

Détection de Rayonnements à Très Basse Température

4^{ième} école d'Automne d'Aussois : Balaruc-les-Bains 14-20 novembre 1999

Statistique Élémentaire

J. Bouchez

DRTBT1999-17

Lorsque le résultat d'une observation ne peut pas être prédit avec certitude, celui-ci est décrit par une variable aléatoire X (uni ou multidimensionnelle) prenant ses valeurs dans un espace Ω .

Les sous-ensembles de Ω , appelés événements, sont munis d'une mesure P (probabilité)

$A \subset \Omega$: $P(A)$ est la probabilité que le résultat X d'une observation $\in A$

Cette mesure P jouit des propriétés suivantes :

$$P(A) \in [0, 1] \quad \forall A$$

$$P(\emptyset) = 0 \quad P(\Omega) = 1$$

$$P(A \cup B) = P(A) + P(B) - P(A \cap B)$$

Si \bar{A} est le complément de A ($A \cup \bar{A} = \Omega$, $A \cap \bar{A} = \emptyset$)

$$P(\bar{A}) = 1 - P(A)$$

Théorème Bayes

On définit la probabilité conditionnelle

$P(A|B)$ c.a.d la probabilité que

↑
si

$x \in A$ sachant que $x \in B$.

$$P(A|B) = \frac{P(A \cap B)}{P(B)}$$

A et B indépendants $\Leftrightarrow P(A|B) = P(A)$

Alors $P(A \cap B) = P(A) \cdot P(B)$

Lois de probabilité pour X

On peut toujours s'arranger pour que X prenne ses valeurs dans \mathbb{R}^k (\mathbb{R} ensemble des réels)

[Exemple : Pile ou face $P_{\text{pile}} = x=1$ Face $= x=0$]

X unidimensionnelle = prend ses valeurs dans \mathbb{R}

Cas discret X prend des valeurs discrètes x_1, x_2, \dots

$$P(X = x_i) = p_i \quad \sum p_i = 1$$

↓
on note $P(x_i)$

Cas continu

On définit $F(x) = P(X < x)$

$F(-\infty) = 0$ $F(+\infty) = 1$ $F(x)$ monotone croissante

$$P(X \in [x, x+dx]) = F(x+dx) - F(x)$$

$$= F'(x) dx = f(x) dx$$

$f(x)$: densité de probabilité

$$f(x) \geq 0 \quad \forall x \quad \int_{-\infty}^{+\infty} f(x) dx = 1$$

Cas multidimensionnel :

$$f(x, y, z, \dots) \geq 0 \quad \forall x, y, z, \dots$$

$$\int f dx dy dz \dots = 1$$

Loi de probabilité réduite

$$f_x(x) = \int f(x, y, z, \dots) dy dz \dots$$

loi de probabilité conditionnelle :

$$f_c(x | y = y_0) = \frac{f(x, y_0)}{\int f(x, y_0) dx}$$

à la p. 100

indépendantes si $f_c(x | y_0) = f_x(x) \forall y_0$

Indépendance \Leftrightarrow Factorisation

$$f(x, y) = f_x(x) f_y(y) \text{ si } x, y \text{ indépendantes}$$

Changement de variable

$y = H(x)$ y variable aléatoire liée fonctionnellement à x

x a pour densité de probabilité $f(x)$
 y " " " " " " " " $g(y)$

si la correspondance $x \leftrightarrow y$ est biunivoque

$$\text{alors } g(y) dy = f(x) dx$$

$$\text{ou } g(y) = \frac{f(x)}{|H'(x)|}$$

Cas multidimensionnel (correspondance biunivoque)

$$f(x, y) dx dy = g(x, y) dx dy$$

$$g(x, y) = \frac{f(x, y)}{\left| \begin{array}{cc} \frac{\partial x}{\partial x} & \frac{\partial x}{\partial y} \\ \frac{\partial y}{\partial x} & \frac{\partial y}{\partial y} \end{array} \right|} \quad \leftarrow \text{Jacobian}$$

Caractéristiques des lois de probⁿ dite

1) Valeur moyenne (ou espérance) notée \bar{x} , $\langle x \rangle$, $E(x)$

Cas discret $\langle x \rangle = \sum p_i x_i$

Cas continu $\langle x \rangle = \int x f(x) dx$

2) Variance (notée σ^2) σ est appelé sigma, écrit quadratique, incertitude, erreur, résolution...

$$\sigma^2 = \langle x^2 \rangle - \langle x \rangle^2 = \langle (x - \bar{x})^2 \rangle$$

3) Covariance

2 variables aléatoires x, y $f(x, y)$ densité de proba.

$$C_{xy} = \langle (x - \bar{x})(y - \bar{y}) \rangle = \int f(x, y) (x - \bar{x})(y - \bar{y}) dx dy$$

Coeff de corrélation $r_{xy} = \frac{C_{xy}}{\sigma_x \sigma_y}$

$$-1 \leq r \leq +1 \quad (\text{Inégalité de Schwarz})$$

$$x \text{ et } y \text{ indépendants} \Rightarrow C_{xy} = 0$$

Réciproque fautive!

4) Matrice de variance-covariance

pour x_1, x_2, \dots, x_k

$$V = \begin{pmatrix} \sigma_1^2 & C_{12} & \dots & C_{1k} \\ C_{12} & \sigma_2^2 & & \\ \vdots & & \ddots & \\ C_{1k} & & & \sigma_k^2 \end{pmatrix}$$

Symétrique

Semi définie > 0

$$X^T V X \geq 0 \quad \forall X$$

Changement de variables:

linéaire $Y_i = \sum a_{ij} X_j \quad Y = M X$

$$V_Y = M V_X M^T$$

quelconque $Y_i = f(x_1, \dots, x_k) \quad M_{ij} = \frac{\partial f}{\partial x_j}$

$$V_Y \sim M V_X M^T$$

Attention, approximation!

(5)

LOIS USUELLES

(+)

1) Exponentielle (des intégration d'une particule)

$$f(t) = \frac{1}{\tau} e^{-t/\tau} \quad t \geq 0$$

$$\langle t \rangle = \tau \quad \sigma_t^2 = \tau^2$$

2) Gaussienne (ou loi normale)

$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} = G(\mu, \sigma^2)$$

$$\langle x \rangle = \mu \quad \sigma_x^2 = \sigma^2$$

3) Binômiale

N observations (N fixé) n résultats $\in A$ $P(A)=p$

$$P(m) = C_N^m p^m (1-p)^{N-m} \quad C_N^m = \frac{N!}{m!(N-m)!}$$

$$\langle m \rangle = Np \quad \sigma_m^2 = Np(1-p)$$

4) Poisson

$$P(m) = e^{-a} \frac{a^m}{m!} = G_a(m)$$

$$\langle m \rangle = a \quad \sigma_m^2 = a$$

Théorèmes de convergence

(8)

1) Loi binômiale N observations, $n \in A$, $P(A)=p$.

$$f = \frac{n}{N} \quad \langle f \rangle = p \quad \sigma_f^2 = \frac{p(1-p)}{N}$$

Quand $N \rightarrow \infty$, la variable aléatoire

f devient "certaine" puisque $\sigma_f^2 \rightarrow 0$ comme $\frac{1}{N}$

2) Théorème Central limit

N mesures indépendantes de la variable aléatoire X de densité $f(x)$, de moyenne μ et de variance σ^2

$$S = \sum_1^N x_i \quad M = \frac{S}{N} \quad R = \frac{S - N\mu}{\sigma\sqrt{N}}$$

