous distribution is

$$L = \exp \left[- \int p(x) dx t \right] \prod_{i=1}^N p(x_i, \alpha) t$$

(ii) Note - that - just as in D4/3 -
if t is unknown we can find it from:

$$\frac{\partial L}{\partial t} = 0 \text{ which gives}$$

$$t = N / \int p(x) dx$$

when , w.r.t remaining parameters,

$$L(\alpha) \propto \prod_{i=1}^N \left\{ p(x_i, \alpha) / \int p(x, \alpha) dx \right\}$$

- the product of normalized probability
distributions.

If t is known, the discussion is again similar to this case in D4/3.

(iii) Binned Data

If we have so many events that I can divide them into K groups of n_k events where p constant has same value for each member of group, then:

$$\lambda = \exp \left\{ - \int p dx t \right\} \prod_{k=1}^K (pt)^{n_k}$$

$$\propto \exp \left\{ - \sum_{k=1}^K \lambda_k t \right\} \prod_{k=1}^K (\lambda_k t)^{n_k} (n_k!)$$

$$\lambda_k = \int_{\text{bin for } k^{\text{'}} \text{th group}} p dx$$

and we return exactly to formulae in D4/3. Further if n_k large for each k , our Poisson distributions become Gaussian and we can use χ^2 method to analyse data.

This is the only case where χ^2 correct in sense of giving best answer. However even if N not so large that p constant in each group, one can still bin data and throw away exact values of x_i for each event - just record what bin it's in. Obviously if p not constant over each bin, this throws away information. However (D&H 4.1) is now exact for L and one chooses bin so that n_B large and χ^2 applicable. This is used in practice because χ^2 is so much easier and cheaper to use. often one can bin data and minimize χ^2 : this gives values for theoretical

parameters which can be used as initial values for expensive full likelihood fit which will then need only a few iterations. One can carry it a bit further - and use a super low grade method (e.g. method of moments - see 24/5) to initialize parameters before χ^2 fit

This is discussed by Solmitz
("histogram method") and Fadie et al.

D4/5 Method of moments

Q This is a simple and, in general, non-optimal method. even here -for >1 parameter - there are two versions (a) ignore correlations (used in practice but wrong) and (b) include correlations (best of a bad lot). of course simplicity can lead to huge savings in computer time and so it should not be sneered at.

The basic idea of method of moments is simple: consider a set of observations x_i , then central limit theorem implies

$$\frac{1}{N} \sum_{i=1}^N f(x_i) \rightarrow \mu$$

$$\text{where } \mu = \int p(x, u) f(x) dx.$$

Here μ is $\mu(\alpha)$ where α is some unknown parameter. So putting

$$\mu(\alpha) = \frac{1}{N} \sum_{i=1}^N f(x_i)$$
 gives one equation for α with solution $\alpha = \alpha^*$: α^* is moment estimate of α . Note:

- (a) One gets a different moment method for each choice of f : from m.l. theorem, all of these are asymptotically worse or perhaps just not better than maximum likelihood method. One must choose moment which is simple and has small expected error.
- (b) If we have to estimate n theoretical parameters α . Then we just get n equations by taking n moments (i.e. n different values of f).

(c) Error in μ is

theoretically $\frac{1}{\sqrt{N}} \sqrt{\int f^2 p(x, \omega) dx - \mu_0^2}$

experimentally $\frac{1}{\sqrt{N}} \sqrt{\frac{1}{N} \sum_{i=1}^N f^2(x_i) - (\frac{1}{N} \sum_{i=1}^N f(x_i))^2}$

using μ_0 to be true value of mean
and $\mu = \frac{1}{N} \sum_{i=1}^N f(x_i)$ to be estimate.

By central limit theorem, for large N ,
experimental estimate will ~~be~~^{have} correct
error. However - if all x_i same
(always true if $N=1$) - then
experimental error = 0, which is
simply wrong. The moral is: use
theoretical estimate if you only
have a few events.

(iii) Example:

$$p(x) = \frac{1}{2}(1+ax) \quad -1 \leq x \leq 1$$

(a) ~~Method~~ Moment method

$$\begin{aligned} \langle x^m \rangle &= \frac{a}{2} \int_{-1}^{+1} x^{m+1} dx \\ &= \frac{a}{m+2} \quad : m \text{ odd.} \end{aligned}$$

So if m_m = moment method using $\langle x^m \rangle$

$$a_m^+ = (m+2) \sum_{i=1}^N x_i^m / N$$

The error² is

$$\begin{aligned} \sigma^2|_m &= \frac{(m+2)^2}{N} \left\{ \int_{-1}^{+1} x^{2m} (1+ax) dx - a^2 / (m+2)^2 \right\} \\ &= \frac{1}{N} \left\{ \frac{(m+2)^2}{2m+1} - a^2 \right\} \end{aligned}$$

which is a minimum for $m=1$

$$\text{when } \sigma^2|_1 = \frac{1}{N} (3-a^2)$$

(b) Maximum Likelihood method

$$\begin{aligned}
 \text{Now } \frac{1}{[N\sigma^2]} \Big|_{m.e.m.} &= \int_{-1}^{+1} dx / \rho \left(\frac{\partial \rho}{\partial a} \right)^2 \\
 &\quad - \text{from our general theory} \\
 &= \frac{1}{2} \int_{-1}^{+1} dx \frac{x^2}{(1+ax)} \\
 &= \frac{1}{2a^2} \left\{ \ln \left[\frac{1+a}{1-a} \right] - 2a \right\} \\
 &= \frac{1}{3} + \frac{a^2}{5} + \frac{a^3}{7} + \dots \\
 &= \frac{1}{3} \quad a=0 \quad \text{which } = \sigma^2 \Big|_1 \text{ value} \\
 &= \infty \quad a=\pm 1. \\
 \text{So } \sigma^2 \Big|_{m.e.m.} &= \sigma^2 \Big|_1 \quad a=0 \\
 &< \sigma^2 \Big|_1 \quad a \neq 0.
 \end{aligned}$$

(iii) Several Parameters

As we have already mentioned, the case of several parameters is treated by taking several moments. Now our moment errors are described by a full second moment matrix for off diagonal terms

come from $\int (f(x) - \langle f \rangle)(g(x) - \langle g \rangle)p(x)dx$
 which is non-zero in general.

It is conventional to ignore correlations:
 the expansion in terms of $X_m^l(\theta, \varphi)$ of
 dσ/dt or an angular decay distribution
 is typical. The expansion functions
 $a_{lm} Y_m^l$ are orthonormal, so

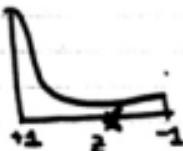
$$W(\theta, \varphi) = \sum_m a_{lm} Y_m^l$$

$$a_{lm} \propto \langle Y_{lm}^* \rangle$$

but still a_{lm} have correlated errors

$$\text{as } \langle Y_{lm} Y_{l'm'} \rangle \neq \langle Y_{lm} \rangle \langle Y_{l'm'} \rangle$$

These correlations in HEP example
 of dσ/dt



express fact that

backward peak (or forward) comes
 from correlated cancellations.

In formula

$$\langle f \rangle = \frac{\sum_{i=1}^N f(x_i)}{N} \quad (*)$$

I used to worry if you should include counting error (\sqrt{N}) in N in error estimate which otherwise only involves the standard deviation of f . If N is a prior fixed, clearly there is no \sqrt{N} error but suppose the events used in (*) are those that fall in a particular subregion (counter) of a bigger sample space (and $f(x)$ is defined over some only over this subregion). Then estimate of probability of falling into counter certainly involves \sqrt{N} error! What about mean of quantity within subregion?

The problem is : what is error in a moment : should it include counting error ? The answer is no - one uses identical formulae to case where number of evals prescribed and not subject to counting error.

Show this in a model problem in 2 variables : x is variable on which moment taken and y is binned. Choose:

$$P(x, y) = p(y) [1 + \alpha(y)x]/2$$

Take a y -bin $y_2 \leq y \leq y_2$

Set $s(y) = 1 : y_2 \leq y \leq y_2$

$= 0 : \text{otherwise.}$

Assume $\alpha(y)$ is constant at x over bin and then method of moments is

$$\alpha = 3 \frac{\sum_{i=1}^N f(x_i, y_i)}{\sum_{i=1}^N g(x_i, y_i)} = \frac{3F}{G}$$

$$\text{where } f(x, y) = x S(y)$$

$$g(x, y) = S(y)$$

Now we can apply F/G theorem from end of (D2) and it is not hard to show that f and g separately have counting errors \sqrt{n} but they are correlated and correlation cancels counting errors, so that α has identical error formula to Standard moment method. (The estimate of $p(y)$ will involve counting error of course).

D5: MinimizationD5/1 : Problem

The basic problem is to minimize $f(\alpha_1, \dots, \alpha_n)$ - a function of n variables $\alpha_1, \dots, \alpha_n$. Equivalently we must find the zero of an arbitrary function $f'(x_1, \dots, x_n)$ of n variables. This arises in previous section when we wish either to minimize χ^2 or maximize the likelihood L . (equivalently to minimize $-\ln L$).

The methods divide according to

- (i) Can calculate f only
- (ii) " " " f, f' and f'' .

(i) is extremely time consuming (e.g. it takes a lot of function evaluations to converge to desired minimum). In all problems I have faced; it is has been simple to

calculate derivatives and so we will only discuss (ii).

D5/2 Derivative Formalism

(i) One great advantage of the χ^2 or maximum likelihood problem is that to calculate the second derivative of the function to be minimized only requires the first derivative of the theoretical functions.

■ First consider this in the case of χ^2 :

$$\chi^2 = \sum_{i=1}^N (t_i - e_i)^2 / \sigma_i^2 \quad (D5.1)$$

where t_i is theory-function of n parameters - we take, wolog, $n=1$ and one parameter α in the following.

$$\frac{\partial \chi^2}{\partial \alpha} = \sum_{i=1}^N 2(t_i - e_i) \frac{\partial t_i}{\partial \alpha} / \sigma_i^2$$

$$\frac{\partial^2 \chi^2}{\partial \alpha^2} = S_1 + S_2$$

$$S_1 = \sum_{i=1}^N 2(t_i - e_i) \frac{\partial^2 t_i}{\partial \alpha^2} / \sigma_i^2$$

$$S_2 = 2 \sum_{i=1}^N \left(\frac{\partial t_i}{\partial \alpha} \right)^2 / \sigma_i^2$$

Now S_1 is $O(\sqrt{N})$: thus regard β as a random variable which is a function of the N random variables $e_1 \dots e_N$. By the Central limit theorem S_1 has mean $\mu = \sum_{i=1}^N \mu_i$ and standard deviation σ

$$\sigma^2 = \sum_{i=1}^N \omega_i^2$$

where $2 \frac{(t_i - \mu_i)}{\sigma_i^2} \frac{\partial t_i}{\partial \alpha^2}$ has mean μ_i and standard deviation ω_i . Now if t_i was evaluated with true theoretical parameters $\mu_i = 0$ and ω_i is some number. So σ^2 is $O(N)$ and S_1 $O(\sqrt{N})$ as desired. meanwhile S_2 is $O(N)$ as it is sum of N numbers. whence - near solution

$$\frac{\partial^2 \lambda^2}{\partial \alpha^2} = S_2 = 2 \sum_{i=1}^N \left(\frac{\partial t_i}{\partial \alpha} \right)^2 / \sigma_i^2$$

As claimed : Second derivatives
of χ^2 only require evaluation of
first derivatives of theory. (Only
true when N large)
(ii) Let's do the same thing for
maximum likelihood method.

$$-\ln L = -\sum_{i=1}^N \ln q_i$$

where q_i are probability of our
 $i=1 \dots N$ events. Assume q_i are normalized

$\int q_i(x) dx = 1$. This is crucial
but our subtle modifications in
e.g. $\Delta 4/3$ (iii) - "generalized" maximum
likelihood method can also be handled.

$$-\frac{\partial^2 \ln L}{\partial \alpha^2} = S_1 + S_2$$

$$S_1 = \sum_{i=1}^N \frac{1}{q_i} \frac{\partial^2 q_i}{\partial \alpha^2}$$

$$S_2 = \sum_{i=1}^N \frac{1}{q_i} \left(\frac{\partial q_i}{\partial \alpha} \right)^2$$

Now $S_{1,2}$ are directly of form suitable to apply central limit theorem : thus they are sums over independent random variables taken from same distribution.

$$\begin{aligned}\langle S_1 \rangle &= N \left\langle \frac{1}{q} \frac{\partial^2 q}{\partial x^2} \right\rangle \\ &= N \int_q \cdot \frac{1}{q} \frac{\partial^2 q}{\partial x^2} dx \\ &= N \frac{\partial^2}{\partial x^2} \left\{ \int q dx \right\} \\ &= N \frac{\partial^2}{\partial x^2} (1) = 0:\end{aligned}$$

So again $\# S_1$ only has standard deviation term and is $O(\sqrt{N})$: similarly S_2 is $O(N)$ and dominates. It follows that exactly same methods can be applied to χ^2 and maximum likelihood fits : in the future we will only mention the formula.

25/3 The method

lets first generalize the above to n parameters - getting a vector F_j of first and matrix m_{jk} of second derivatives

$$F_j = \frac{\partial \chi^2}{\partial \alpha_j} = \sum_{i=1}^N \frac{(t_i - e_i)}{\sigma_i^2} \frac{\partial t_i}{\partial \alpha_j} \quad (25.2)$$

$$m_{jk} = \frac{1}{2} \frac{\partial^2 \chi^2}{\partial \alpha_j \partial \alpha_k} = \sum_{i=1}^N \frac{1}{\sigma_i^2} \frac{\partial t_i}{\partial \alpha_j} \frac{\partial t_i}{\partial \alpha_k}$$

Note that - in our approximation - m is positive semi definite - which is also a technical advantage.

Now given any guess $\alpha = \alpha_0$

$$\chi^2(\alpha) = \chi^2(\alpha_0) + \sum_j (\alpha - \alpha_0)_j F_j$$

$$+ \sum_{jk} (\alpha - \alpha_0)_j (\alpha - \alpha_0)_k m_{jk}$$

Adopt a vector notation with $\Delta \alpha_j = \alpha_j - \alpha_{j0}$

$$\chi^2(\underline{\alpha}) = \chi^2(\underline{\alpha}_0) + \Delta \underline{\alpha}^T \underline{F} + \frac{1}{2} \Delta \underline{\alpha}^T M \Delta \underline{\alpha} \quad (D5.3)$$

minimum condition is

$$\frac{\partial \chi^2(\underline{\alpha})}{\partial \alpha_j} = 0 \quad \text{or}$$

$$M \Delta \underline{\alpha} = -\underline{F} \quad (D5.4)$$

so we have new guess

$$\underline{\alpha} = \Delta \underline{\alpha} + \underline{\alpha}_0 = \underline{\alpha}_0 - M^{-1} \underline{F}$$

for true minimum. Note that as M positive definite in (D5.4), this is indeed a minimum as long as Second order Taylor expansion valid.

(a) Error Calculation

when we are at a minimum, $\underline{F} = 0$ and (D5.3) becomes

$$\chi^2(\underline{\alpha}) = \chi^2(\underline{\alpha}_0) + \frac{1}{2} \Delta \underline{\alpha}^T M \Delta \underline{\alpha}$$

or remembering that the likelihood

$$\mathcal{L} = \exp(-\frac{1}{2} \chi^2)$$

$$\mathcal{L} = \mathcal{L}(\alpha) \propto \exp(-\frac{1}{2} \Delta_{\alpha}^T M \Delta_{\alpha}) \quad (05.5)$$

(05.5) is expected Gaussian form for \mathcal{L} . According to D4, M is just inverse of error matrix I^{-1} .

$$\langle (\alpha^* - \hat{\alpha})_j (\alpha^* - \hat{\alpha})_k \rangle = M^{-1}_{jk}$$

where α^* = estimate from minimizing χ^2
 $\hat{\alpha}$ = true parameters.

In particular "error" σ_j in α_j is

$$\sigma_j = \sqrt{(M^{-1})_{jj}} \quad (05.6)$$

and the other elements in M^{-1} are the much forgotten correlation parameters. Note that $\Delta \alpha_j = \sigma_j$ all other $\Delta \alpha_i = 0$, produces a χ^2 change of 1.

D5-4 The Real World

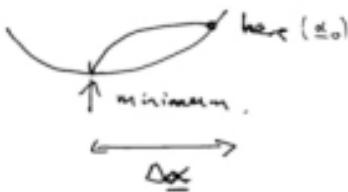
For a real problem, the iterative method described in D5-1 always diverges: one must invent various artistic devices to ensure procedure converges.

The simplest reliable algorithm is called Marquardt's method and this is intuitively easy to understand

Consider χ^2 as a function of several parameters.



χ^2 will have a bunch of hills and valleys and we are trying to discover the location of the lowest point in valley. If our assumptions are correct then χ^2 is just a nice quadratic function and we have found a step that exactly moves to the minimum.



Unfortunately, this is not usually correct as the step $\underline{\chi}_0 \rightarrow \underline{\chi}_0 + \Delta \underline{\chi}$ maybe too large to support the assumption, that terms in $\frac{\partial^k (\chi^2)}{\partial \underline{\chi}^k}$, $k > 3$ are zero, to be correct.

Consider the case when there are two variables.

$$\chi^2 = \chi_0^2 + \text{linear terms}$$

$$\quad \quad \quad + A_1 (\alpha_1 - \alpha_1^*)^2 + A_2 (\alpha_2 - \alpha_2^*)^2$$

$$\quad \quad \quad + B_1 (\alpha_1 - \alpha_1^*)^3 + B_2 (\alpha_1 - \alpha_1^*)^2 (\alpha_2 - \alpha_2^*) + \dots$$

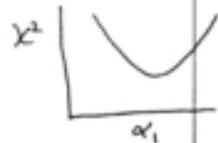
+ ...

where we have taken the special case where there is no term proportional to $(\alpha_1 - \alpha_1^*)(\alpha_2 - \alpha_2^*)$; Such a term only confuses the current discussion and we will soon see how to put it back in.

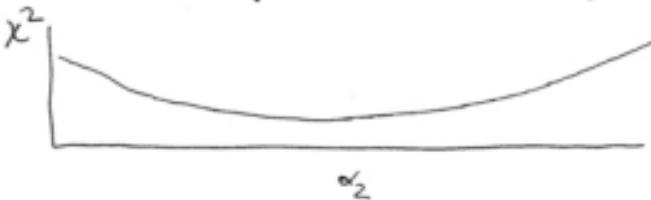
The difficulty (lack of convergence of iterative scheme) occurs when

$$A_1 \gg A_2$$

so that as a function of α_1 , we get



and as a function of α_2 , we get



with a valley with steep walls in α_1 and gentle walls in α_2 direction.

The predicted shift from eqn (B5.4), $\Delta\alpha_2 \sim \frac{1}{A_2}$ and is typically large compared to $\Delta\alpha_1 \sim \frac{1}{A_1}$.

Unfortunately, there is no reason to believe that this large shift is well estimated as assumes that lead to expansion break down.

Terms like $B_2 (\Delta\alpha_1)^2 \Delta\alpha_2$ are quite likely to be.

We see that estimate is false and one will typically jump out of valley & find

$$\chi^2 (\underline{\alpha}_0 + \Delta\underline{\alpha}) \gg \chi^2 (\underline{\alpha}_0)$$

The method will diverge!

In the above example, the estimate $\hat{\alpha}_1$ was valid and not spoilt by $\hat{\alpha}_2$ misestimate. In a real example, there will be mixing between $\hat{\alpha}_1$ and $\hat{\alpha}_2$ due to errors. ~~$\frac{\delta \hat{x}}{\delta \hat{x}_1, \delta \hat{x}_2}$~~ . Both shifts $\delta \hat{x}_1, \delta \hat{x}_2$ will be large and specimen. There is an intuitively obvious solution of this. Consider the second order expansion term.

$$(\hat{\alpha} - \hat{\alpha}_0)^T M (\hat{\alpha} - \hat{\alpha}_0)$$

M is a symmetric positive semi-definite matrix. M can be diagonalized

$$M \rightarrow P^T D P \quad \text{nxn matrix.}$$

$$D = \begin{bmatrix} \lambda_1 & & \\ & \ddots & 0 \\ 0 & & \lambda_N \end{bmatrix}$$

λ_j are eigenvalues of M
 P are eigenfunctions

$$\text{Put } \underline{\beta} = \underline{P} \underline{\alpha}$$

Then $\chi^2 = \chi_0^2 + \text{linear terms}$

$$+ \lambda_1 (\beta_1 - \beta_{01})^2 + \lambda_2 (\beta_2 - \beta_{02})^2 + \dots$$

Now the shifts $\Delta \beta_j$ can be calculated independently and do not affect each other.

Assume that we order the eigenvalues

$$\lambda_1 > \lambda_2 > \dots > \lambda_n > 0.$$

Then clearly the first few shifts $\Delta \beta_j$ corresponding to $\lambda_1, \dots, \lambda_L$, say are reliable and the remainder $j=L+1, \dots, n$ unreliable.

Suppose we ignore the predictions for these last and set

$$\begin{aligned}\Delta \beta_j &= \text{Predicted value } j \leq L \\ &= 0 \quad j > L.\end{aligned}$$

We need a criterion for choosing L .
 But subject to this, this method is intuitively reasonable.

In a program that I used for many years successfully, I chose L by the largest value such that

$$\lambda_L \geq r\lambda_1 \text{ and } \lambda_{L+1} < r\lambda_1.$$

where initially $r=0.1$ and r was increased or decreased based on experience i.e. if the resultant $X^T(\beta_0 + \alpha\beta)$ was indeed smaller (as predicted) or larger than $X^T(\beta_0)$.

This method is equivalent to replacing M by the regularized version.

$$M = P^T \begin{bmatrix} \lambda_1 & & & 0 \\ & \lambda_2 & & \\ & & \ddots & 0 \\ 0 & & & \lambda_L \\ & & & 0 \end{bmatrix} P$$

I note that in using this method, I got to look at values of the λ_j . It was not uncommon for

$$\lambda_{\max}/\lambda_{\min} = \lambda_1/\lambda_n \sim 10^7 \text{ for } n \sim 30$$

You can ask, why one gets such large ratios. Intuitively the eigenvectors of M with small eigenvalue correspond to linear

combinations of a_j that are badly determined so very small λ 's say that your theoretical formalism is redundant or alternatively the data inadequate.

For instance in fitting a cross section $d\sigma/d(\cos\theta)$ to $\sum_{j=1}^n a_j P_j(\cos\theta)$

then often the coefficients of the higher order Legendre functions are poorly determined and if you choose n to be "too big", then small eigenvalues will develop. In this case, one can cure the problem by reducing the top index n but in general no such simple solution exists.

Margaret's method is a simpler technique but with similar goals. One replaces

$$M \rightarrow M + \beta I$$

where I is the identity matrix.

So means we replace M by

$$M = P \begin{bmatrix} \lambda_1 + Q & & \\ & \ddots & \\ & & \lambda_n + Q \end{bmatrix} P$$

Clearly if $\lambda_L \sim Q$, then the two formalism's are roughly equivalent

Eigenvalues with eigenvalues $> Q$, are shifted by same prediction while those with eigenvalue $< Q$, have their shift drastically reduced.

We need an iterative method to estimate Q which again is varied based on experience. Q is increased when predicted $X^*(\beta_0 + \delta\beta) > X^*(\beta)$ and is otherwise decreased or left unchanged.

The natural "Scale" for Q is given by

$$\frac{1}{n} \text{Tr } M = \frac{1}{n} \sum_{j=1}^n \lambda_j$$

the average eigenvalue.

Marquardt's method is less precise than the diagonalization method but much faster; finding eigenvalues and vectors is time consuming.

I have used both successfully

Diagonalization method for "large" individual fits with many parameters.

I used Marquardt's method in the analysis of some 10^6 femtobolt events where in each event I fitted known source shapes to measured photon energy measurements in a Segmented Lead-Sцинillator Sandwich. Here $n \approx 10$ and the large number of fits (10^6) required a fast method.

§7 χ^2 , Bartlett's S function

When the maximum likelihood method reduces to χ^2 , it is relatively easy to judge goodness of fit.

$$\chi^2 = \sum_{i=1}^n \left(\frac{e_i - t_i}{\sigma_i^2} \right)^2 \quad (*)$$

where t_i depend on m theoretical parameters. Then the result of minimizing (*) wrt m theoretical parameters should be distributed according to χ^2 distribution with $n-m$ degrees of freedom. This function is plotted on p. 64 of EDJRS. If k is number of degrees of freedom.

$$p(\chi^2) = \frac{\left(\chi^2/k\right)^{k/2-1}}{2^{k/2} \Gamma(k/2)} \exp(-\chi^2/k) \quad (*)$$

$$\langle \chi^2 \rangle = k$$

and Standard deviation is $\sqrt{2k}$.

Consider χ^2 with one degree of freedom

$$\chi^2 = \frac{(x-t)^2}{\sigma^2}$$

$$\text{where } p(x) = \frac{1}{\sigma\sqrt{2\pi}} \exp\left(-\frac{(x-t)^2}{2\sigma^2}\right)$$

$$\begin{aligned} p(\chi^2) &= p(x) \frac{dx}{d\chi^2} = \frac{p(x) \sigma^2}{2(x-t)} \cdot 2 \\ &= \frac{\exp(-\chi^2/2) \cdot 2}{\sqrt{2\pi} 2 \sqrt{\chi^2}} \end{aligned}$$

↑
gives same
value of χ^2

which is value on (*) using $B(1/2) = \sqrt{\pi}$

One can derive χ^2 distribution for k degrees of freedom in two ways

- i) Matthews and Walker : direct calculation as above noting that $\sum_{i=1}^k x_i^2 - k$ means 0, std. deviation 1 is τ^2 where τ is radial distance in k dimensional space. Then factors multiplying $\exp(-\chi^2/2)$ in (*) are just volume element in k dimensions
- ii) Eade et al. Show characteristic function of χ^2 with one degree of freedom is

$$\int_0^{+\infty} \exp(it\chi^2) p(\chi^2) d(\chi^2) = \sqrt{1-2it}$$

∴ characteristic function for χ^2 with

k degrees of freedom is $(1 - 2/t)^{-k/2}$ and this can be inverse Fourier transformed to give (*).

Introduce the characteristic function $F(\chi^2)$ by

$$F(\chi^2) = \int_0^{\chi^2} p(\chi^2) d(\chi^2) = \frac{1}{2\Omega(k)} \int_0^{\chi^2} (u_{k/2})^{(k/2)-1} e^{-u/2} du$$

F has probability distribution $w(F)$

$$w(F) = p(\chi^2) \frac{d(\chi^2)}{dF} = 1$$

so F is uniformly distributed between 0 and 1.

We prefer to define $C(\chi^2) = 1 - F(\chi^2)$ which is called the confidence level

$C(\chi^2) = 1$ is $\chi^2 = 0$ "good"

$= 0$ is $\chi^2 = \infty$ "bad"

or "low confidence"

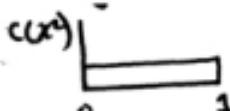
Suppose you conduct a series of experiments which gives you a sequence of χ^2 values

$$\chi^2 |_{1,2} \dots$$

An example in high energy physics would come from analysing a set of events each containing particle tracks

- each track will have an individual χ^2 . A given experiment would generate many million such χ^2 values!

If one can do the "cuts" right, the resulting distribution in $c(\chi^2)$ should be flat

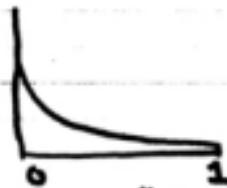


but if you have "background"
i.e. events that do not satisfy
claimed hypothesis, then one will
have an excess of events at
 $\chi^2 \gg k$ or small $c(\chi^2)$

$C(\chi^2)$ will ~~randomly~~ look like



The plot of $C(\chi^2)$ is a good way of estimating background. Unfortunately one often gets $C(\chi^2)$ plots that look like



with no clear "flat" piece. This maybe because your "signal" is very small. It could be because you got your errors wrong (!) e.g. maybe the true σ_i is are

all larger than ones used in calculating χ^2 . In this case it makes sense to scale χ^2 ($\chi^2_{\text{new}} = c \chi^2_{\text{old}}$) until we get a nice flat distribution for $C=1$ and analyse as on top of g^{19} .

If we expand χ^2 near its minimum - with for convenience - just one theoretical parameter.

$$\chi^2 = \chi_0^2 + \alpha(\lambda - \lambda_0)^2 \quad (**)$$

where λ_0 is our estimate of λ . Then probability distribution of λ is just

$$\propto \exp(-\chi^2/2)$$

$$\propto \exp\left(-\frac{\alpha(\lambda - \lambda_0)^2}{2}\right)$$

i.e. $1/\sqrt{\lambda}$ is standard deviation of estimates. Returning to (**), we find that if we ~~want~~ put $\lambda = \lambda_0 \pm 1/\sqrt{\lambda}$ (i.e. at its one standard deviation value), χ^2 is increased by 1. This χ^2 increase "rule" is common method of finding "error" in estimate λ even when the Taylor expansion (**) is not a good approximation.