Quand $N \rightarrow \infty$

• $M \rightarrow \mu$ "certainement" $\sigma_M^2 \rightarrow 0$ comme $\frac{1}{N}$

• R a une densité de proba $G_N(R)$

$$G_N(R) \xrightarrow{N \rightarrow \infty} G(0, 1) \text{ Gaussienne}$$

• S a une densité de proba $h_N(S)$

$$h_N(S) \xrightarrow{N \rightarrow \infty} G(N\mu, N\sigma^2)$$

COMPLEMENTS COURS 1

LOI MULTIGAUSSIENNE (MULTINORMALE)

k variables X_i de variance-covariance V

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_k \end{pmatrix} \quad \bar{X} = \begin{pmatrix} \bar{x}_1 \\ \vdots \\ \bar{x}_k \end{pmatrix}$$

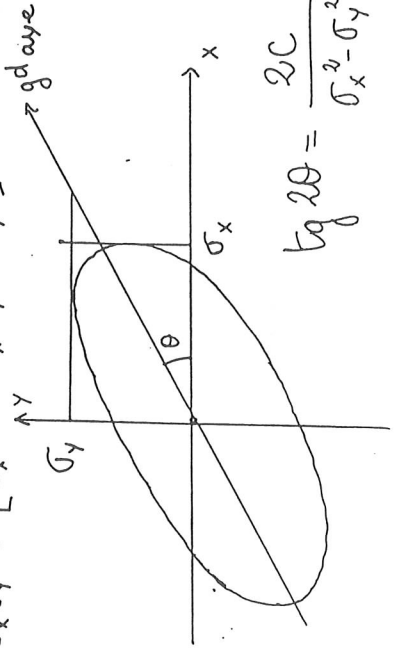
$$f(x_1, \dots, x_k) = 2\pi^{-\frac{k}{2}} (\det V)^{-\frac{1}{2}} \exp \left\{ -\frac{1}{2} (X - \bar{X})^T V^{-1} (X - \bar{X}) \right\}$$

les iso contours sont des k -ellipses

$$k=2 \quad V = \begin{pmatrix} \sigma_x^2 & c \\ c & \sigma_y^2 \end{pmatrix} \quad \bar{X} = \bar{Y} = 0 \quad \text{pour simplifier}$$

$$f(x, y) = \frac{1}{2\pi \sqrt{\sigma_x^2 \sigma_y^2 - c^2}} \exp \left\{ -\frac{1}{2} \frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 \sigma_y^2 - c^2} \left[\frac{x^2}{\sigma_x^2} - 2c \frac{xy}{\sigma_x^2 \sigma_y^2} + \frac{y^2}{\sigma_y^2} \right] \right\}$$

$$\frac{\sigma_x^2 \sigma_y^2}{\sigma_x^2 \sigma_y^2 - c^2} \left[\frac{x^2}{\sigma_x^2} - 2c \frac{xy}{\sigma_x^2 \sigma_y^2} + \frac{y^2}{\sigma_y^2} \right] = 1 \quad : \text{ ellipse}$$



$$\text{tg } 2\theta = \frac{2c}{\sigma_x^2 - \sigma_y^2}$$

SOMME DE VARIABLES ALEATOIRES INDEPENDANTES

X de loi $g(x)$ } + X et Y indépendantes
 Y de loi $h(y)$ }

$S = x + y$ S a pour densité de probabilité

$$f(s) = \int g(s-y) h(y) dy = g * h$$

Cas particuliers intéressants :

Loi de Poisson $P_a * P_b = P_{a+b}$

Loi Binomiale $B_{N_1, p} * B_{N_2, p} = B_{N_1 + N_2, p}$

Gaussienne $G(N_1, \sigma_1^2) * G(N_2, \sigma_2^2) = G(N_1 + N_2, \sigma_1^2 + \sigma_2^2)$

Théorème "Central-limit"

$$f(x) \quad \bar{x} = \mu, \quad V_x = \sigma^2$$

$$f^{*N} \xrightarrow{N \rightarrow \infty} G(N\mu, N\sigma^2)$$

①

ESTIMATEURS

Définitions X de loi $f(x; \theta_0)$ θ_0 paramètre inconnu
 N observations x_i

$t = h(x_1, \dots, x_n)$ est un estimateur de θ_0 non biaisé si $\langle t \rangle = \theta_0$

t est un estimateur convergent de θ_0 si $\langle t_N \rangle = \theta_0 + b_N$ $b_N \xrightarrow{N \rightarrow \infty} 0$ comme $\frac{1}{N}$

$\sigma_{t_N}^2 \xrightarrow{N \rightarrow \infty} 0$ comme $\frac{1}{N}$

- t_N sera dit optimal si
 - il est non biaisé
 - il est de variance minimum

On peut montrer qu'il existe une borne inférieure à la variance d'un estimateur

- t_N sera dit efficace s'il est optimal et si sa variance vaut cette borne inférieure
- ⚠ Les estimateurs efficaces n'existent pas toujours

ESTIMATEURS USUELS

I. Méthode des moments

X suit $f(x; \theta_0)$ - On considère $Y = a(X)$

$E(Y) = \int a(x) f(x; \theta_0) = g(\theta_0)$

N observations x_i :

$Z = \frac{1}{N} \sum_{i=1}^N a(x_i)$

- Z est un estimateur non biaisé et convergent de $g(\theta_0)$, pas nécessairement optimal
- $g^{-1}(z)$ est un estimateur biaisé et convergent de θ_0 , pas nécessairement optimal

Estimateur utilisé du fait de sa simplicité

$\hat{\theta} = g^{-1}(z) \quad V_{\hat{\theta}} \rightarrow \frac{1}{N} \left[\frac{\partial g}{\partial \theta}(\theta_0) \right]^{-2} V_a(\theta_0)$

en pratique, estimés en $\hat{\theta}$

Critère de Darmois

- Si 1) $f(x; \theta_0) = \exp \{ \alpha(x) a(\theta) + \beta(x) + c(\theta) \}$
 2) le domaine de variation de x ne dépend pas de θ

Le maximum de vraisemblance à échantillon fini donne une estimation efficace de

$$\kappa(\theta_0) = \frac{c'(\theta_0)}{a'(\theta_0)} \quad c'(\theta_0) = \frac{dc(\theta)}{d\theta} \Big|_{\theta=\theta_0}$$

[Se généralise au cas où θ est multidimensionnel]

III Maximum de vraisemblance

X suit $f(x; \theta_0)$ θ_0 paramètre inconnu.
 On estime θ_0 par $\hat{\theta}$ défini par

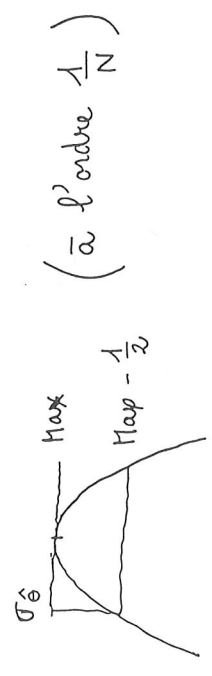
$$\log L = \sum \log f(x_i; \hat{\theta}) \quad \text{maximum}$$

- Estimateur asymptotiquement efficace
- A échantillon fini, biaisé et non nécessairement optimal

$$V(\hat{\theta}) \xrightarrow{N \rightarrow \infty} \frac{1}{N} \left[-E \left(\frac{\partial^2 \log L}{\partial \theta^2} \right) \Big|_{\theta=\theta_0} \right]^{-1}$$

en pratique

$$V(\hat{\theta}) \sim \left[\frac{\partial^2 \log L}{\partial \theta^2} \Big|_{\theta=\hat{\theta}} \right]^{-1}$$



Moindres carrés

k observations x_i de moyennes $\mu_i(\theta_0)$ et de variances $\sigma_i^2(\theta_0)$ θ_0 paramètre inconnu

$$\chi^2 = \sum_1^k \frac{(x_i - \mu_i(\theta))^2}{\sigma_i^2(\theta)} \text{ minimum en } \theta = \hat{\theta}$$

$\hat{\theta}$ est un estimateur convergent de θ_0

A échantillon fini, il est biaisé et non nécr optimal

$$V(\hat{\theta}) \xrightarrow{k \rightarrow \infty} 2 \left[\frac{\partial^2 \chi^2}{\partial \theta^2} \right]^{-1}_{\theta = \theta_0} \sim 2 \left[\frac{\partial^2 \chi^2}{\partial \theta^2} \right]^{-1}_{\theta = \hat{\theta}}$$

$\chi^2 \rightarrow \chi^2_{\min} + 1$ pour avoir une estimation des dérivées secondes

Rem Si les x_i suivent des lois gaussiennes, $\text{Min}(\chi^2)$ et Max log LL sont des estimateurs équivalents

Si 1°) $\mu_i = A_i \theta + B_i$

2°) $V(\theta)$ indépendante de θ

alors
$$\chi^2 = \sum (x_i - \mu_i(\theta)) V_{ij}^{-1} (x_j - \mu_j(\theta))$$

donne un estimateur $\hat{\theta}$ de θ qui est

- 1) soluble analytiquement
- 2) non biaisé
- 3) optimal



L'ajustement d'histogrammes au le contenu de chaque bin n_i suit une loi de poisson de moyenne $\bar{n}_i(\theta_0)$ ne dépend pas aux hypothèses du modèle linéaire, car les variances $\bar{n}_i(\theta_0)$ dépendent du (ou des) paramètres

~~4) χ^2 dépend~~

χ^2 linéaire, suite.

Si de plus les variables x_i suivent une

loi multigaussienne =

- 1) $\chi^2(\theta_0)$ suit une loi de χ^2 à k degrés de lib.
- 2) $\chi^2(\hat{\theta})$ suit une loi de χ^2 à $(k-r)$ degrés de liberté (r dimension de θ)

La loi de χ^2 à N degrés de liberté est la loi
suivre par $\sum_1^N x_i^2$ lorsque chaque x_i
suit une loi $G(0,1)$

Loi tabulée

$$\langle \chi_N^2 \rangle = N$$

$$V_{\chi_N^2} = 2N$$

$$f(\chi_N^2) = \frac{1}{2} \left(\frac{\chi_N^2}{2} \right)^{\frac{N-1}{2}} e^{-\frac{\chi_N^2}{2}} \frac{1}{\Gamma(\frac{N}{2})}$$

- 3) $\hat{\theta}$ a une distribution (multi)gaussienne

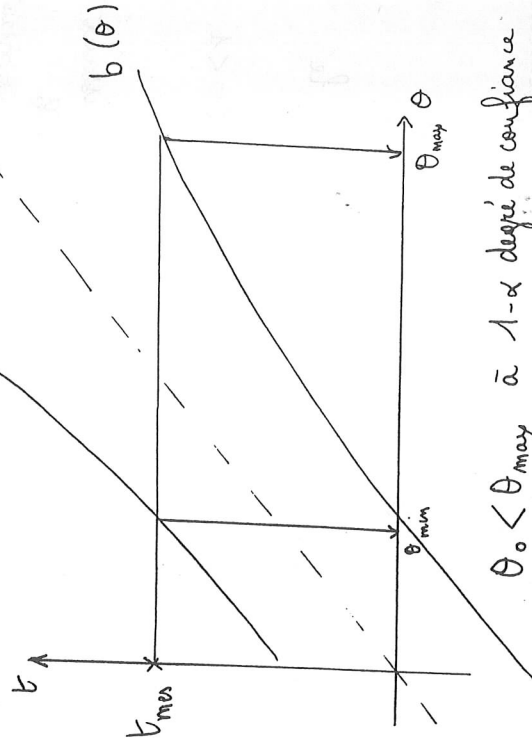
INTERVALLES DE CONFIANCE

t estimateur de θ_0 , de densité de probabilité connue

$$f(t|\theta)$$

On définit $b(\theta) \rightarrow P(t < b|\theta) = \alpha$
 $B(\theta) \rightarrow P(t > B|\theta) = \alpha$.

Par exemple $\alpha = 5\%$



$\theta_0 < \theta_{max}$ à $1-\alpha$ degré de confiance

$\theta_0 > \theta_{min}$ " "

$\theta_{min} < \theta_0 < \theta_{max}$ à $1-2\alpha$ degré de confiance

Ces 3 assertions concernent le comportement des variables aléatoires θ_{min} et θ_{max}

② Méthode de Casimir-Feldmann

RÉSUMÉ COURS n°3

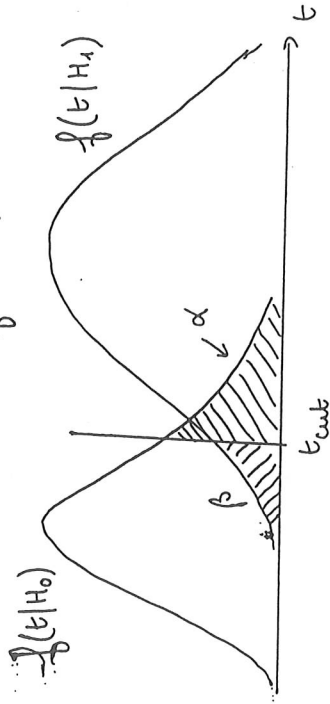
A) Tests d'hypothèses simples

2 hypothèses H_0 et H_1 complètement spécifiées

t fonction des observations x_i suit la

loi de probabilité $f(t|H_0)$ quand H_0 est vrai

$f(t|H_1)$ " H_1 "



On rejette H_0 si $t > t_{crit}$:

Perte α (H_0 rejeté quand H_0 vrai)

Contamination β (H_0 accepté quand H_1 vrai)

$1 - \beta$ est la puissance du test

Théorème

Le test le plus puissant (Neyman-Pearson) est

le rapport des vraisemblances

$$t = \frac{f(x|\theta_1)}{f(x|\theta_0)}$$

$X = \{x_1, \dots, x_n\}$

B) Test d'hypothèses paramétriques

observations $\{x_1, \dots, x_n\} = X$

$f(x_i|\theta)$ $\theta \in \Omega$ espace des paramètres

$$H_0: \theta \in \mathcal{V} \subset \Omega$$

$$H_1: \theta \in \mathcal{V}' \subset \Omega$$

$$\mathcal{V} \cap \mathcal{V}' = \emptyset$$

On utilise pour accepter ou rejeter H_0 le

$$t = \frac{\text{Max}_{\theta \in \mathcal{V}} f(x|\theta)}{\text{Max}_{\theta \in \mathcal{V}'} f(x|\theta)} \quad 0 \leq t \leq 1$$

Δ Le test n est plus puissant le plus puissant

Si H_0 consiste à fixer la valeur de ~~estimation~~ ^{paramètre} θ et si $\mathcal{V} \cup \mathcal{V}' = \Omega$, $-\ln t$ suit asymptotiquement une loi de χ^2 à k degrés de liberté.