This can be generalized to Bartlett's S function which is defined as

$$S(\lambda) = \frac{1}{\sqrt{I(\lambda)}} \frac{\partial \ln L}{\partial \lambda}$$

$I(\lambda)$ is information $\langle \left(\frac{\partial}{\partial \lambda} \ln L \right)^2 \rangle$

$$I(\lambda) = N \int \left[\frac{\partial}{\partial \lambda} \ln P(x, \lambda) \right]^2 P(x, \lambda) dx$$

Can be calculated
from theoretical forms.

Number of events.

From (**) what follows we have

$$\frac{\partial}{\partial \lambda} \ln L = -\frac{2\alpha}{2}(\lambda - \lambda_0) \quad (L = \exp(-\frac{x^2}{2}))$$

and using our earlier ~~asymptotic~~
result $\alpha = \sqrt{I(\lambda)}$, we find

$$\begin{aligned} S(\lambda) &= \sqrt{I(\lambda)} (\lambda_0 - \lambda) \\ &= \sqrt{\alpha} (\lambda_0 - \lambda) \end{aligned}$$

$\therefore \lambda = \lambda_0 \pm 1$ standard deviation
is just $S(\lambda) = \pm 1$. In fact
for any n , $S(\lambda)$ has mean 0
and variance 1 and so the
rule $S(\lambda) = \pm 1$ ($\pm l$ for l standard
deviations) is an appropriate
generalization of that to increase
 χ^2 by $\pm 1(l)$.

§8 Hypothesis Testing

We wish to test a specific hypothesis H_0 (referred to as null hypothesis) against a set of possible other hypotheses. In goodness of fit analysis, we test H_0 (= hypothesis that data described by claimed theory with parameters found from maximizing L) against "all" other hypotheses. If we are trying to find type of a particle then H_0 might be particle is an electron, H_1 that is a pion and there could be only possibilities. We are given

our Standard N observations
 (x_1, \dots, x_n) on which to make our decision. We do this by constructing yet another statistic Y which is some function of x_1, \dots, x_n . We assume that Y has distribution $P_Y^{(1)}$ if hypothesis 1 is true ($1 = 0$ or 1 if 2 hypotheses). Y could be vector valued but we will take scalar case here. In the rest of all possible worlds, one can find a Y such that the values taken by it under two hypotheses are disjoint i.e.

$$\int P_Y^{(0)}(y) P_Y^{(1)}(y) dy = 0 \quad (*)$$

Unfortunately (*) is not always attainable but one should always try to minimize the probability overlap integral. We define τ as region of y in Y Space where we will deem ~~H_0~~ hypothesis (i) to be true and (ii) to be false. But in the remaining region of Y Space, $R-\tau$, (R total Y range) we will take (i) to be true, (ii) to be false. This scheme is characterized by 2 constants:

$\alpha = \underline{\text{level of significance or size}}$

$$\text{of test} = \int_{\tau} p_Y^{(0)}(y) dy$$

α is chance that (i) will be

declared false when it is in fact true.

$\beta = \int_{R-r}^{\infty} p_y^{(1)}(y) dy$ is the contamination

i.e. chance that ω will be declared true when it is in fact false. Clearly it depends on the application without the balance between efficiency $1-\alpha$ and contamination β .

g) Student's t distribution

x_i are drawn from a normal distribution of unknown mean and standard deviation.

$H^{(0)}$ is mean = μ^0

$H^{(1)}$ is mean = $\mu^{(1)}$

$$\text{Define } t = \frac{(\bar{x} - \mu^{(1)})\sqrt{n}}{\bar{\sigma}}$$

$$\text{where } \bar{x} = \frac{1}{n} \sum_{i=1}^n x_i$$

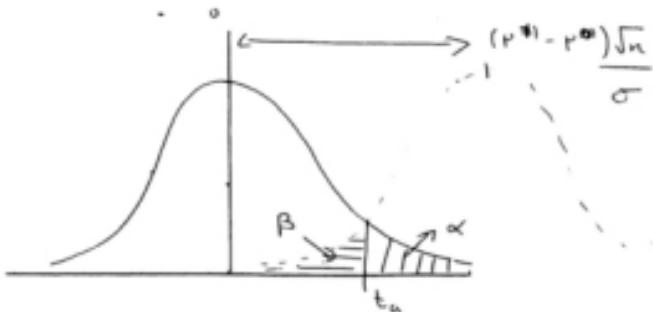
$$\bar{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})^2$$

Under hypothesis (i), t follows a so called Student's t -distribution with $n-1$ degrees of freedom. ($f = n-1$). The probability distribution

$$p(t) = \frac{\Gamma(\frac{f+1}{2})}{\sqrt{\pi} \Gamma(\frac{f}{2})} \frac{1}{(1+t^2/f)^{(f+1)/2}}$$

$f = 1$ is Cauchy distribution

$f \rightarrow \infty$ it becomes gaussian with mean 0 and standard deviation 1.



One could say that

0 is true if $t < t_a$

1 is true if $t > t_a$.

The corresponding efficiency / contamination are hardly calculable from cumulative probability functions of the Student's t distribution. As discussed on page 168 of Brandt,
(we average \bar{G})
 this method can be adapted to test if 2 sets of measurements have the same mean (e.g. Set 1: Smokers Set 2: non smokers)

g16 Robust Estimation

See Andrews et al. "Robust Estimates of Location", Princeton University Press
"Robustness in Statistics", edited by R.L. Launer, G.N. Wilkinson, Academic Press (1979).

In the maximum likelihood method we found the method that had the smallest possible standard deviation (asymptotically at least). However as we emphasized in the Study of non parametric ~~estimates~~ estimates, it is quite possible (and indeed probable) that systematic error from incorrect theoretical form of probability is large. Even when we know main functional dependence (Bret Wigner for a resonance, say) there are often "small" contributions to the

probability distribution (e.g. background resonance) which must be estimated as polynomial/phase space)

which must be parameterized in an ad hoc fashion. It is hard to evaluate distortion due to effect of an incorrect parameterization of these unknown contributions. There is thus a great deal of interest in the study of so called robust techniques which are insensitive to errors in our knowledge about the underlying distribution. Statisticians are very interested in such techniques but (apart from non parametric methods we have already discussed) there has been little use/experience of them.

the type of
unhard problems that physicists
are want to find. Further the derivation
of robust methods is rather intuitive
and not as mathematically precise
as those of classical statistics.

The first book cited has an
instructive study of a very simple
problem; give a univariate sample
 $\{x^{(i)} : i=1, N\}$ from some distribution,
find an estimate of the mean μ of
the underlying probability distribution.
The ~~simple~~ natural method is to
estimate μ from the arithmetic mean

$$\mu_{\text{AM}}^* = \frac{1}{N} \sum_{i=1}^N x^{(i)}$$

As long as the underlying distribution
has a finite standard deviation this
will have an error $\propto 1/\sqrt{N}$ which

$\rightarrow 0$ as $N \rightarrow \infty$. Further the method is unbiased for any N

$$\text{i.e. } \langle p_{\text{est}}^* \rangle = \mu.$$

If the $x^{(i)}$ are distributed according to either the Gaussian

$$p(x) \propto \exp(-\frac{(x-\mu)^2}{2\sigma^2})$$

or exponential form

$$p(x) \propto \exp(-x/\mu)$$

then we know from maximum likelihood method that the arithmetic mean is the best method. However for other distributions, the arithmetic mean is unlikely to be best. For a trivial example, suppose I measured $y = x^2$ where x had an exponential distribution. Then the best estimate of the mean of

y is obviously

$$y_{\text{mean}}^* = \left(\frac{1}{N} \sum_{i=1}^N (y^{(i)})^{1/2} \right)^2$$

and not $y_{\text{mean}}^* = \sqrt{\frac{1}{N} \sum_{i=1}^N y^{(i)}}$

The problem of estimating μ is hardest for distributions $p(x)$ that have ~~the~~ long tails. Then with a small probability we will generate events from the tail but these will be way off the mean and if included in calculation of mean will distort results in a finite sample. One method that is more robust than the arithmetic mean, is the trimmed mean.

$$\mu_{\text{trim}}^* = \frac{\sum_{i=pN+1}^{N(1-p)} x^{(i)}}{N(1-2p)}$$

where ρN events are left off each end. As long as ρ is fixed < 0.5 , for the special case of distributions μ_{trim}^*

Symmetric about their mean, the trimmed μ^* will also have mean μ and error $\propto \frac{1}{\sqrt{N}}$. For fixed ρ the variance of μ_{trim}^* will be larger than μ_{AM}^* for Gaussian-like distributed $x^{(i)}$ but as you investigate distributions broader than Gaussians you will find μ_{trim}^* becoming smaller than that of μ_{AM}^* . A special case of μ_{trim}^* is the median (i.e. middle value of $x^{(i)}$ when they are numerically ordered) gotten with $\rho = 0.5$.

An interesting class of robust estimators are based on the modification of the maximum likelihood equations.

Thus let $z^{(i)} = \frac{x^{(i)} - \mu}{\sigma}$

and suppose that $p(x, \theta)$ is just a function of z . (True for Gaussian case!). The likelihood equation for μ is

$$\frac{\partial}{\partial \mu} \ln L = 0$$

$$\text{or } \sum_{i=1}^N \frac{\partial}{\partial \mu} \ln p\left(\frac{x^{(i)} - \mu}{\sigma}\right) = 0$$

Now - as we are assuming p to be an even function of z i.e. symmetric about the mean - we write the above as

$$\sum_{i=1}^N f\left(\frac{x^{(i)} - \mu}{\sigma}\right) = 0 \quad (*)$$

where f is an odd function of its argument.

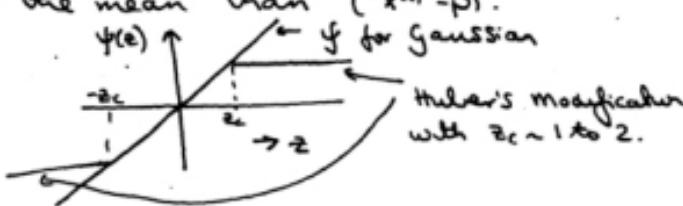
For the case of Gaussian p , f is just $\frac{1}{\sigma^2} (x^{(i)} - \mu)$ and $(*)$ is

$$\frac{1}{\sigma^2} \sum_{i=1}^N (x^{(i)} - \mu)^2 = 0$$

which gives the arithmetic mean as the maximum likelihood estimate

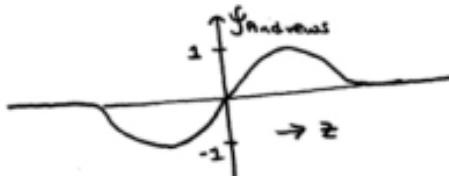
of the mean. Now we invent a function of which is less sensitive to $x^{(i)}$'s far from the mean than $(x^{(i)} - \mu)$.

e.g.



Many forms for f have been proposed. The most successful was that due to Andrews with

$$f(z) = \begin{cases} \sin(z/2\pi) & |z| < 2\pi \\ 0 & |z| \geq 2\pi \end{cases}$$



Note z is $(x - \mu) / \sigma$ where σ supplies natural scale for problem. We do not need a very good estimate for σ but

it should be robust! A simple choice is $\sigma = \text{median } |x^0 - \mu|$ which is applied iteratively as value of μ is refined.

A large number of tests are ^{described} ~~done~~ in just reference to this section.

Let $N = \text{Gaussian Distribution} - \sigma = 4$

$C = \text{Cauchy}$

$M = \text{Gaussian + uniform random variable between 0 and 1.}$

Here is a sample of their results for 10 observations

Method	N	C	M	← Distribution
mean	1	27.314	5.255	
median	1.37	3.66	7.33	
Andrews	1.08	4.66	9.29	

The above table gives variances in units of which I am not certain; but the relative values within a column are independent of this! We see very clearly that the arithmetic mean is hopeless for

broad distributions. The median/are very good in these cases ~~but~~ with the Andrews being better in Gaussian case (or only 4% more than best estimate) but less robust in other examples.

Note that all the methods except the arithmetic mean require all the $x^{(i)}$ being in memory so that we can order them. So the computer complexity is much higher for the robust estimates. The book gives computer programs for these cases all the robust estimators.

Ph 129 - 1990

G. Fox

MONTE CARLO INTEGRATION

This expands discussion on page 23 (Section (iii) of Section 1.7) of type written notes

E 5/4 Monte Carlo Integration

- (i) As we will find out later in the course (when we discuss Simpson's Rule and Gaussian Integration), ALL integration methods can be written as

$$I = \int_a^b f(x) dx = \sum_{i=0}^n f(x_i) w_i$$

The $n+1$ points x_i and weights w_i are determined

E5/4 Monte Carlo Integration

(ii) "All" integration methods can be written as

$$I = \int_a^b g(x)f(x)dx = \sum_{i=0}^n w_i f(x_i)$$

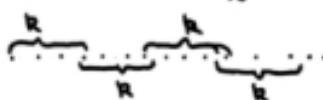
where the $n+1$ points x_i and $n+1$ weights w_i are determined

in sundry ways to optimize various properties -as described earlier. If we increase n , the dependence of error is

Error $\propto 1/n^2$ iterated Trapezoidal

$1/n^4$ iterated Simpson

$1/n^m$ iterated m -order Gauss



$k=2$ Trapezoidal.

$k=3$ Simpson.

$k=m$ Gauss.

with iterated groups of k points

Taking $g(x) = 1$ for convenience, we write

$$I = (b-a) \int_{-\infty}^{+\infty} f(x) p(x) dx$$

$$p(x) = \frac{1}{(b-a)} \quad a \leq x \leq b$$

$$= 0 \quad \text{outside} \quad [a, b]$$

Now consider the Monte Carlo which writes $I = \langle \tilde{f} \rangle$

$$\langle \tilde{f} \rangle = \int_a^b \tilde{f}(x) dx / b-a \quad \tilde{f} = (b-a)f$$

I is mean of $f(x)$ where x is uniformly distributed random variable between a and b .

Defining

$$\sigma^2 = \int_a^b [\tilde{f}(x) - I]^2 dx / b-a$$

we know from central limit theorem that

$$\frac{1}{n} \sum_{i=1}^n \tilde{f}(x_i) \rightarrow I \quad (\text{E5.ii})$$

with error σ/\sqrt{n} .

(E5.ii) is of the same form as our previous formulae : all the w_i are equal, while the positions x_i are random. However it seems a bad idea as $\text{error} \rightarrow 0$ much slower than even iterated trapezoidal.

But wait! Consider an integral in k dimensions - for other methods, we must make a grid in each direction with $n^{1/k}$ points for each coordinate. Thus error for Trapezoidal is $\propto 1/n^{1/k}$ which has a slower n dependence as k increases. However Monte-Carlo method still keeps its σ/\sqrt{n} behaviour and for $k > 5$, converges faster than Trapezoidal method. The other methods appear better than Monte-Carlo until larger k . However, in practice, this is not so because there the implicit

condition $n^{1/k} \gg 1$ before the error formulae even apply. For instance $n^{1/k} = 10$, $k = 5$ gives ridiculous values of n .