Propriété

③ Tests de "qualité d'ajustement"

H_0 spécifiée

H_1 non spécifiée

On veut "savoir" si les données sont compatibles avec H_0 -

La vraisemblance ne donne lieu à aucun test utile

1 - Données x_i binées dans un histogramme :

On fait un test de χ^2 (Pearson)

$$\chi^2 = \sum_{i=1}^k \frac{(m_i - N p_i)^2}{N p_i} \quad \text{où } p_i \text{ sont les probas. du bin } i \text{ pour l'hyp. } H_0$$

χ^2 suit une loi de χ^2 à $k-1$ degrés de liberté

(~ vrai dès que $m_i > 5 \forall i$)

Choix optimal du binning $p_i = \frac{1}{k} \forall i$

$P(\chi^2) = \int_{k-1}^{\infty} \chi^2 dx$ est le degré de confiance de l'ajustement

si $P(\chi^2) < \alpha$, on rejette H_0 avec une perte α

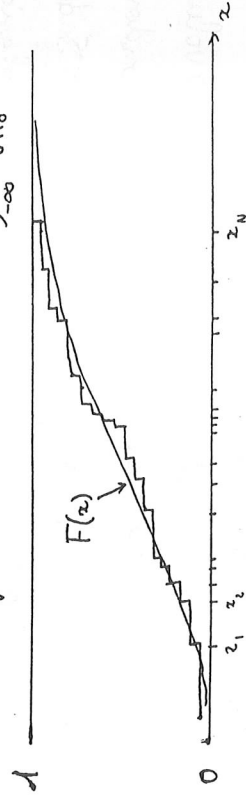
Si H_0 dépend de r paramètres inconnus, on minimise le χ^2 sur ces r paramètres et on teste H_0 par χ^2_{\min} qui suit asymptotiquement une loi de χ^2 à $k-r-1$ degrés de liberté

② Données non binées

N Observations x_i tirées selon $f_{H_0}(x)$, x unidimensionnel

On compare $S_N(x)$ fonction escalier de marche $\frac{1}{N}$ en chaque x_i

à la fonction cumulative $F(x) = \int_{-\infty}^x f_{H_0}(u) du$



Test de Smirnov :

$$W_N^2 = \int_{-\infty}^{+\infty} (S_N - F)^2 f_{H_0}(x) dx$$

W_N^2 suit une loi de proba indépendante de f tabulable

$$E(W_N^2) = \frac{1}{6N} \quad \sigma^2(W_N^2) = \frac{4N-3}{180N^3}$$

Vrai seulement si H_0 complètement spécifiée (pas de paramètres inconnus)

$$K_N = \max_x |F(x) - S_N(x)|$$

K_N suit asymptotiquement ($N > 180$) une loi universelle
[à condition que H_0 soit complètement spécifiée]

- ③ Degré de confiance conjoint de 2 tests
indépendants de H_0 , de degrés de confiance
 α_1 et α_2

$$\alpha = \alpha_1 \alpha_2 (1 - \ln \alpha_1 \alpha_2)$$

Exemple $\alpha_1 = 10\%$
 $\alpha_2 = 15\%$ } $\Rightarrow \alpha \approx 7.8\%$

Bibliographie

En français

- ① C. Fourgeaud, A. Fuchs
"Statistique" DUNOD 67

(Collection universitaire de mathématiques)

- ② Ecole d'été de Physique des Particules de Cif, 1977.

"Statistique appliquée à la Physique"

(L. Behr, D. Drijard, B. Sadoulet, B. Scher)

Edition IN2P3

En anglais

- ① M. Kendall and A. Stuart

"The advanced theory of statistics" 3 volumes

Charles Griffin & company limited, London

- ② W.T. Eadie, D. Drijard, F.E. James, M. Roos, B. Sadoulet

"Statistical methods in experimental physics"

North Holland Publishing Company, Amsterdam-Louis

- ③ Byron E. ROE

"Probability & Statistics in Experimental Physics"

Springer-Verlag

- ④ "Statistics for nuclear and particle physicists"

L. Lyons (1985)

PROBABILITY, STATISTICS, AND MONTE CARLO

1. PROBABILITY

1.1 General

If x is the outcome of an observation, we define the probability of x as the relative frequency with which x occurs out of a (possibly hypothetical) large set of similar observations. If x may take any value from a continuous range, we write f(x;θ) dx as the probability of observing x between x and x + dx. The function f(x;θ) is the probability density function (p.d.f.) for the random variable x, which may depend upon a parameter θ. If x can take on only one of a set of discrete values (e.g., the non-negative integers), then f(x;θ) is itself a probability, but we still refer to it as a p.d.f. The p.d.f. is always normalized to unit area (unit sum, if discrete). Both x and θ may have multiple components and are then usually written as column vectors. If θ is unknown and we wish to estimate its value from a given set of data x, we may use statistics (Section 2).

The cumulative distribution function F(x) expresses the probability that x ≤ a.

F(x) = ∫_{-∞}^x f(x) dx (1.1)

Here and in what follows, if x is discrete-valued, the integral is replaced by a sum. The endpoint a is expressly included in the integral or sum. Then 0 ≤ F(x) ≤ 1, F(x) is nondecreasing, and Prob(a ≤ x ≤ b) = F(b) - F(a). If x is discrete, F(x) is flat except at allowed values of x, where it has a discontinuous jump equal to f(x).

Any function of random variables is itself a random variable, with (in general) a different p.d.f. The expectation value of any function u(x) is

E[u(x)] = ∫_{-∞}^∞ u(x) f(x) dx (1.2)

The expectation value is said to exist only if it is finite. For x and y any two random variables, E[(x + y)] = E[x] + E[y]. For c and k constants, E[(cx + k)] = cE[x] + k.

The nth moment of a distribution is given by

α_n = E[x^n] (1.3a)

and the nth moment about the mean by

m_n = E[(x - α_1)^n] (1.3b)

The most commonly used are the mean and variance:

μ ≡ α_1 (1.4a)

σ^2 ≡ Var(x) ≡ m_2 = α_2 - μ^2 (1.4b)

The mean is the location of the "center of mass" of the distribution of x and the variance is a measure of the square of its width. Note that Var(cx + k) = c^2 Var(x).

Any odd moment about the mean is a measure of skewness; the simplest of these is the dimensionless coefficient of skewness

γ_1 = m_3/σ^3 (1.4c)

In addition to the mean, another useful indicator of the x location near which most of the probability is likely to concentrate is the median x_med. This is that value of x such that F(x_med) = 1/2, i.e., exactly half of the probability lies above and half lies below x_med.

For a given sample of events, x_med is that observed x such that half the events have larger x and half have smaller x (as closely as possible, not counting any that have the same x as the median). If this lies between two observed x values, the sample median is set by convention to be halfway between them. If the p.d.f. for x has the form f(x - μ) and μ is both mean and median, then for a large number of events N the variance of the median approaches 1/(4Nf'(μ)), provided f'(μ) > 0.

Let x and y be two random variables with joint p.d.f. f(x, y). The marginal p.d.f. of, for example, x, expressing the p.d.f. for x with y unobserved, is

f(x) = ∫_{-∞}^∞ f(x, y) dy (1.5)

and similarly for f_2(y). If y is fixed, the conditional p.d.f. for x given the fixed y is given by

f_1(x|y) = f(x, y)/f_2(y) (1.6)

The x mean is

μ_x = ∫_{-∞}^∞ ∫_{-∞}^∞ x f(x, y) dy dx = ∫_{-∞}^∞ x f_1(x) dx (1.7)

and similarly for y. The correlation between x and y is a measure of the dependence of one on the other:

ρ_{xy} = E[(x - μ_x)(y - μ_y)]/σ_x σ_y ≡ Cov(x, y)/σ_x σ_y (1.8)

where σ_x, σ_y are defined in analogy with Eq. (1.4b); it can be shown that -1 ≤ ρ_{xy} ≤ 1. The symbol "Cov" represents the covariance of x and y, a 2-variable analogue to the variance, Eq. (1.4b). Two random variables are independent if and only if

f(x, y) = f_1(x) f_2(y) (1.9)

If x and y are independent then ρ_{xy} = 0; the converse is not necessarily true except for Gaussian-distributed x and y. If x and y are independent, E[u(x)v(y)] = E[u(x)]E[v(y)] and Var(x + y) = Var(x) + Var(y); otherwise, Var(x + y) = Var(x) + Var(y) + 2Cov(x, y) and E[xy] does not factor.