We cannot study it here, but this example shows that the problem of multidimensional integrals is not a trivial generalization of that in one dimension. You can pose general problem of test method in k dimensions : the answer is surely not either Monte-Carlo or Gaussian for each co-ordinate. There is a feeble account in Saeson and Keller. (chapter 7, p.352)

iii Complicated Domains

One advantage of the Monte Carlo method is that it can ~~not~~ very easily take account of complicated integration domains.

Suppose we wished to calculate the integral



over the inside of the circle

$$I = \int_{x^2+y^2 \leq 1} g(x,y) dx dy.$$

One way of doing this is to put

$$\begin{aligned} g'(x,y) &= g(x,y) & x^2+y^2 \leq 1 \\ &= 0 & x^2+y^2 > 1 \end{aligned}$$

and calculate

$$I = \int_0^1 dz \int_{-1}^1 dy g'(x,y)$$

which is an elementary Monte-Carlo integral with x and y generated independently between -1 and 1.

Actually the brute force method (generate x uniformly between -1 and 1, y uniformly between $\pm\sqrt{1-x^2}$) is more efficient - in the language of (5.11) the second method has a smaller variance.

However this is an important technique for I is often impossible to find ranges analytically while

It is easy to decide if a given (x, y) lie within range of integration.

(iii) Importance Sampling

of course it is obvious from (ES.ii) that monte Carlo integrator is better, the smaller is σ . Using definition of σ^2 , (or remembering its name "variance") it follows that method is best when $f(x)$ is slowly varying over domain.

[for instance if $f(x) \equiv \text{constant}$, $\sigma=0$ and method is exact]. Importance Sampling ("sampling where integrand important") is a technique for reducing variation of $f(x)$.

$$\text{Thus if } I = \int f g(x) dx$$

Change variables $y = g(x)$

$$dy = g'(x) dx$$

$$\text{put } f \frac{g(x)}{g'(x)} = F(y)$$

$$\text{Then } I = \int F(y) dy$$

we now apply Monte Carlo method
generating y uniformly.

We try to choose $q'(x)$ so $f(x)/q'(x)$
has a smaller standard deviation
than $f(x)$.

Best choice

$$q'(x) = f(x)$$

when

- 1) $F(y)$ is constant with zero standard deviation
- 2) Probability of getting a particular x value is $\propto q'(x) = f(x)$
i.e. one "chooses" x values to congregate where function f is large

This is of course not realistic
as to find $q'(x) = f(x)$ requires

$$q(x) = \int_a^x f(x) dx$$

i.e. Solving integral

In general one chooses $q'(x)$ so
that

- a) $q'(x)$ "approximately" $f(x)$
- b) $\int q'(x)$ can be performed.

Example

$$f(x) = e^{-x} \otimes \text{a mess}$$

where mess slowly varying

choose $q'(x) = e^{-x}$

$$q(x) = 1 - e^{-x}$$

Then $F(y) = \text{"a mess"}$

One can apply Monte Carlo methods to Sums as well as continuous integrals.

e.g. Suppose we have

$$\int dx_a dx_b dx_c \dots \sum_{l=1}^2 f_l(x_a, x_b, x_c \dots)$$

Then generate

x_a with random number

x_b

x_c

Generate another random number

r $0 \leq r \leq 1$ uniformly

if ($r < 0.5$) $l = 1$

else $l = 2$

and obvious modifications for

other discrete summations

Accept/Reject method

We wish to generate x according
to $p(x)$ $x_1 \leq x \leq x_2$

$$p_{\max} = \max_x p(x)$$

Let w_{\max} be $\geq p_{\max}$ (any number)

Generate x uniformly in $[x_1, x_2]$ ←
" " " in $[0, 1]$

if $r > \frac{p(x)}{w_{\max}}$ ignore x
 and Start again

$\leq \frac{p(x)}{w_{\max}}$ accept x

Metropolis Method (Journal of Chemical Physics 21, 1087 (53))

Given a point x_i — initially chosen at random, choose a new point x_j as follows: choose x_{i+1} in the neighbourhood of x_i uniformly.

In the $\int p(x) dx$

if $p(x_{i+1}) > p(x_i)$, accept it,
and, set $x_{i+1} = x_{i+1}$

if $p(x_{i+1}) < p(x_i)$, generate a random number r uniformly in $[0,1]$,

if $r < p(x_{i+1})/p(x_i)$ set $x_{i+1} = x_{i+1}$
 $r \geq " "$ " $x_{i+1} = x_i$

Eventually this produces a set of points distributed "according to $p(x)$ ".

reasons:

- 1) The neighbourhoods $\Delta x - x_i \leq x_t \leq x_i + \Delta x$ eventually cover the whole Space
- 2) Take any two points x_i, x_j
Suppose $p(x_i) > p(x_j)$

$$\Pr(i \rightarrow j) \text{ is } \frac{p(x_j)}{p(x_i)}$$

$$\Pr(j \rightarrow i) \text{ is } 1$$

i.e. $\frac{\Pr(i \rightarrow j)}{\Pr(j \rightarrow i)} = \frac{p(x_j)}{p(x_i)}$

In "equilibrium"

$$N_i \Pr(i \rightarrow j) = N_j \Pr(j \rightarrow i)$$

$$\therefore \frac{N_i}{N_j} \propto p(x_i)/p(x_j)$$

- So we have generated x_i according to $p(x_i)$
i.e. exact importance Sampling according

Philosophy of Standard Experimental Monte Carlo

One is evaluating integrals of the form

$$I(\Theta) = \int \Theta(x) p(x) dx$$

Here vector x is in

parameter space of experimental events e.g. in high energy physics
 x labels momenta of final state particles

$\rightarrow p(x)$ is probability that a given x
 is formed e.g. in high energy physics, $p(x)$ is

Square of modulus of

production amplitude. $\Theta(x)$ is a particular experimental observable e.g. energy into a calorimeter of a particular geometry,

Typically $p(x)$ is fixed by the theory but one wishes to calculate I for several choices of Θ as one explores different facets of the

experiment. Thus although the result of any given I depends on σ or of product Θp , to get good results for many different Θ 's ~~means~~ suggest that we minimize the variance of $p(x)$. There are three methods

1) use "importance Sampling": i.e. change variables to y where the

$$\text{Jacobian } \frac{\partial(x)}{\partial(y)} = J(x)$$

$$dx = J(x) dy \quad \text{and we try}$$

to choose $J(x)$ so that variance of $J(x)p(x)$ is minimized - the best choice is of course $J(x) = (p(x))^{-1}$.

2) Metropolis Method.

3) Uniformize "weights".

$$\text{Suppose } W = \max_x [p(x)]$$

- Then having generated x uniformly choose a random number r between 0 and 1

if $r < \frac{p(x)}{w}$ keep x

$r > \frac{p(x)}{w}$ reject it.

As events are accepted with probability $\propto p(x)$, they are "distributed according to $p(x)$ " -

i.e. $I(\Theta) = \frac{1}{N} \sum_{i=1}^N \Theta(x_i)$: x_i chosen as above.

This method is particularly good if it is very hard to calculate $\Theta(x)$. This happens say if response of particular part of apparatus is hard to calculate. (It might involve tracking particles through a magnetic field etc). Then it is "waste" to do a lot of work for events that

are only going to get a small weight
 $p(x)$ in the final analysis.

In hard problems, W can often not be precalculated and must be found from a preliminary examination of the integrand. Also the acceptance rate may be very poor e.g. a fraction $\langle p(x) \rangle / W$ is kept and this can be small ($.01$ is quite typical in problems I have worked on).

Lattice Monte Carlo

A typical "theoretical" calculation in condensed matter or high energy physics involves integrals of form

$$I = \int \exp[-H/kT] \otimes dS_1 \dots dS_m$$

where S_i is spin (gluon field) at site i . Typically the sites are in a regular lattice in 2, 3 or 4 dimensions. This integral would involve some $10^4 \rightarrow 10^6$ dimensions and so must be evaluated by Monte Carlo.

The Hamiltonian H is proportional to size of system and so the weighting function $\exp[-H/kT]$ is very small and one must be very careful to chose one integration points appropriately.

Assume you have a particular choice
 $S_1^{(0)} \dots S_m^{(0)}$ with

$$H^{(0)} = H(S_1^{(0)} \dots S_m^{(0)})$$

Then we try to move a small way from $S_i^{(0)}$ by just changing one of the spins

$$S_k^{(0)} \rightarrow S_k^{(1)}$$

$S_i^{(0)}$, $i \neq k$, unchanged.

Typically this produces a change $\Delta H \sim kT$ as kT is energy per site (upto numerical constants) - remember classical result of energy of $\frac{1}{2}kT$ per degree of freedom.

One usually one of two methods to change S_k

(a) metropolis - generate $S_k^{(1)}$ uniformly about $S_k^{(0)}$ and accept or reject change based on $r = \exp \left[- [H(\dots S_k^{(1)} \dots) / kT - H(\dots S_k^{(0)} \dots) / kT] \right]$

if $r > 1$ (energy decreased), accept $S_k^{(1)}$

if $r < 1$ (energy increased), accept $S_k^{(1)}$ with probability r .

(b) Heaviside - Sometimes one can use "exact" method to generate a variable according to a given probability distribution i.e. if distribution function

$$F_k(x) = \int dx_k \exp(-H(\dots s_k\dots)/kT)$$

can be calculated and inverted. Here we view H as a function of s_k with other spins fixed.

~

In each case (a) or (b), spins are generated according to distribution $\exp[-H/kT]$ and so values of observables are calculated as

$$\langle \Theta \rangle = \frac{1}{\# \text{Spin Configurations}} \sum_{l=1}^N \Theta(s_1^{(l)}, s_m^{(l)})$$

$= N$

If we had not been careful about generating the spin values

but say. choose them uniformly, then

$$\langle \Theta \rangle = \frac{\sum_{i=1}^N \exp[-H[S^{(i)}]/kT] \Theta(S^{(i)})}{\sum_{i=1}^N \exp[-H[S^{(i)}]/kT]}$$

This is correct but numerically hopeless. The wide variation in exponential (remember H proportional to size) will give a huge standard deviation in $\langle \Theta \rangle$.

Generation of Random Numbers

See

Knuth "The Art of Computer Programming"

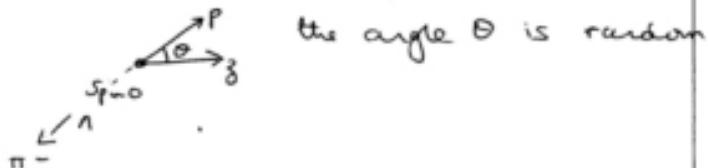
Volume 2, "Semi Numerical Algorithms"

Press "Numerical Recipes : The Art of Scientific Computing" - chapter 7

T.W. Chin Computer Physics Communications
Vol 47, 129 (1987).

Knuth is a reasonable pedagogical description although it is a little out of date - Chin describes newer algorithms which are imported in some cases. Press implements algorithms at about the level of Knuth.

One needs specialized hardware to generate true random numbers e.g.



for successive decays of a bunch of n 's!
There are a lot of difficulties here; not

only must you measure Θ in an unbiased fashion (impossible), to find a random number with 32 bits of precision would require several separate Θ values - maybe each gives 4 bits of information and so 8 values are needed!

Rather one usually a deterministic method which gives "essentially" random numbers. These are called pseudorandom. Most "standard" library random number generators use the so-called linear congruential generators which generate a sequence of integers I_j by

$$I_{j+1} = [a I_j + c] \bmod(m)$$

Here m is called the modulus
 a " " " multiplier
 and c " " " increment

This is deterministic - given any seed I_0 , one generates an identical sequence I_0, I_1, I_2, \dots

There are most m possible random numbers but there may be less. One could generate a cycle if

$$I_{J_1} = I_{J_2} \quad J_1 > J_2$$

Then $I_{J_1+1} = I_{J_2+1}$ etc.

also $I_{J_1-1} = I_{J_2-1}$

and in particular $I_{J_1-J_2} = I_0$.

If the cycle has a length M , then imagine our calculation of π

$$\pi = \int_{-1}^{+1} \int_{-1}^{+1} dx dy \Theta[1-x^2-y^2]$$

Then random numbers after the M^{th}
are no longer independent and so ~~the~~
the use of central limit theorem
breaks down.

Our error is $\propto \max \left[\frac{1}{\sqrt{N}}, \frac{1}{M} \right]$

)]
number cycle length
of events

The references give various formal criteria to guide the choice of a c and m. Unfortunately I don't see any of these references as being very good as they don't test and compare some key choices.

In my group we have used

$$m = 2^{31}$$

$$a = 1103515245$$

$$c = 12345$$

with effective C code

$$\text{newran} = \lceil a * \text{oldran} + c \rceil \text{ MASK}$$

where MASK has lower 31 bits = 1

all higher bits = 0.

We get a floating point number between 0 and 1 by:

$$\text{floatran} = \text{newran} / 2147483648.0$$

This is very simple to use especially in C which allows you to use machine capabilities effectively.

The Shift register methods are not discussed (correctly) in Knuth and may well be better as they are as easy to calculate and have longer period.

Chiu has implemented a particular Simple method.

$$I_j = I_{j-q} \oplus I_{j-p}$$

where \oplus = Exclusive OR

$$\begin{array}{l} 1+1=0 \\ 1+0=0+1=1 \\ 0+0=0 \end{array}$$

This \oplus operates independently on each bit of the numbers I_{j-q}, I_{j-p}

The period of this could be as large as $2^p - 1$ which is big for Chiu's choice $p = 250, q = 147$.

Initialization

The Shift register method is unpopular because it is newer and because it is harder to initialize! However Chiu describes a reasonable way using a congruential generator.

Typically a random number package

includes a call to a routine with a
name like ranset

A call ranset (seed)

Seeds $I_0 = \text{Seed}$ for the congruential method

If you want a truly random start,
a standard "dodge" is to set the Seed
to the binary image (remember I_0 is
just a (32 bit) integer) of the System clock.

One may wish to redo a calculation
with the same initial Seed so use

clock (seed)

print (binary, seed)

ranset (seed)

Correlations.

Often poorly chosen random
number generators exhibit correlations
between neighbouring random numbers

This is well illustrated by
tests by Chin of integer generator
Supplied by Borland with Turbo Pascal

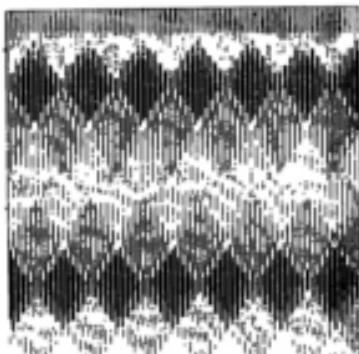


Fig. 2. More than 1000000 random numbers generated by Turbo-Pascal v3.0 and plotted as (x, y) pairs.