In a change of continuous random variables from (x, y) to (x_1, x_2, ..., x_n), with p.d.f. f(x_1, ..., x_n), to (y_1, y_2, ..., y_n), a one-to-one function of the x's, the p.d.f. g(y_1, ..., y_n) is found by substitution for (x_1, ..., x_n) in f followed by multiplication by the absolute value of the Jacobian of the transformation:

g(y) = f[u_1(y), ..., u_n(y)]|J| (1.10)

The functions u_i express the reverse transformation x_i = u_i(y) for i = 1, ..., n, and |J| is the absolute value of the determinant of the square matrix J_{ij} = ∂x_i/∂y_j. Such transformations must always preserve the number of random variables, n. To transform to fewer variables, first perform (1.10) and then use Eq. (1.5) to eliminate unwanted variables. If the transformation from x to y is not one-to-one, the situation is more complex and a unique solution may not exist. To change variables for discrete random variables simply substitute; no Jacobian is necessary because in that case f is a probability rather than a probability density. If f depends upon a parameter set θ, we can change to a different parameter set φ = φ(θ) by simple substitution; no Jacobian is used.

1.2 Characteristic functions [1]

The characteristic function φ(u) associated with the p.d.f. f(x) is essentially its Fourier transform, or the expectation value of exp(iux),

φ(u) = E[e^{iux}] = ∫_{-∞}^∞ e^{iux} f(x) dx (1.11)

It is sufficiently useful to deserve special attention, and several of its properties follow.

We note from Eqs. (1.3a) and (1.11) that the nth moment of the distribution f(x) is given by

i^{-n} d^n φ / du^n |_{u=0} = ∫_{-∞}^∞ x^n f(x) dx = α_n (1.12)

As a result, it is often easy to calculate all the moments of a distribution defined by φ(u) even when the inversion is not available. If f_1(x) and f_2(x) have characteristic functions φ_1(u) and φ_2(u), then the characteristic function of the weighted sum ax + by is φ_1(ax)φ_2(by).

Let the (partial) characteristic function corresponding to the conditional p.d.f. f_1(x|z) be φ_1(u|z), and the p.d.f. of z be f_2(z). The characteristic function after integration over the conditional value is

φ(u) = ∫ φ_1(u|z) f_2(z) dz (1.13)

Suppose we can write φ_2 in the form

φ_2(u|z) = A(u)g(z) (1.14)

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

1.3.4 Normal or Gaussian distribution (continuous)

The Gaussian distribution is

f(x; μ, σ^2) = 1 / (σ√2π) e^{-(x-μ)^2/2σ^2}, -∞ < x < ∞ (1.24)

E[x] = μ; Var(x) = σ^2 (1.25)

The characteristic function of a Gaussian p.d.f. with mean m and variance σ^2 is

φ(u) = e^{iμu - 1/2σ^2 u^2} (1.26)

so the Gaussian is that unique distribution for which all semi-invariants beyond the second vanish.

For x and y independent and normally distributed, z = x + y obeys

f(z; μ_x + μ_y, σ_x^2 + σ_y^2)

The integrated probability for z to fall in the range μ - σ to μ + σ is 0.683. Other measures of width commonly encountered are: probable error (central region containing 0.50 of the probability) = μ ± 0.67σ; mean absolute deviation; E[|x - μ|] = 0.80σ; rms deviation = σ; half-width at half-maximum = 1.18σ.

The Gaussian gets its importance in large part from the central limit theorem: if a continuous random variable x is distributed according to any p.d.f. with finite mean and variance, then the sample mean x̄_n of n observations of x will have a p.d.f. that approaches a Gaussian as n increases. Therefore the end result Σ x_i ≡ n x̄_n of a large number of small fluctuations x_i will be distributed as a Gaussian, even if the x_i themselves are not.

The cumulative distribution (1.1) for a Gaussian with μ = 0 and σ^2 = 1 is given by the error function, erf(a), through the following ugly relation:

F(a; 0, 1) = 0.5 [1 + erf(a/√2)] (1.27)

The function erf(a) is tabulated in Ref. 2 and is available as a FORTRAN function on many computers [caution: other definitions of erf(a) are sometimes used]; for mean μ and variance σ^2 replace a by [(a - μ)/σ].

For x̄ a set of n (not necessarily independent) Gaussian random variables x_i arranged into a column vector, their joint p.d.f. is the multivariate Gaussian:

f(x̄; μ̄, V) = 1 / ((2π)^n |V|)^{1/2} × exp[-1/2(x̄ - μ̄)^T V^{-1} (x̄ - μ̄)], |V| ≠ 0 (1.28a)

where V is the covariance matrix of the x_i's, V_{ij} = Var(x_i) and V_{ij} = E[(x_i - μ_i)(x_j - μ_j)] ≡ ρ_{ij} σ_i σ_j, and |V| is the determinant of V. The quantity ρ_{ij} is the correlation coefficient for x_i and x_j; |ρ_{ij}| ≤ 1. For n = 2, this becomes

f(x_1, x_2; μ_1, μ_2, σ_1, σ_2, ρ) = 1 / (2π σ_1 σ_2 √(1 - ρ^2)) × exp{-1/2(1 - ρ^2) [(x_1 - μ_1)^2 / σ_1^2 - 2ρ(x_1 - μ_1)(x_2 - μ_2) / (σ_1 σ_2) + (x_2 - μ_2)^2 / σ_2^2]} (1.28b)

The special case σ_1 = σ_2 and ρ = 0 is called the Rayleigh distribution. If V is singular, there is a linear relation among some variables; in this case one usually wants to eliminate completely dependent variables and work in a smaller number of dimensions. The marginal distribution of any x_i is a Gaussian with mean μ_i and variance V_{ii}. V is n × n, symmetric, and positive definite. Therefore for any vector x̄, the quadratic form x̄^T V^{-1} x̄ = c traces an n-dimensional ellipsoid as x̄ varies for any given c > 0. If X_1 = (x_1 - μ_1)/σ_1, then c is a random variable obeying the χ^2(n) distribution, which is discussed in the following section. The probability that X̄ corresponding to a set of Gaussian random variables x̄, lies outside the ellipsoid characterized by a given value of c (= X̄^2) is given by Eq. (1.31) and may be read

of probability distributions, on the Web

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

from Fig. 1. For example, the "s-standard-deviation ellipsoid" occurs at $\epsilon = s^2$. For the two-variable case ($n = 2$) the point \bar{X} lies outside the one-standard-deviation ellipsoid with 61% probability, so both X_1 and X_2 lie inside the ellipsoid with 39% probability. This assumes that μ_1 and σ_1 are correct. For $X_1 = z_1/\sigma_1$, the ellipsoids of constant χ^2 have the same size and orientation but are centered at $\bar{\mu}$. The use of these ellipsoids as indicators of probable error is described in Sec. 2.4.1.

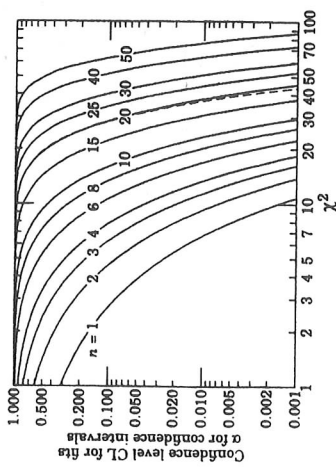


Fig. 1. χ^2 confidence level vs χ^2 for n degrees of freedom, as defined in Eq. (1.31). The curve for a given n expresses the probability that a value at least as large as χ^2 will be obtained in an experiment; e.g., for $n = 10$, a value $\chi^2 \geq 18$ will occur in 5% of a very large number of experiments. For a fit, CL is a measure of goodness-of-fit in that a good fit to a correct model is expected to yield a low χ^2 (Sec. 2.3.3). For a confidence interval, α measures the probability that the interval does not cover the true value of the quantity being estimated (Sec. 2.4). The dashed curve for $n = 20$ is calculated using the approximation of Eq. (1.32).

It is a characteristic of the multivariate Gaussian that $\rho_{ij} = 0$ is necessary and sufficient for z_i and z_j to be independent. For a given covariance matrix V , there always exist nonsingular $n \times n$ matrices H such that $HH^T = V$; H is usually upper or lower triangular in the most efficient algorithms. Then $\bar{z} = H^{-1}(\bar{x} - \bar{\mu})$ is a vector of n independent Gaussian random variables with zero mean and with covariance matrix equal to the identity.

1.3.5 The χ^2 distribution (continuous)
If z_1, \dots, z_n are independent Gaussian distributed random variables, the sum $z = \sum_{i=1}^n (z_i - \mu_i)^2 / \sigma_i^2$ is distributed as a χ^2 with n degrees of freedom $[\chi^2(n)]$.

$$f(z; n) = \frac{1}{2^n \Gamma(n/2)} z^{n/2-1} e^{-z/2}, \quad z \geq 0; \quad (1.29)$$

$$E(z) = n; \quad \text{Var}(z) = 2n. \quad (1.30)$$

Under a linear transformation to n dependent Gaussian variables x_i , the χ^2 at each transformed point retains its value; then $z = \bar{x}^T V^{-1} \bar{x}$, as in the previous section. For a set of z_i , each of which is $\chi^2(n_i)$, $\sum z_i$ is a new random variable which is $\chi^2(\sum n_i)$.

Fig. 1 shows the Confidence Level (CL) obtained by integrating the tail of the function given in Eq. (1.29) for n degrees of freedom:

$$CL(\chi^2) = \int_{\chi^2}^{\infty} f(z; n) dz; \quad (1.31)$$

this area is shown schematically in Fig. 2. It is equal to 1.0 minus the cumulative distribution function $F(z = \chi^2, n)$. It is useful in

evaluating the consistency of data with a model (see Sec. 2). The CL is the probability that a random repeat of the given experiment would observe a worse χ^2 , assuming the correctness of the model. It is also useful for confidence intervals for statistical estimators (Sec. 2.4), when one is interested in the unshaded area of Fig. 2.

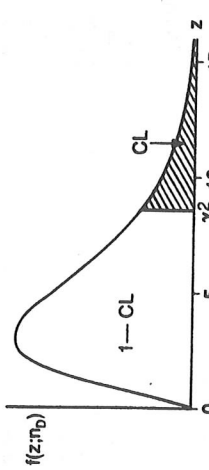


Fig. 2. Schematic illustration of the confidence level integral given in Eq. (1.31).

Since the mean of the χ^2 distribution is equal to the number of degrees of freedom, one expects to obtain $\chi^2 \approx n$ in a "reasonable" experiment. While caution is necessary because of the skewness of the distribution, the "reduced $\chi^2 \equiv \chi^2/n$ " is therefore a useful quantity. Figure 3 shows χ^2/n for useful CL's as a function of n . It contains the same information as Fig. 1, but is easier to read.

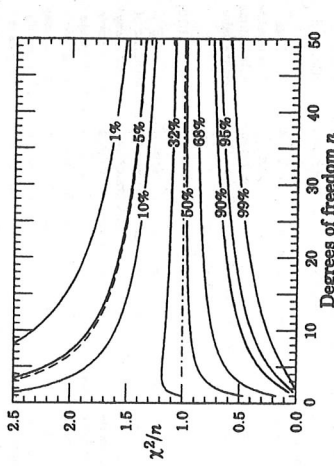


Fig. 3. Confidence limits as a function of the "reduced" $\chi^2 \equiv \chi^2/n$ and the number of degrees of freedom n . Curves are labeled by the probability of a measurement resulting in a value of χ^2/n greater than that given on the y axis; e.g., for $n = 10$, a value $\chi^2/n \geq 1.8$ will occur in 5% of a very large number of experiments. The dashed curve for CL = 5% is calculated using the approximation of Eq. (1.32).

It is commonly stated that for large n the CL is approximately given by [1.7]

$$CL \approx \frac{1}{\sqrt{2n}} \int_0^{\chi^2/n} e^{-x^2/2} dx, \quad (1.32)$$

where $y = \sqrt{2x^2} - \sqrt{2n-1}$. This approximation was used to draw the dashed curves in Fig. 1 (for $n = 20$) and Fig. 3 (for CL = 5%). However, all of the functions and their inverses are now readily available in standard mathematical libraries (such as IMSL, used to generate these figures), and so the approximation (and even such figures and tables) plays only a secondary role in practical problems.

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

1.3.6 Student's t (continuous)
Suppose that z and z_1, \dots, z_n are independent and normal with mean 0 and variance 1. We then define $t = z / \sqrt{z^2/n}$, and

$$t = z / \sqrt{z^2/n}. \quad (1.33)$$

The variable z thus belongs to a $\chi^2(n)$ distribution. Then t is distributed according to a Student's t distribution with n degrees of freedom:

$$f(t; n) = \frac{1}{\sqrt{\pi}} \frac{\Gamma[(n+1)/2]}{\Gamma(n/2)} \left(1 + \frac{t^2}{n}\right)^{-\frac{n+1}{2}} \quad (1.34)$$

and

$$E(t) = 0 \quad \text{for } n > 1; \quad \text{Var}(t) = \frac{n}{n-2} \quad \text{for } n > 2. \quad (1.35)$$

Here $\Gamma(k)$ is the gamma function, equal to $(k-1)!$ if k is an integer. Student's t distribution resembles a Gaussian distribution with wide tails. As $n \rightarrow \infty$, the distribution approaches a Gaussian, and if $n = 1$, the distribution is Cauchy, or Breit-Wigner. The mean is finite for $n > 1$ and the variance is finite for $n > 2$, so for $n = 1$ or $n = 2$, t does not obey the central limit theorem.

As an example, consider the sample mean $\bar{x} = \sum x_i/n$ and the sample variance $s^2 = \sum (x_i - \bar{x})^2 / (n-1)$ for normally distributed random variables x_i with unknown mean μ and variance σ^2 . The sample mean has a Gaussian distribution with a variance σ^2/n , so the variable $(\bar{x} - \mu) / \sqrt{\sigma^2/n}$ is normal with mean 0 and variance 1. Similarly, $(n-1)s^2 / \sigma^2$ is independent of this and is χ^2 distributed with $n-1$ degrees of freedom. The ratio

$$t = \frac{(\bar{x} - \mu) / \sqrt{\sigma^2/n}}{s / \sqrt{n-1}} = \frac{\bar{x} - \mu}{s / \sqrt{n}} \quad (1.36)$$

distributes as $f(t; n-1)$. The unknown true variance σ^2 cancels, and t can be used to test the probability that the true mean is some particular value μ .