Adjacent random numbers considered as (x, y) pairs

Other generators are better!

Numerical Recipes has a program that removes any such problems

Essentially set up a list $R(i)$ of L random numbers.

Then choose $j \quad 1 \leq j \leq L$

Select desired random number = $R(j)$

Replace $R(j)$ by a new random number
~~(perhaps)~~

These correlations are particularly important in accurate calculations of condensed matter phase transitions as you are measuring correlations between local spins and so these must have uncorrelated random generators.

Gaussian Random Numbers:

$$g = \frac{\sum_{i=1}^N u_i - N/2}{\sigma \sqrt{N}}$$

where $\sigma^2 = 1/12$ and u_i are uniformly distributed between 0 and 1, is gaussianly distributed for large N . This for $N \approx 10$ is sometimes used in computers.

More convenient is often,

$$\begin{aligned} g_1 &= (-2 \ln u_1)^{1/2} \cos 2\pi u_2 \\ g_2 &= (-2 \ln u_1)^{1/2} \sin 2\pi u_2 \end{aligned}$$

where u_i are uniform in $[0,1]$ and then g_i are independent random numbers gaussianly distributed with zero mean and unit variance. You can prove this by taking polar co-ordinates in

$$\exp [-\frac{1}{2}(g_1^2 + g_2^2)] dg_1 dg_2.$$

$$g_1 = r \cos \theta$$

$$g_2 = r \sin \theta$$

$$\exp(-\frac{1}{2}[g_1^2 + g_2^2]) dg_1 dg_2$$

$$= \exp(-\frac{1}{2}r^2) r dr d\theta$$

$$= d(\exp(-\frac{1}{2}r^2)) d\theta$$

$$\therefore u_1 = \exp(-\frac{1}{2}r^2) \quad u_2 = \theta$$

are uniform

$\Rightarrow g_1, g_2$ Gaussian.

Calculations for General Distributions

- i) Example: exponential

Suppose x has distribution e^{-x} for $x = 0$ to ∞ .

$$e^{-x} dx = d(e^{-x})$$

$y = e^{-x}$ is uniform.

$\therefore x = -\ln y$ where y is a uniform random variable 0 to 1

- ii) More generally

$$p(x) dx = d(F(x))$$

Then $w = F(x)$ is uniformly distributed

$$w = \int_{-\infty}^x p(x) dx$$

* could be different here but if as above $p(x) = 0 \quad x \leq 0$

You must be able to invert integral to be able to do this

$$\text{For } p(x) = e^{-x}$$

$w = 1 - e^{-x}$ is essentially as above because $e^{-x} = 1-w$ as same as $e^{-x} = y$ as y and $1-w$ are both

uniform between 0 and 1.

iii) Interpolation:

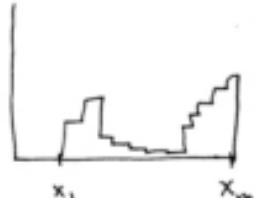
Suppose $p(x)$ is given by a table

$$x_1 \quad p(x_1)$$

$$x_2 \quad p(x_2)$$

$$\vdots$$

$$x_m \quad p(x_m)$$



in Say equally spaced intervals

$$x_k - x_{k-1} = \delta x.$$

Then you form cumulative distribution

$$x_i + \delta x \quad p(x_i) \quad F(x_i + \delta x)$$

$$x_i + 2\delta x \quad p(x_i) + p(x_L)$$

$$\vdots$$

$$x_{m+1} \quad \sum_{i=1}^m p(x_i) = 1 \text{ if normalized}$$

Then generate random number r
 $0 \leq r \leq 1$.

Then use inverse interpolation to find x

$$F(x) = r$$

Then x is distributed according to $p(x)$

(iv) Accept/reject

Then $p_{\max} = \max_x p(x)$ $x_1 \leq x \leq x_2$

Let w_{\max} be any number $> p_{\max}$

Generate x uniformly in $[x_1, x_2]$

Generate r uniformly in $[0, 1]$

If $r > \frac{p(x)}{w_{\max}}$, ignore x and start again

$\leq \frac{p(x)}{w_{\max}}$, accept x .

(v) One can combine this with Technique in (i).

Suppose $f(x)dx = d(F(x))$

and we can invert $F(x)$ as in (i)

Let w_{\max} be any number $> \max_x p(x)/f(x)$

Then generate x in $[x_1, x_2]$ distributed according to $f(x)$ using technique in (i)

for use (iv) with $p(x) \rightarrow p(x)/f(x)$

to "correct" $f(x)$ to $p(x)$

Clearly one chooses $f(x)$ to have

a distribution "near" $p(x)$ but with a simple enough form that one can find and invert it.

- (ii) Accept-reject is sometimes an important technique in event generation.

Suppose we generate a set N of events i with weight w_i representing the the production probability. Then if we want to calculate some simple quantity q_i , the best is to take

$$\langle Q \rangle = \frac{\sum_{i=1}^N q_i w_i}{\sum_{i=1}^N w_i} \quad (1)$$

An alternative is to use acceptance-rejection to find a set M of events with "weight 1".

$$m < N$$

then

$$\langle Q \rangle = \frac{1}{m} \sum_{j=1}^m q_j \quad (2)$$

when is (2) superior to (1).

If you are just calculating $\langle Q \rangle$, then (1) ~~uses~~ all the information you have and is "best".

However if you are doing a

lot of work to calculate q_i , or maybe there are many different quantities q_i , the accept-reject technique is best. Thus it is clearly a waste to spend a lot of time calculating q_i for an event whose weight w_i is small e.g. $0.001 w_{\max}$!

NATIONAL
23-680

NATIONAL INSERTABLE-TAB INDEXES ENABLE YOU TO
MAKE YOUR OWN SUBJECT ARRANGEMENT, USING PLAIN
INSERTS ON WHICH TO WRITE YOUR OWN CAPTIONS.

The Beaded edge on tab makes it easy to insert captions

Made in U. S. A.

Maximum likelihood is Best! (See Brandt p138-143, EDJRS chapters 5 and 7)

Consider a Statistic $S(x_1 \dots x_n)$ where x_i are independent random numbers drawn from the same distribution. Suppose we are using S to estimate some parameter λ characterising the distribution. We can assume that $\langle S \rangle \sim \lambda$ (if $\langle S \rangle \sim \lambda^3$ then replace S by $\sqrt[3]{S}$ etc.)

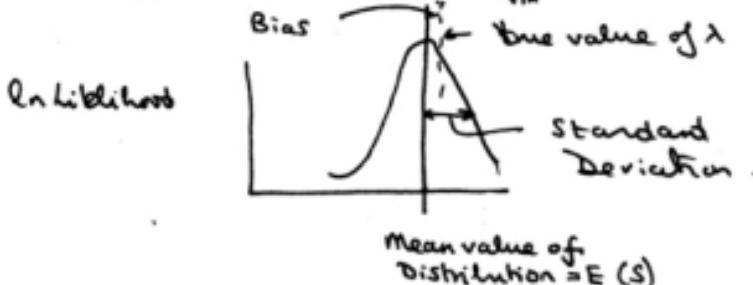
The bias

$$B(\lambda) = E(S) - \lambda$$

This should be distinguished from the variance $\sigma^2(S)$
 $\sigma^2(S) = \langle (S - E(S))^2 \rangle$

For the maximum likelihood method

$$B(\lambda) \sim \frac{1}{n} \quad \text{and} \quad \sigma(S) \sim \frac{1}{\sqrt{n}}$$



As the statistical error is \sqrt{n} times bigger than the bias, one usually ignores it. The small maximum likelihood method has asymptotically small bias and the minimum possible variance. We can only bound σ^2 variance for unbiased estimators (Statistics S) as for instance, the Statistic $S=0$ has zero variance but obviously a bad bias (unless $\lambda=0$!).

Let $P(x_i, \lambda)$ be probability distribution of the x_i .

$$\lambda + B(\lambda) = \int S \exp(\ln L) dx_1 \dots dx_n$$

where $L = \prod_{i=1}^n P(x_i, \lambda)$

Differentiating wrt λ , we get

$$\begin{aligned} 1 + B'(\lambda) &= \int S \left[\frac{\partial}{\partial \lambda} (\ln L) \right] \exp(\ln L) dx_1 \dots dx_n \\ &= \langle S \left[\frac{\partial}{\partial \lambda} (\ln L) \right] \rangle \end{aligned}$$

$$\begin{aligned} \text{as } \langle \frac{\partial}{\partial \lambda} (\ln L) \rangle &= \frac{\partial}{\partial \lambda} \int \exp(\ln L) dx_1 \dots dx_n \\ &= \frac{\partial}{\partial \lambda} \frac{1}{\prod_{i=1}^n P(x_i, \lambda)} \end{aligned}$$

by normalizer of probabilities \nearrow

$$= 0$$

we can replace S by $S - \text{any constant}$
or in particular $S - \langle S \rangle$

$$\therefore I + B'(\lambda) = \langle [S - \langle S \rangle] \frac{\partial}{\partial \lambda} (\ln L) \rangle \quad (9)$$

Now for any random variables α, β
we have $\langle \alpha \beta \rangle^2 \leq \langle \alpha^2 \rangle \langle \beta^2 \rangle$ $(*)$

This is ~~of course~~ obvious because for
any constants a, b

$$\langle (a\alpha + b\beta)^2 \rangle = a^2 \langle \alpha^2 \rangle + 2ab \langle \alpha \beta \rangle + b^2 \langle \beta^2 \rangle \geq 0$$

This implies $(*)$.

Applying this inequality to (9)
we have

$$[I + B'(\lambda)]^2 \leq \sigma^2(S) \langle \left[\frac{\partial}{\partial \lambda} (\ln L) \right]^2 \rangle$$

$$\text{or } \sigma(S) \geq \sqrt{\frac{I + B'(\lambda)}{\langle \left[\frac{\partial}{\partial \lambda} (\ln L) \right]^2 \rangle}} \quad (CR)$$

We set the denominator as $\sqrt{I(\lambda)}$
where the "information" $I(\lambda)$ is given by:

$$I(\lambda) = \int \left[\sum_{i=1}^n \frac{\partial}{\partial \lambda} (\ln P(x_i, \lambda)) \right] \left[\sum_{j=1}^n \frac{\partial}{\partial \lambda} (\ln P(x_j, \lambda)) \right] L(x_1, x_2, \dots, x_n) dx_1 \dots dx_n$$

The off diagonal terms are just

$$\langle \frac{\partial}{\partial \lambda} (\ln P(x_i, \lambda)) \rangle \langle \frac{\partial}{\partial \lambda} (\ln P(x_j, \lambda)) \rangle$$

$= 0$ by same argument as on the bottom of ML2.

$$\therefore I(\lambda) = n \underbrace{\int \left[\frac{\partial}{\partial \lambda} (\ln P(x, \lambda)) \right]^2 P(x, \lambda) dx}$$

This can also be written as

$$- \int \frac{\partial^2}{\partial x^2} (\ln P(x, \lambda)) P(x, \lambda) dx$$

$I(\lambda)$ is information on λ contained in measurement x_i . Notice that I has intuitively reasonable property (for name to be "information") of being

proportional to the number of observations.
 (CR) on bottom of MLS (the Cramer Rao bound) states that for an unbiased statistic

$$\sigma(S) \geq \frac{1}{\sqrt{I(\lambda)}}.$$

Relevance to maximum likelihood Method

The maximum likelihood statistic S is given by the equation

$$\begin{aligned} \frac{\partial}{\partial S} \ln \prod_{i=1}^n L(S, x_1, x_2, \dots, x_n) &= 0 \\ \text{or } \sum_{i=1}^n \frac{\partial}{\partial S} \ln P(x_i, S) &= 0. \end{aligned} \quad (1)$$

Expand $\ln P(x_i, S)$ about $S = \lambda$

$$\begin{aligned} \ln P(x_i, S) &= \ln P(x_i, \lambda) + (S - \lambda) \frac{\partial}{\partial \lambda} \ln P(x_i, \lambda) \quad (2) \\ &\quad + \frac{1}{2} (S - \lambda)^2 \frac{\partial^2}{\partial \lambda^2} \ln P(x_i, \lambda) + \dots \end{aligned}$$

we will return to this

Substitute (2) into (1), to find:

$$S - \lambda = \frac{\sum_{i=1}^n \frac{\partial}{\partial \lambda} \ln P(x_i, \lambda)}{\sum_{i=1}^n -\frac{\partial^2}{\partial \lambda^2} \ln P(x_i, \lambda)}$$

This is of our standard form

$$\frac{\sum_{i=1}^n f(x_i)}{\sum_{i=1}^n g(x_i)} \quad f = \frac{\partial}{\partial \lambda} \ln P(x, \lambda) \\ g = -\frac{\partial^2}{\partial \lambda^2} \ln P(x, \lambda)$$

$$\langle S - \lambda \rangle = \frac{\langle f \rangle}{\langle g \rangle} = 0 \quad (\langle f \rangle = 0 \text{ by } \text{argued on bottom of ML2}).$$

$$\begin{aligned} \langle (S - \lambda)^2 \rangle &= \frac{1}{n \langle g \rangle^2} \langle (f - \frac{\langle f \rangle g}{\langle g \rangle})^2 \rangle \\ &= \frac{\langle f^2 \rangle}{n \langle g \rangle^2} \end{aligned}$$

$$\text{but } \langle f^2 \rangle = \langle g \rangle = \frac{I(\lambda)}{n} \quad (\text{see middle of ML4}).$$

$$\therefore \langle (S - \lambda)^2 \rangle = \frac{1}{I(\lambda)}. \quad (**)$$

- i.e. the maximum likelihood method is
- Unbiased — See (*)
 - Optimal i.e. Variance = Cramer Rao bound (**).