The distribution (1.34) is written such that n is not required to be an integer. A Student's t distribution with noninteger $n > 0$ is useful in certain applications.

1.3.7 The gamma distribution (continuous)
If a process generating events as a function of x (e.g., space or time) satisfies conditions (a)-(c) of the Poisson distribution, then the x distance from an arbitrary starting point (which may be some particular event) to the k th event is belongs to a gamma distribution:

$$f(x; \lambda, k) = \frac{\lambda^k x^{k-1} e^{-\lambda x}}{\Gamma(k)}, \quad 0 < x < \infty. \quad (1.37)$$

$\Gamma(k)$ is the gamma function, equal to $(k-1)!$ if k is an integer. The Poisson parameter μ is λ per unit x :

$$E(x) = k/\lambda; \quad \text{Var}(x) = k/\lambda^2. \quad (1.38)$$

The special case $k = 1$ is called the exponential distribution. A sum of k exponential random variables x_i is distributed as $f(\sum x_i; \lambda, k)$. Eq. (1.37) allows $k > 0$ to be noninteger. If $\lambda = 1/2$ and $k = n/2$, the gamma and $\chi^2(n)$ distributions are identical.

2. STATISTICS

2.1 General

A probability density function with known parameters enables us to predict the frequency with which a random variable will take on a particular value (if discrete) or lie in a given range (if continuous). In parametric statistics we have the opposite problem of estimating the parameters of the p.d.f. from a set of actual observations.

We refer to the true p.d.f. as the population; the data form a sample from this population. A statistic is any function of the data, plus known constants, which does not depend upon any of the unknown parameters. A statistic is a random variable if the data have random errors. An estimator is any statistic whose value is intended as a meaningful guess for the value of an unknown parameter; we denote estimators with hats, e.g., $\hat{\theta}$.

Often it is possible to construct more than one reasonable estimator. Let θ represent the true value of a parameter to be estimated; θ is a vector if there is more than one parameter. Then if $\hat{\theta}$ is an estimator for θ , desirable properties for $\hat{\theta}$ are: (a) Unbiased, bias $b = E(\hat{\theta}) - \theta$, where the expectation value is taken over a hypothetical set of similar experiments in which $\hat{\theta}$ is constructed the same way. The bias may be due to statistical properties of the estimator or to systematic errors in the experiment. If we can estimate the average bias b we usually subtract it from $\hat{\theta}$ to obtain a new $\hat{\theta} \equiv \hat{\theta} - b$. However, b may depend upon θ or other unknowns, in which case we usually try to choose an estimator which minimizes its average size. (b) Minimum variance; the minimum possible value of $\text{Var}(\hat{\theta})$ is given by the Rao-Cramr-R-Frechet bound:

$$\text{Var}_{\min} = [1 + \partial b / \partial \theta]^2 / I(\theta); \quad (2.1)$$

$$I(\theta) = E \left[\frac{\partial}{\partial \theta} \ln f(x; \theta) \right]^2$$

The sum is over all data and b is the bias; if any, the x_i are assumed independent and distributed as $f(x_i; \theta)$, and the allowed range of x must not depend upon θ . The ratio $\epsilon = \text{Var}(\hat{\theta}) / \text{Var}(\theta)$ is the efficiency. An efficient estimator (with $\epsilon = 1$) exists only for certain cases. The square root of the variance expresses the expected spread of $\hat{\theta}$ about its average value, as would be observed in a large number of repeats of the same measurement. (c) Minimum mean-squared error (mse); $\text{mse} = E[(\hat{\theta} - \theta)^2] = V(\hat{\theta}) + b^2$. The mse combines the error due to any bias quadratically with the variance, which expresses only the spread about $E(\hat{\theta})$, as distinct from θ , the true value. (d) Robust; a robust estimator is not sensitive to errors in our assumptions, e.g., to departures from the assumed p.d.f. due to such factors as noise.

These criteria (and others) allow us to evaluate any procedure for obtaining $\hat{\theta}$. In many cases these criteria conflict. The bias, variance, and mse may depend on the unknown θ . In this case the optimum prescription for $\hat{\theta}$ may depend on the range in which we assume θ to lie.

Following are techniques in common use for obtaining estimators and their standard errors $\sigma(\hat{\theta}) = \sqrt{\text{Var}(\hat{\theta})}$. When the conditions of the central limit theorem are satisfied, the interval $\hat{\theta} \pm \sigma(\hat{\theta})$ forms a 68.3% confidence interval. This is a random interval in that its endpoints depend upon the randomly sampled data; its meaning here will be taken to be that in 68.3% of all similar experiments the interval will include the true value θ . One should be aware that in most practical cases the central limit theorem is only approximately satisfied and accordingly confidence intervals which depend on that are only approximate. Confidence intervals are discussed in Section 2.4 below.

2.2 Data with a common mean

(1) Suppose we have a set of N independent measurements y_i assumed to be unbiased measurements of the same unknown quantity μ with a common, but unknown, variance σ^2 resulting from measurement error. Then

$$\hat{\mu} = \frac{1}{N} \sum_{i=1}^N y_i; \quad (2.2)$$

$$\hat{\sigma}^2 = \frac{1}{N-1} \sum_{i=1}^N (y_i - \hat{\mu})^2 = \frac{N}{N-1} (E(y^2) - \hat{\mu}^2) \quad (2.3)$$

are unbiased estimators of μ and σ^2 . The variance of $\hat{\mu}$ is σ^2/N . If the common p.d.f. of the y_i is Gaussian, these statistics are independent. Then, for large N , the variance of $\hat{\sigma}^2$ is $2\sigma^4/N$. If the y_i are Gaussian or N is large enough that the central limit theorem applies, then $\hat{\mu}$ is an efficient estimator for μ . Otherwise $\hat{\mu}$ is sometimes subject to large fluctuations, e.g., if the p.d.f. for y_i has long tails. In this case the median of the y_i may be a more robust estimator for μ , provided the median and mean are expected to lie at the same point in the p.d.f. for y . For Gaussian y , the median has asymptotic (large- N) efficiency

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

7/N ~ 0.64. The Student's t distribution provides an example in which there are large tails. In this case, for large N the efficiency of the sample mean relative to the sample mean is (oo, oo, 1.62, 1.12, 0.96, 0.80, 0.64) for (1, 2, 3, 4, 5, 8, oo) degrees of freedom.

If theta is known, mu is given in Eq. (2.2) is still the best estimator for mu. If mu is known, substitute it for mu in Eq. (2.3) and replace N-1 by N to obtain a somewhat better estimator theta-hat.

(2) If the yi have different, known, variances sigma_i^2, then mu-hat = 1/w sum w_i y_i (2.4)

is an unbiased estimator for mu with smaller variance than Eq. (2.2), where w_i = 1/sigma_i^2 and w = sum w_i. The variance of mu-hat is 1/w.

2.3.1 General

"From a theoretical point of view, the most important general method of estimation so far known is the method of maximum likelihood. [1] We suppose that a set of independently measured quantities X came from a p.d.f. f(X; theta), where theta is an unknown set of parameters. The method of maximum likelihood consist of finding the set of values of theta, theta-hat, which maximizes the joint probability density for all the data, given by

L(theta-hat) = product f(X_i; theta-hat) (2.5)

where L is called the likelihood. It is usually easier to work with ln L, and since both are maximized for the same set of theta-hat, it is sufficient to solve the likelihood equation

partial ln L / partial theta = 0 (2.6)

The solution is called the maximum likelihood estimate of theta-hat. The importance of the approach is shown by the following proposition, proved in Ref. 1:

If an efficient estimate theta-hat of theta exists, the likelihood equation will have a unique solution equal to theta-hat.