This is only true in the limit $n \rightarrow \infty$.

let us examine this in more detail

Bias in Maximum Likelihood Method for finite n

There are two corrections to the formula $\langle S - \lambda \rangle = 0$ derived on the previous page. The first is the truncation involved in (2). In fact:

$$\begin{aligned} & \left[- \sum_{i=1}^n \frac{\partial}{\partial \lambda^2} \ln P(x_i, \lambda) \right] (S - \lambda) \\ &= \sum_{i=1}^n \frac{\partial}{\partial \lambda} \ln P(x_i, \lambda) + \underbrace{\frac{1}{2} (S - \lambda)^2 \sum_{i=1}^n \frac{\partial^3}{\partial \lambda^3} \ln P(x_i, \lambda)} \end{aligned}$$

Here we can replace $(S - \lambda)^2$ by $\frac{1}{n} I(\lambda)$

and $\sum_{i=1}^n \frac{\partial^3}{\partial \lambda^3} \ln P(x_i, \lambda)$ by $n \langle \frac{\partial^3}{\partial \lambda^3} \ln P(x, \lambda) \rangle$ with an error of $O(\frac{1}{n})$ as the each

of random variables $(S - \lambda)^2$ and $\sum_{i=1}^n \frac{\partial^3}{\partial \lambda^3} \ln P(x_i, \lambda)$ is peaked about its mean with error of $O(\frac{1}{n})$.

e.g. in our standard example

$$Fg = (n \langle f \rangle + \cancel{n \ln f}) (n \langle g \rangle + \cancel{n \ln g})$$

with $\langle \hat{f} \rangle = \langle \hat{g} \rangle = 0$.

$$\therefore Fg = n^2 \langle f \rangle \langle g \rangle + O(n \sqrt{n}) \text{ term with zero mean}$$

and $O(n)$ term with non-zero mean

Consider $\frac{f}{g}$ with $\langle f \rangle = 0$

$$\begin{aligned}\langle \frac{f}{g} \rangle &= \left\langle \frac{\sum_i f(x_i)}{n \langle g \rangle + \sum_i [g(x_i) - \langle g \rangle]} \right\rangle \\ &\sim \frac{1}{n \langle g \rangle} \left\langle \sum_i f(x_i) \left[1 - \frac{1}{n \langle g \rangle} \sum_i [g(x_i) - \langle g \rangle] \right] \right\rangle + \dots \\ &= \frac{-1}{n^2 \langle g \rangle} \underbrace{\left\langle \sum_i f(x_i) \sum_j [g(x_j) - \langle g \rangle] \right\rangle}_{\substack{\text{irrelevant as multiplies term with zero} \\ \text{only } i=j \text{ counts as } x_i \text{ independent}}} \\ &= \frac{-1}{n^2 \langle g \rangle^2} \left\langle \sum_{i,j} f(x_i) g(x_j) \right\rangle \\ &= -\frac{1}{n^2 \langle g \rangle^2} \langle fg \rangle\end{aligned}$$

i.e. when f and g are correlated with we get a $O(1/n)$ correction to mean.

Apply this to with $f = \frac{\partial}{\partial \lambda} \ln P(x, \lambda)$

$$g = -\frac{\partial^2}{\partial x^2} \ln P(x, \lambda).$$

we finally get

$$\begin{aligned} B(\lambda) = \langle (S - \lambda) \rangle &= + \frac{n}{2I^2(\lambda)} \left\langle \frac{\partial^3}{\partial \lambda^3} \ln P(x, \lambda) \right\rangle \quad (*) \\ &+ \frac{n}{I^2(\lambda)} \left\langle \frac{\partial}{\partial \lambda} \ln P(x, \lambda) \frac{\partial^2}{\partial \lambda^2} \ln P(x, \lambda) \right\rangle \end{aligned}$$

as $I(\lambda) \propto n$, we see that $B(\lambda) \sim \frac{1}{n}$.
The above can be written in a different
form using

$$\begin{aligned} \frac{\partial}{\partial \lambda} \left\langle \frac{\partial^2}{\partial x^2} \ln P(x, \lambda) \right\rangle &= \frac{\partial}{\partial \lambda} \int \frac{\frac{\partial^2}{\partial x^2} \ln P(x, \lambda)}{P(x, \lambda)} dx \\ &= \left\langle \frac{\partial^3}{\partial x^3} \ln P(x, \lambda) \right\rangle + \left\langle \frac{\partial}{\partial \lambda} \ln P(x, \lambda) \frac{\partial^2}{\partial x^2} \ln P(x, \lambda) \right\rangle \end{aligned}$$

but I think (*) is more convenient.

Note that all the terms in (*) can
be estimated from a given data sample
and so bias is ^{partly} corrected for. As $B(\lambda)$
 $\sim \frac{1}{n}$ and estimate has statistical error $\sqrt{\frac{1}{n}}$,
one can in fact find $\hat{\lambda}$ with a $\hat{\sigma}_{\lambda}$

with a bias $\sim \frac{1}{n\sqrt{n}}$ in a given experiment. Normally one does not do this because statistical error is of $O(\frac{1}{\sqrt{n}})$ and outweighs bias.

Example:

$$\text{if } P(x, \mu, \sigma) = \frac{1}{\sqrt{2\pi}\sigma} \exp[-(x-\mu)^2/2\sigma^2]$$

where x is observed and we wish to find μ, σ then maximum likelihood estimates are

$$S(\mu) = \sum_{i=1}^n x_i / n$$

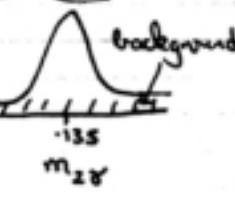
$$S(\sigma^2) = \sum_{i=1}^n (x_i - S(\mu))^2 / n$$

but in fact we know that $S(\sigma^2)$ is biased and unbiased formula is $\frac{n}{(n-1)}$ times $S(\sigma^2)$. This correction is $O(\frac{1}{\sqrt{n}})$ by above theorem and indeed can be calculated from it. (after we do trivial extension to ~~several~~ a vector of theoretical parameters $\underline{\lambda}$).

Before going on to goodness of fit tests, we will discuss non parametric representations of data. We will restrict ourselves here to the one dimensional case. As we will later see the multidimensional case is often reduced to this!

A very good reference is chapter 4 of SLAC-176, 1974 notes of Jerry Friedman at a CERN Computer School.

Non parametric Data Representation

Given a mass spectrum, 

Given a parametric summary of this is to assume that the number of counts (probability) as a function of mass can be fitted

to a Gaussian plus a polynomial background.

$$p(m) = c \exp [-(m-m_0)^2/2\sigma^2] + \sum_{i=0}^N a_i (m-m_0)^i$$

i.e. the data is summarized by the $4+N$ parameters $c, m_0, \sigma, \{a_i\}$. often this is the best approach but it is only possible if you know an appropriate form for $p(m)$. Even in the above case, we had to guess that smooth background could be adequately parameterized by a polynomial (of degree N chosen by trial and error).

A non parametric representation of this data is a histogram i.e.

i.e. divide mass range into equal intervals of size Δm and count the number n_i of counts in the i th bin $m_i + (i-1)\Delta m \rightarrow m_i + i\Delta m$ for $i = 1 \dots (m_2 - m_1)/\Delta m$. As long as the bin size is small compared with the spread in the data, the numbers n_i contain all the information of the data (see D4140). In a way that is clumsier but free of ad hoc (and possibly wrong) assumptions about the functional form of $p(m)$.

Let us now approach the problem more formally: We have measurements $\{x_i \mid i=1..N\}$ where the x_i are independent random

variables distributed according to some unknown probability density function $p(x)$. We wish to form estimates $\hat{p}(x) = \phi_N(x_1, \dots, x_N)$ of $p(x)$.

§1: Histogram method

Define $g_i(x)$ to be 1 in i th bin and zero elsewhere. $\mathbb{1}(g_i(x))$ is called an indicator function.
 i.e. if range ~~a~~ $a \rightarrow b$ and
 $\Delta x = (b-a)/m$, then

$$g_i(x) = \begin{cases} 1 & a + (i-1)\Delta x \leq x \leq a + i\Delta x \\ 0 & \text{elsewhere} \end{cases}$$

Then we define our estimate of $p(x)$ as

$$\hat{p}_{\text{Histogram}}(x) = \frac{1}{N} \sum_{i=1}^m \left(\sum_{j=1}^N g_i(x_j) \right) g_i(x)$$

In each bin; \hat{P}_{ii} will for large N tend to its mean = $\int_{bin} p(x) dx$ with a standard deviation that this mean $\pm \sqrt{\frac{\int_{bin} p(x) dx}{N}}$. (See

argument on 333^{1/4} for a more accurate analysis based on correct binomial distribution).

This method has the following disadvantages

- (1) AS $N \rightarrow \infty$, $\hat{P}_{ii}(i)$ does not $\rightarrow p(x)$ but rather to $\sum_{l=1}^m g_l(x) \int_{within bin} p(x) dx$
This can only be cured by letting bin size $\rightarrow 0$ as $N \rightarrow \infty$.
- (2) AS bin size uniform some bins

will have lots of entries (where $p(x)$ big) and others a small number. One would like to automatically choose larger bin size where $p(x)$ small.

S2 Orthogonal Functions

A generalization of g_1 is to replace the indicator functions $g_i(x)$ by general orthogonal functions $y_i(x)$

$$\int y_i(x) y_j(x) dx = \delta_{ij}$$

$$\hat{f}_0(x) = \sum_{i=1}^m \left[\int y_i(x) p(x) dx \right] y_i(x)$$

$$= \frac{1}{N} \sum_{i=1}^m \left(\sum_j y_i(x_j) \right) y_i(x).$$

Clearly the histogram formulae are special cases of this where

$g_i(x)$ are just indicator functions.
 By the central limit $\hat{P}_0(x)$ has
 a nice limit as $N \rightarrow \infty$ ($\frac{1}{\sqrt{N}} \sum g_i(x_i)$)
 $\rightarrow \int g_i(x) p(x) dx$ with error $\propto \frac{1}{\sqrt{N}}$)
 although the coefficients of $g_i(x)$
 are no longer statistically independent.
 Again $\hat{P}_0(x)$ does not $\rightarrow p(x)$
 unless you take $M \rightarrow \infty$ as $N \rightarrow \infty$
 so that all terms are included
 in $g_i(x)$ expansion of $p(x)$.

g3 Empirical Cumulative Distribution

Define

$$\begin{aligned}\hat{F}(x) &= 0 & x < x_1 \\ &= i/N & x_i \leq x \leq x_{i+1} \\ &= 1 & x \geq x_n\end{aligned}$$

where x_i have been ordered in

increasing size. This is estimate from our data of the cumulative distribution $F_{\hat{F}}^{(x)} = \int_{-\infty}^x p(x') dx'$

as $\hat{F}(x) = \frac{1}{N} \sum_{j=1}^N \Theta(x - x_j)$, we see from central limit theorem that $\hat{F}(x)$ does indeed $\rightarrow F(x)$ as $N \rightarrow \infty$.

This empirical distribution will be very useful in goodness of fit test. Now we use it to find another estimate of $p(x)$

§ 4 The Rosenblatt estimator

We just define $\hat{P}_R(y)$ as the naïve derivative of $\hat{F}(y)$

$$\hat{P}_R(y) = \frac{1}{2h} [\hat{F}(y+h) - \hat{F}(y-h)] \quad (4.1)$$

Actually this is just number of points in bin $y-h \leq x \leq y+h$ divided by bin size. So this is really no different from histogram method except we have a bin of size $2h$ centered at each point x . Like the other estimators to be discussed $\hat{P}_R(y) \approx \hat{P}_R(y')$ ($y \neq y'$) are strongly correlated (they use many of the same x_j in their calculation) whereas neighbouring bins of a histogram where

uncorrelated.

One can find the mean and variance of $\hat{P}_R(y)$ as

$$\begin{aligned}\langle \hat{P}_R(y) \rangle &= \frac{1}{2h} [F(x+h) - F(x-h)] \quad (4.2) \\ &= p(y) + \frac{h^2}{6} p''(y) + \dots \quad (4.1a)\end{aligned}$$

$$\langle (\hat{P}_R(y) - p(y))^2 \rangle = \frac{p(y)}{2hN} + \frac{h^4}{36} |p''(y)|^2 \quad (4.3)$$

These ~~estimates~~^{results} are derived on writing

$$\hat{P}_R(y) = \frac{1}{2hN} \sum_{i=1}^N f_0\left(\frac{y-x_i}{h}\right)$$

where $f_0(x) = 1$ for $-1 \leq x \leq 1$
 $= 0$ otherwise.

① is just standard $\sigma \propto \sqrt{N}$ errors we found for in histogram case.

② comes from $p''(y)$ term in $\langle \hat{P}_R(y) \rangle$. Notice these $p''(y)$ terms would be absent if we had

calculated behavior of $\hat{P}_E(y)$ wrt
 bin integral $\frac{1}{2h} [F(x+h) - F(x-h)]$:
 we would have found similar
 terms for the histogram method if
 we had collected viewed it as
 an estimate of $p(y)$ and not the
 bin integral.

From (4.3) we want $h \rightarrow 0$
 while $N \rightarrow \infty$. we can minimize
 (4.3) for fixed $p(y), p''(y)$ by choosing

$$h = C N^{\alpha} \quad (4.4)$$

Putting this into ~~left~~^{right} hand side of
 (4.3) and ~~requiring~~ $\frac{\partial}{\partial C}$ and $\frac{\partial}{\partial \alpha} = 0$,
 we find $\alpha = -\frac{1}{5}$ (So that both
 terms on r.h.s go like $N^{-4/5}$) and

$$C = \left[\frac{9p(y)}{21p''(y)^2} \right]^{1/5} \quad (4.5)$$

G5 Parzen Estimates

These are just like the Rosenthal case except that ψ defined below (4.3) is replaced by any function whose integral is 2 from $-\infty$ to $+\infty$.

Possible choices are

$$\psi_0(x) = \exp(-|x|)$$

$$\text{or } = \frac{2}{\sqrt{2\pi}} \exp\left(-\frac{x^2}{2}\right)$$

These generalized estimates have very similar properties to Rosenthal's.

g6 k-th nearest neighbour

This method allows h to be automatically determined and also to vary with x . We could apply method more generally but here we will use it on Rosenblatt's method. We choose h so that interval $y-h \leq x \leq y+h$ has exactly k points (i.e. k out of set of N values). Here k is fixed ahead of time although it will depend on N . This choice leads to formula

$$\hat{p}_k(y) = k(N)/2N h(\text{say}, k) \quad (6.1)$$

where h is determined for each y as above.

Qualitatively this method is reasonable i.e. it gets large where p_{KL} is small and it overcomes one of the major problems with the histogram method.

We can show that

$$\langle \hat{P}_{KL}(y) \rangle = p(y) + \frac{1}{24} \left[\frac{k(N)}{N} \right]^2 \frac{p''(y)}{p^2(y)} \quad (6.2)$$

$$\langle (\hat{P}_{KL}(y) - p(y))^2 \rangle = \frac{p^2(y)}{k(N)} + \left[\frac{p''(y)}{24 p^2(y)} \left(\frac{k(N)}{N} \right)^2 \right]^2 \quad (6.3)$$

There are just ~~(6.2)~~ 6.3 with $2h$ replaced by $\frac{k(N)}{N p(y)}$ as one would expect. The qualitative origin of the two terms in (6.2, 3) is just as for the Rosedalek case.

We need $R(N) \rightarrow \infty$ but
 $\frac{R(N)}{N} \rightarrow 0$
and we can estimate best $R(N)$
just as for h to get

$$\left| \frac{R(N)}{N} \right|_{\text{best}} = \left\{ \frac{144 p^6(x)}{|p''(x)|^2 N} \right\}^{1/5}$$

and $\frac{R(N)}{N}$ behaves in same way
 $(N^{-1/5})$ as h did in Rosenblatt case.

As discussed by Freedman,
this method fails where you hit
a boundary when finding $h(\cdot, y, x)$.
He recommends switching to a
Simple Rosenblatt near the edges
of a plot. Freedman also compares
all the methods on a typical
problem and finds the

ranking

Nearest neighbour > Rosenblatt >

Parzen (gaussian \hat{f}_0) > Histogram

The nearest neighbour method has not caught on because of

- a) complexity of calculation. (you must order x_i for quick computing)
- b) correlation between estimates $\hat{P}_k(y)$ at neighbouring y .

and data is age). Also described in Bravais (p. 160 onwards) is the Fisher-Snedecor F distribution which can be used to test if two variances are equal.

G20 Likelihood Ratio Method

We wish to test if a particular theoretical parameter λ has a particular value $\lambda^{(0)}$. Now $H^{(1)}$ is all other values of λ . we find

$T = \text{Value of likelihood maximized over all values of } \lambda \div \text{Value of likelihood calculated for } \lambda = \lambda^{(0)}$.

Note that in this ratio ~~all~~ the dubious undefined parts of L

cancelled out.

The value of T is a generalization of χ^2 test ($-2 \ln T = \chi^2_{\text{minimized}} - \chi^2$ calculated for $\lambda = \lambda^{(0)}$) but and clearly $T \approx 1$ is a ~~success~~ implies $H^{(0)}$ is acceptable. we need to be able to calculate distribution in T when guess $\lambda^{(0)}$ is true value to be able to decide if a particular T is probable if $H^{(0)}$ true. This calculation is not in general very easy! (except in limit of large n when h becomes Gaussian in λ). Bressler does the calculation for a gaussian distribution and derives the Student's t distribution

from this more general point of view.

531 Methods based on Cumulative Distribution

These avoid "binning" necessary in χ^2 method but are only ~~useful~~ in 1 dimension. We start with the empirical cumulative distribution

$$\hat{F}(x) = \frac{1}{n} \sum_{j=1}^n \Theta(x_j - x)$$

We wish to test if the x_j are drawn from a particular distribution $p(x)$ with

$$F(x) = \int_{-\infty}^x p(y) dy.$$

The Smirnov-Cramer-Von Mises test forms

$$\omega^2 = \int_{-\infty}^{+\infty} [\hat{F}(x) - F(x)]^2 p(x) dx$$

Rather remarkably, the distribution of ω^2 is independent of $p(x)$ * and one can show that

$$\langle \omega^2 \rangle = \frac{1}{6n}$$

$$\sigma_{\omega^2}^2 = \frac{4n-3}{180n^3}$$

- * This is not intuitively obvious to me but can be shown because
 - i) all moments of $F(x)$ only depend on $F(x)$
 - ii) then we can let $z = F(x)$ in integral to remove F dependence
(note $dz = p(x)dx$).

One can also use

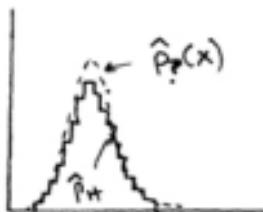
$$D = \max_{-\infty \leq x \leq \infty} |F(x) - \hat{F}(x)|$$

For large n , D has and this is discussed on p. 270 of ED 3RS. Again the distribution of D is independent of $F(x)$.

G12 Rootograms

In a normal histogram, where the i th bin contains N_i events, the error is proportional to $\sqrt{N_i}$. The variation of the size of error with population makes it hard to judge by eye the goodness of fit of some estimate $\hat{p}_*(x)$ to the histogram $\hat{p}_+(x)$.

[↑] maybe from maximum likelihood



One way around this is to plot $\hat{\alpha}_i = \sqrt{N_i}$ for

$$\delta \hat{\alpha}_i = \sqrt{N_i + S_{N_i}} - \sqrt{N_i}$$

$$\sim \frac{S_{N_i}}{\sqrt{N_i}} = 1.$$

So a plot of the α_i has an error independent of the number of counts and so the goodness of fit is easier to judge. Such a plot is called a nomogram. Friedman suggests plotting the residuals.

$$r(x) = \sqrt{\hat{p}_+(x)} - \sqrt{\hat{p}_-(x)}$$

as being even clearer. [$r(x)$ should have mean zero and an error independent of x ; i.e. if $\hat{p}_n(x) = cN(x)$ then standard deviation of $r(x)$ is $\pm \sqrt{c}$].

We found on one experiment where we were plotting electromagnetic shower pulse heights (so fitting them to predetermined forms) that we got a greater dynamic range $\sqrt{P_{\min}} \rightarrow \sqrt{P_{\max}}$ so was visually far better than $P_{\min} \rightarrow P_{\max}$ (e.g. in our case $P_{\min} \approx 10$ cts $P_{\max} \approx 1000$ cts and so a linear histogram had a dynamic range ~ 100 whereas the histogram only varied by a factor of 10).

g13 Multivariate ~~one~~ mapping methods

The previous techniques for non parametric representation do not easily generalize to ~~multiple~~ more than 1 dimension ^{when} if the observable x_i is a vector in d dimensions.

Typically one tackles vector valued observables by forming functions $y = f(\underline{x})$ and using univariate techniques on y . Most histogramming packages also allow one to examine correlations (using a "Scatter-plot") between two such f 's i.e. $y = f_1(\underline{x})$ and $z = f_2(\underline{x})$

$$\text{and } z = f_2(\underline{x})$$

if \underline{x} corresponds

to particle parameters
 $(f_a, f_b, f_c \text{ for 3 particles } a, b, c)$
 in a 3-particle final state



then perhaps

$$y \text{ is } (f_a + f_b)^2$$

$$z \text{ is } (f_b + f_c)^2$$

and Scatterplot is usually called the Dally plot.

One can choose the function $f(\mathbf{x})$ either by limited insight (e.g. one uses one's physics knowledge to suggest that mass² in the Dally plot are useful variables) or from statistical analysis of the actual distribution of the x_i . We will only discuss the latter methods here.

A) Principal Components

Here one chooses the general linear form $y = \underline{\alpha} \cdot \mathbf{x}$. For each α_j , we calculate the standard deviation of the univariate distribution of y .

$$\text{if } \underline{\underline{V}}_{kk} = \frac{1}{N} \sum_{l=1}^N \underline{x}_k^{(l)} \underline{x}_k^{(l)} - \left[\frac{1}{N} \sum_{l=1}^N \underline{x}_k^{(l)} \right] \left[\frac{1}{N} \sum_{l=1}^N \underline{x}_k^{(l)} \right]$$

is the usual co-variance matrix of the observations $\underline{x}^{(l)}$, then.

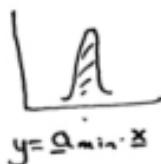
$$\sigma(y) = \sqrt{\underline{a}^T \underline{\underline{V}} \underline{a}} \quad (*)$$

The principal axis method gives the projection axis \underline{a} that maximizes $\sigma(y)$. This clearly corresponds to the eigenvector of largest eigenvalue for $\underline{\underline{V}}$. Continuing this process we deduce that the natural axes for examining our data are just the eigenvectors of the correlation matrix $\underline{\underline{V}}$.

e.g. y distributions vary from



$$y = a_{\max} \cdot z$$



$$y = a_{\min} \cdot z$$

largest standard deviation \rightarrow Smallest.

Now the above choices are not obviously best even when one restricts oneself to the linear form. (Nonlinear choices allow even freedom of course).

B: Projection Pursuit. (Friedman and Tukey,
IEEE transactions on Computers C-23, 881 (74))

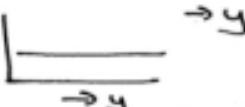
The idea here is that we wish to look for separated clusters and the above principal axis method emphasized separation but not clustering. So we try to maximize

$$I(\hat{a}) = \sigma(\hat{a}) d(\hat{a})$$

where $\sigma(\hat{a})$ is given by (*) on G37. $d(\hat{a})$ is to be chosen to be larger in situations like



than for



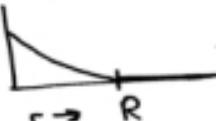
Friedman and Tukey choose

$$d(\hat{a}) = \sum_{ij=1}^n f(r_{ij})$$

where r_{ij} is distance from x_i to x_j projected on axis \hat{a}

$$r_{ij}^2 = |\hat{a} \cdot (x_i - x_j)|^2$$

and f is any smooth function of the form f



$$f = 0 \rightarrow R$$

and we search for structure of size R . Then d is big if a lot of points lie within distance R of each other

Maximizing $I(\hat{a})$ is a task can be performed by standard nonlinear maximization methods - it cannot be solved analytically as for choice of $d(\hat{a})$ as quantity to maximize.

The method can be generalized to more than one direction \hat{a} . Thus one replaces $I(\hat{a})$ by

$$I(\hat{a}, \hat{b}) = \sigma(\hat{a}) \sigma(\hat{b}) d(\hat{a}, \hat{b})$$

where in d one calculates

$$r_{ij}^2 = |\hat{a} \cdot (\mathbf{x}_i - \mathbf{x}_j)|^2 + |\hat{b} \cdot (\mathbf{x}_i - \mathbf{x}_j)|^2$$

and take \hat{a}, \hat{b} to be orthogonal directions.

c) Non parametric Representation of Multivariate Data (Friedman et al. SLAC-PUB-2116, 2466).

We can generalize the above ideas to least-squares estimate of the multivariate probability distribution (as discussed earlier in this course for the univariate case). We don't all the tools necessary for this yet (we need to discuss criteria for comparing

two samples from two possibly different multivariate distributions. However we can discuss the slightly simpler problem of "smoothing" i.e. ~~given~~ we suppose that a "response" $y^{(i)}$ is associated with each of our random variables $x^{(i)}$. We wish to find an estimate of the supposed functional relation:

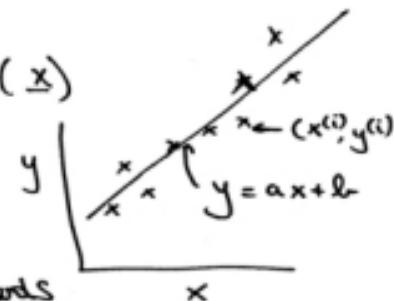
$$\hat{y} = \varphi(x)$$

In the 1D case,

there are various

well known ~~forms~~ methods

(least squares fit \approx regression analysis, Splines) which we will hopefully discuss. To generalize to the multidimensional



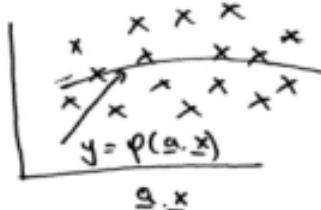
case, we suppose that

$$y = \varphi(\underline{a} \cdot \underline{x})$$

for each \underline{a} , we apply the univariate smoothing algorithm to the points $(\underline{a} \cdot \underline{x}^{(i)}, y^{(i)})$. In general this will not look very good

as y is multivalued

as a function of $\underline{a} \cdot \underline{x}$



So we choose \underline{a} to minimize this:

i.e. to minimize Scatter about the Smooth i.e. \underline{a} is vector that minimizes

$$\sum_{i=1}^n (y^{(i)} - \varphi(\underline{a} \cdot \underline{x}^{(i)}))^2 \quad (*)$$

Again this needs a nonlinear minimizer.

Calling this vector $\underline{a}^{(i)}$ we replace

$$y^{(i)} \rightarrow y^{(i)} - \varphi_{\min}(\underline{a}^{(i)} \cdot \underline{x}^{(i)})$$

and repeat procedure until the residual sum (*) is small.

S14 Minimal Spanning Trees: (mST)

See J. L. Bentley and J. H. Friedman

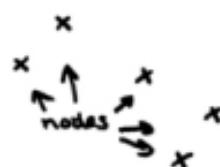
Ieee Trans on Comp C-27, 97 (1978)

J. Dofan, "A Cluster Algorithm for the
Study of Jets in High Energy Physics;
SLAC - 2623 (1980).

Given a bunch of points in a multidimensional Space on which a metric has been defined i.e. d_{ij} is distance between $x^{(i)}$ and $x^{(j)}$. Then a the mST is gotten by joining points ~~x⁽ⁱ⁾~~ (called nodes) by edges in such a way that

- (i) There is a connected path between any two nodes
- (ii) The sum of the lengths of the edges is the smallest possible

e.g. in

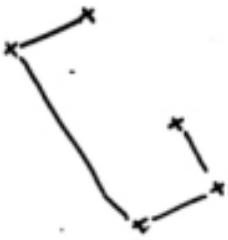


with distance as euclidean metric

Then
MST
is:



and



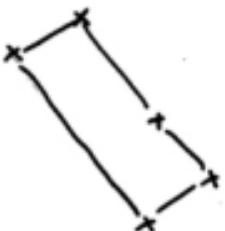
is not an MST
as the Σ lengths of
connected edges is
longer than optimal
solution.

and



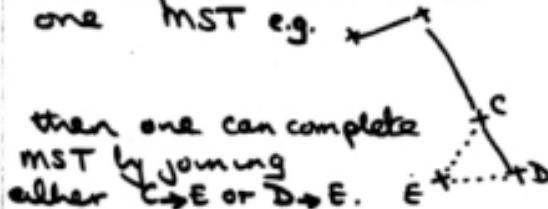
is not an MST as
there is no path to
A

and



is not an MST as
there is an
unnecessary edge.

It can be shown that an MST
always exists and is essentially
unique (i.e. only if some of the d_{ij}
are equal can there be more than
one MST e.g.

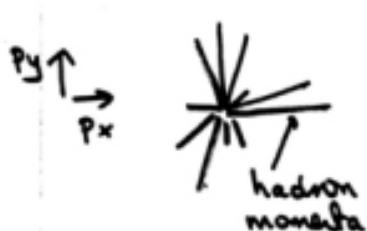
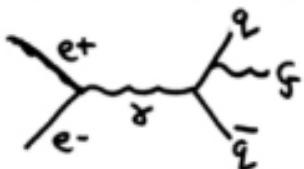


if $d_{CE} = d_{DE}$

then one can complete
MST by joining
either C-E or D-E. E

.....D

MST's have many uses that don't concern us. Obvious examples are connecting a telephone network with minimal amount of wire and of connecting electrical components in a chip. Zahn (Ieee Trans on Computers 20, 68(71)) discusses use of MSTs to find general patterns. Dorfan (See D43) uses the MST to find jets in e^+e^- collisions.



appears as a bunch of hadrons loosely clustered around original parton directions.
If we establish a metric
 $d_{ij} = |p_i - p_j|$ of particle i wrt j' (where i

is that one of i, j which is lower momentum and j' is higher momentum particle).

(Dorfan uses a different metric: ours is motivated by the limited p_T of particles w.r.t original jet axis). Then MST will naturally connect particles from the same jet and when one "jumps" from one jet to another one will use an edge which is longer than normal. Such long edges (called bridges) can be recognized and the collection of hadrons naturally decompose into jets. The same technique is used by Zahn to break a general bunch of points into isolated clusters (part of the "pattern recognition" problem).

We will use MSTs later in a slightly different as they provide a natural ordering of the points in a multidimensional space.

Note that Bentley and Friedman show that MSTs can be found in the typical $N \log N$ sorting time.

G105 Comparison of Multivariate Distributions

We now return to the problem discussed in G10.11 for the univariate case.

i.e. This problem has two variants.

(i) Given two sets of points $\{x^{(1,j)} : j=1..N_1\}$ $\{x^{(2,j)} : j=1..N_2\}$, are these two sets drawn from the same probability distribution?

(ii) Given one set of points $\{x^{(l)} : l=1..N\}$ - is it drawn from a given probability distribution $p(x)$? This second case can be reduced to the first if one uses a Monte Carlo method to generate a set of points according to $p(x)$. In this case $p(x)$ could be the probability distribution found from the maximum likelihood method.

Pr: k'th Nearest Neighbour Comparison of two Sets

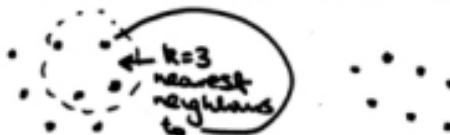
~~Steps for k-th Nearest Neighbour Comparison~~

Join our two sets into a sample of size $N_1 + N_2$ but remember whether from which of two sets each point comes from. Fix k .

Now look at each point and find the k points from joint sample that are nearest to the point. Record the number $n_1(k)$ of points from set 1 that are ~~nearest~~ among our k nearest neighbours.

Consider the N_1 values $S_1: \{n_1(k)\}$ of ~~Set 1 counts among k-th nearest~~ neighbors of set 1 points : and N_2 values $S_2: \{n_2(k)\}$ of ~~Set 2 counts~~ among k th nearest neighbors of set 2 points. Then if the two sets ~~are~~ $\{S^{(k)}\}$ are really the same, the two integer sets S_1, S_2 are also the same. The point is that we have swapped multivariate comparison to a univariate comparison. The latter can be performed by simple modifications of the methods described in Pg 21. (The latter is phrased for case of comparing sample with its supposed distribution but the modification for comparing two samples is straightforward). We can discuss how the integers in the sets S_i should be distributed.

If the vector sets $\{\underline{x}^{(1,i)}\}$ bear no resemblance to each other



Set 1:

Set 2:

then nearest neighbors to a set 1 point are all set 1 and to a set 2, all set 2.

$$\text{Then } S_1 \text{ is } \begin{pmatrix} k_1, k_1 & \dots & k_n \\ 0, 0 & \dots & 0 \end{pmatrix}$$

$$S_2 \text{ is } \begin{pmatrix} 0, 0 & \dots & 0 \\ k_1, k_1 & \dots & k_n \end{pmatrix}$$

~~and in fact they agree even though the original vector sets were dissimilar~~
and the integer sets are very dissimilar.

If the vector sets are in fact from the same distribution, then the distribution of $\gamma_2(k)$ is - near both

a by set 1 and set 2 point multinomial corresponding to k trials with probability $N_1/(N_1+N_2)$ of success. This is not quite right because the po set 1 points in different neighbourhoods overlap; one would have to do a Monte Carlo calculation to find \approx the distribution of $n_2(k)$ for a particular multivariate distribution $p(\underline{x})$. The binomial form is \approx usually good and one can compare $S_1: \{n_2(k) \text{ for } k \text{ Set 1 centers}\}$ directly with this distribution using standard univariate methods for comparing a sample and an hypothesized distribution (§1a).

B: k'th Nearest Neighbour Comparison between a Sample Set and a Supposed distribution $p(x)$ (problem 2 on p.547)

As we discussed

We can use the volume of each k'th nearest neighbour set to construct a nice univariate statistic for if we define

$V^{(k)}$ = volume of k'th nearest neighbour set (taken as a Sphere) about point x_i ^{P.T.O.}

$$\text{and } V^{(k)} = \int_{\text{k nn. volume about } x_i^{(k)}} \frac{1}{p(x')} \text{ Pactual}(x') dx'$$

$$= \frac{1}{N} \sum_{i=1}^N \frac{1}{p(x_i^{(k)})} \begin{matrix} \frac{1}{N} \text{ point in volume} \\ 0 \text{ point outside volume} \end{matrix}$$

($p(x)$ is hypothesised distribution:
Pactual is true distribution)

$$= \frac{1}{N} \sum_{\substack{\text{R+ points incl. center} \\ \text{inside nn volume}}} \frac{1}{p(x^{(k)})}$$

Sphere whose
+ More precise seems average of
volume is average of k and k+1th
nearest neighbour Spheres.

Now $\{v^{(1)}\}$ $\{v^{(2)}\}$ can be compared by our univariate methods again.

This method ~~only~~ only tests p near our sample points $x^{(1)}$; if p is very large in regions where $x^{(1)}$'s do not exist, then one will only indirectly see this through the normalization of p . This effect is not present in the alternate method of generating points according to $p(x)$ and apply the two sample comparison methods.

Note this nn methods ~~are~~^{is} sensitive to metric used to define volume. It can be shown that it is best to transform co-ordinates so that covariance matrix is the unit matrix.

C: Methods Based on MST

Consider again two samples of size N_1 and N_2 . Then the univariate Smirnov test for equality of their underlying distribution can be phrased as follows:

(i) Order the points by ascending value irrespective of the sample they come from.

(ii) For each i , $1 \leq i \leq N_1 + N_2$ Set

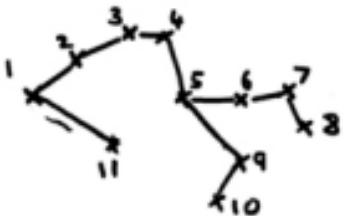
$$d_i = \frac{\left(\# \text{Set 2 points whose position in list is } \leq i \right) - \left(\# \text{Set 2 points whose position in list is } \leq i \right)}{N_2}$$

$$D = \max_i |d_i|.$$

We generalize this to the multivariate case by taking our vectors $\underline{x}^{(i)}$ and constructing the Minimal Spanning Tree (MST) of the pooled data. We can use this to order the data as follows;

(1) Find the "beginning (or end)" of the tree.
 This is ~~done~~ found by taking each node and finding the minimum length like breadth of a path to each other node b . Then the path of maximum length joins the beginning and end nodes of the tree.

(2) Order the points on tree by starting at beginning and traverse tree hierarchically - choosing nearest node whenever there is a choice



This shows a typical ranking.

(3) Calculate d_i and D just as on page 953.

An alternate univariate method is the Wald-Wolfowitz run test. we start as in 55 (i) but simply count the number of "runs" i.e. number of consecutive sequences of the same label.

If $N_1 = 2$ $N_2 = 1$ there are three orderings this labels \rightarrow

Set to which part belongs	$\begin{matrix} 2 & 2 & 1 \\ x & x & x \\ 2 & 1 & 2 \\ x & x & x \end{matrix}$	* runs = 2
	$\begin{matrix} 2 & 2 & 1 \\ x & x & x \\ 1 & 2 & 2 \\ x & x & x \end{matrix}$	3
	$\begin{matrix} 2 & 2 & 1 \\ x & x & x \\ 1 & 2 & 2 \\ x & x & x \end{matrix}$	2

values of ordered points \rightarrow

If the samples come from the same distribution, then the distribution of the * runs can be calculated by combinatorics - independent of the distributions. The mean number of runs is

$$\text{just part } \frac{1}{N} + \sum_{k=2}^N \frac{N_1}{N} \frac{N_2}{N-1} + \frac{N_2}{N} \frac{N_1}{N-1}$$

$k-1$ th type : k th type 2 $k-1$ k th type 2 k th type 1.

$$N = N_1 + N_2 .$$

by counting chance of a change in
Set # at k'th point of pooled data.

$$\langle \text{runs} \rangle = 1 + \frac{2N_1 N_2}{N_1 + N_2}$$

$$= 2^{1/2} \text{ in example on } 55.$$

We can clearly calculate standard deviation ~~standard~~ ($\sigma = \frac{2N_1 N_2 (2N_1 N_2 - N)}{N^2 (N-1)}$) of a distribution which becomes as usual Gaussian for large N_1, N_2 .

This method is even easier to generalize to the multivariate case. Take the MST as before but there is no need to rank points. Rather remove all edges than link ~~points~~ nodes (points $\mathbf{x}^{(i)}$) from different samples. Then the number of connected subtrees that remain is exactly analogous to the number of runs.

D: Estimate of a multivariate distribution

Given $\{\mathbf{x}^{(i)} \mid i=1..N\}$, we wish to estimate the underlying distribution $p(\mathbf{x})$. Here is a possible technique

- Transform co-ordinates to the principal axes so that the covariance matrix of the $\mathbf{x}^{(i)}$ is diagonal. We invent a first guess

$$t(\mathbf{x}) = t_1(x_1) t_2(x_2) \dots t_d(x_d) \quad (*)$$

where \mathbf{x} is a d dimensional vector. The product form $(*)$ is consistent with a diagonal covariance matrix but the latter only implies $(*)$ if each t_i is Gaussian. We can find $t_i(y)$ by projecting our sample onto each (principal) axis and using our (non parametric) univariate methods.

Now (*) may be good enough - use the methods described in B to find out if not, by

$$P(S) = t(\underline{x}) r(\underline{x})$$

and form for each point $\underline{x}^{(i)}$ in Sample

$$\underline{x}^R = \int_{\text{k'th nearest neighbour}}^{\frac{1}{t(\underline{x})} \text{ Pactual } (\underline{x}) d\underline{x}} V_R(\underline{x}^{(i)})$$

where Pactual is actual probability distribution. If Pactual(\underline{x}) = $p(\underline{x})$

$$R = \int_{V_R(\underline{x}^{(i)})} r(\underline{x}) d\underline{x} = V_k r(\underline{x}^{(i)})$$

where we have assumed that $r(\underline{x})$ is slowly varying over volume. We deduce that

$$r(\underline{x}^{(i)}) = \frac{1}{V_R(\underline{x}^{(i)})} \cdot \frac{1}{N} \sum_j \frac{1}{t(\underline{x}^{(j)})}$$

$k+1$ points $\underline{x}^{(j)}$ in
k'th nearest neighbour volume to $\underline{x}^{(i)}$

Now we have associated a "response" $r(\underline{x}^{(i)})$ to each point $\underline{x}^{(i)}$. We can use the techniques described in g13 to smooth (parameterize) $r(\underline{x}^{(i)})^+$ and so we have now obtained our final estimate $t(\underline{x}) r(\Delta)$ to the underlying distribution. We did have to assume that $r(\underline{x})$ was smooth over the nearest neighbor volume but even if it wasn't the data sample does not ^{have} information in it to tell structure over scales small compared to this volume.

* use $\log r$ to keep parameterization for r positive!

S17 Multivariate Decisions

Suppose we want to classify a vector random variable x into one of two (you replace two by any integer) possibilities. We discussed the univariate case in S8. This problem crops up often e.g. You wish to classify a particle as either π or electron (here x represents momentum measurement and pulse heights in various parts of a calorimeter). Alternatively you wish to classify a candidate track as either good or bad. Here x is a such things as χ^2 for segments of tracks making up candidate, # chamber hits in each segment, information on whether hits are already used in another good

track.

The general problem takes as data training vectors which are sets $\{x^{(1,i)} \mid i=1..N_1\}$, $\{x^{(2,i)} \mid i=1..N_2\}$ which are known to belong to be type I and II respectively. On the basis of these we are required to design the optimal algorithm to ~~separate~~^{classify} an arbitrary \underline{x} as type I or II. Note that if $\int p^{(1)}(\underline{s}) p^{(2)}(\underline{s}) d\underline{s}$ is nonzero then there is an intrinsic ambiguity in at least some regions of \underline{x} space.

A k'th Nearest Neighbour

Given \underline{x} , find the k nearest neighbours to \underline{x} from the pooled data $x^{(1,i)}, x^{(2,j)}$ of training vectors.

Classify x as type I/II if there are more type I/II points in this nearest neighbour set.

Comments:

- (i) If $N_1 \neq N_2$, one should weight events e.g. if $N_1 > N_2$ only weight type I events as of weight N_2/N_1 .
- (ii) One should (a) choose R and (b) transform co-ordinates to change metric for nearest neighbour distance calculation; by requiring these choices to minimize misclassification of vectors in training samples i.e. to minimize ^{that} # of type II points within nearest neighbour volume of type I points.

I believe this method if applied carefully should work well but it is computationally quite tricky (a) minimization alone is lengthy as nonlinear minimization (b) care needed near edges of data space in nearest neighbour

methods. (c) Well known but lengthy to find nearest neighbours.

B: Projection method

J. H. Friedman Isse T of C p 406 April 1971.

In the univariate case, suppose types I/II have cumulative distribution $F_{I,II}(x)$. Suppose we make a simple cut $x < x^*$ as I and $x > x^*$ as type II. ($F_1(x^*) > F_2(x^*)$). Then given a type I error chance of misclassification is $1 - F_1(x^*)$ and for a type II chance is $F_2(x^*)$. Thus average misclassification probability is $\frac{1 - F_1(x^*) + F_2(x^*)}{2}$ and is minimized by choosing x^* to maximize

$$D_{1,2}(x^*) = |F_1(x^*) - F_2(x^*)|.$$

We now describe Friedman's multivariate generalization of this.

Look at each component π_j of π and find $D_j = \max \{ F_I(x_j) - f_I(x_j) \}$. Here the cumulative distributions are estimated as in §3 from training sample. Choose j with maximum D_j . From previous page, a decision based on this j will have the smallest misclassification. The cut $x > x_j^*$ will now divide sample into 2 and we can repeat procedure and so build a binary tree with each node choosing a component with the smallest misclassification. This tree terminates if either (i) a given sample only contains training events of one

type. (ii) There are less than a certain preset number k of events in ~~the~~ ^{sample}.

Having built tree with training vectors one can use it to make a very fast decision on any given x as to its category. This is described for π/e case in a memo by J. Dofan (Mark II Internal note, 1979).