In evaluating theta-hat, it is important that any normalization factors in the f's which involve theta-hat be included. However, we will only be interested in the maximum of L and in ratios of L at different theta-hat's; hence any multiplicative factors which do not involve the parameters we want to estimate may be dropped; this includes factors which depend on the data but not on theta-hat.

If the solution to Eq. (2.6) is at a maximum, partial ln L / partial theta will have negative slope in its vicinity. In many practical problems, one often uses nonlinear algorithms for finding the maximum, and must be alert in various possibilities for error: (a) Eq. (2.6) may yield a minimum, therefore one must check the second derivative; (b) there may be more than one maximum--one must try to find the global maximum; (c) the global maximum may lie at a boundary of the physical region, in which case Eq. (2.6) will not find it.

If an unbiased, efficient estimator exists, this method will find it. If partial ln L / partial theta is linear in the vicinity of the root, an efficient estimator is guaranteed; other efficient cases are discussed in the literature. For large data samples, the central limit theorem will usually assure this condition in some significant neighborhood of zero; hence the estimator is usually efficient in that case, provided certain conditions are met (e.g., that the solution does not lie on a boundary). In this case, the neighborhood of the maximum ln L is a downward-curving parabola and L is proportional to a Gaussian. The results of two or more experiments may be combined by forming the product of the L's, or the sum of the ln L's.

Under a one-to-one change of parameters from theta-hat to phi-hat(theta-hat), the maximum likelihood estimate is simply phi-hat = phi(theta-hat), given the solution for theta-hat. That is, the maximum likelihood solution for phi-hat is found by simple substitution of theta-hat into the transformation equation. It is possible that the new solution phi-hat will be a biased solution for the true value of phi even if theta-hat is not biased, and vice-versa. In the asymptotic limit (of large amounts of data) both phi-hat and phi-hat will (usually) converge to unbiased solutions, but at different rates.

Except in special cases like the least-squares method, the value of the likelihood function at the solution does not necessarily tell us whether the final fit was a sensible description of the data or not. To evaluate this, one may: (a) prepare histograms of the data projected onto various axes and make chi^2 (or other) comparisons with the fitted model projected upon the same axes; and/or (b) do numerous Monte Carlo simulations of the experiment under the hypothesis that the fitted parameters are correct, fit each of these, and compare the experimental likelihood (or ln L) with those obtained from these simulations. If the experimental likelihood is lower than that of some agreed-upon fraction of these results, one should question the appropriateness of the p.d.f. f. At the same time one can check for bias in the solution.

2.3.2 Error estimates

The covariance matrix V may be estimated from

V_nm = - (partial^2 ln L / partial theta_n partial theta_m) | theta-hat^-1 (2.7)

If partial ln L / partial theta is linear, the "expectation" operation in Eq. (2.7) has no effect because the second derivative of ln L is constant. Otherwise, it may be approximated by taking the average of the quantity in square brackets over a range of theta_n and theta_m near the solution. For complex cases it may be more practical to evaluate s-standard-deviation errors from the contour

ln L(theta-hat) = ln L_max - s^2 / 2, (2.8)

where ln L_max is the value of ln L at the solution point (compare with chi^2(theta-hat) = chi^2_min + 1 and the discussion in the least-squares case below). The extreme limits of this contour parallel to the theta_n axis give an approximate s-standard-deviation confidence interval in theta_n. These intervals may not be symmetric and they may even consist of two or more disjoint intervals. This procedure gives one standard-deviation errors in theta_n equal to sqrt(V_nm) of Eq. (2.7) if the estimator is efficient. If it is not efficient, the level of confidence implied by the value of s is only approximate.

2.3.3 Method of least squares

By far the most common case of the maximum likelihood approach is the method of least squares. We suppose a set of N measurements at points x_i. The ith measurement y_i is assumed to be chosen from a Gaussian distribution with mean F(x_i; theta-hat) and variance sigma_i^2. Then

-1/2 ln L = chi^2 = sum (y_i - F(x_i; theta-hat))^2 / sigma_i^2 (2.9)

Finding the set of parameters theta-hat which maximizes L is equivalent to finding the set which minimizes chi^2.

At the outset it should be said that the method of least squares is sometimes applied in cases where the distribution is not Gaussian or is not known to be Gaussian. In such cases it can still be used, but it is then not a special case of the maximum likelihood method, and the theorems having to do with that approach no longer apply.

In many practical cases one further restricts the problem to the situation in which F(x_i; theta-hat) is a linear function of the theta_m's,

F(x_i; theta-hat) = sum a_nm f_n(x_i) (2.10)

where the f_n are k linearly independent functions (e.g., 1, x, x^2, ... range of x. We require k <= N, and at least k of the x_i must be distinct. We wish to estimate the linear coefficients a_nm. Later we will discuss the nonlinear case.

If the point errors epsilon_i = y_i - F(x_i; theta-hat) are Gaussian, then the minimum chi^2 will be distributed as a chi^2 random variable with n = N - k degrees of freedom. We can then evaluate the goodness-of-fit (confidence level) from Figs. 1 or 3, as per the earlier discussion. The confidence level expresses the probability that a worse fit would be obtained in a large number of similar experiments under the assumptions that: (a) the model y = sum a_nm f_n is correct and (b) the errors epsilon_i are Gaussian and unbiased with variance

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

If y is not linear in the fitting parameters a_nm, or if the errors epsilon_i depend upon y and therefore on a_nm, the solution vector may have to be found by iteration of Eqs. (2.12)-(2.14) or Eq. (2.15b). The same results may be obtained by numerical techniques from the sum of squares, chi^2, directly, if we have a reasonable first guess theta-hat_0 for the solution vector:

theta-hat = theta-hat_0 - (partial^2 chi^2 / partial theta^2)^-1 | theta-hat_0 (2.19a)

and

V_2 = 2 (partial^2 chi^2 / partial theta^2)^-1 | theta-hat (2.19b)

where partial^2 chi^2 / partial theta is a k-element vector whose mth element is partial^2 chi^2 / partial theta_m, partial^2 chi^2 / partial theta^2 is a k x k matrix with mth element partial^2 chi^2 / partial theta_m partial theta_n, and all derivatives are to be evaluated at the points indicated. If "chi^2" is a true chi^2, the second-derivative matrix is independent of theta-hat; therefore the shape of the chi^2 as a function of theta-hat is a paraboloid and Eq. (2.19a) will give the solution immediately. Otherwise one may need to iterate Eq. (2.19a) to arrive at a solution (Newton-Raphson method).

Note that in Eq. (2.15b), one needs only a matrix proportional to V_2 to find theta-hat. Hence, for example, if the variances sigma_i^2 of the errors are unknown but assumed equal and independent, and E(epsilon_i) = 0, one can still solve for theta-hat. These can be estimated from the residuals, r_i = y_i(x_i) - y_i, where y_i(x_i) is the fitted curve at x_i, because study of the r_i enables one to estimate V_2. In addition, the residuals can be used to look for evidence of bias such as trends in the data not incorporated in the model [3].

Note that the errors on the solution theta-hat are independent of the value of chi^2 at minimum--they depend only upon the shape about the minimum. Eq. (2.19b) implies that one-standard-deviation limits on the elements of theta-hat are given by the set of theta-hat such that

chi^2(theta-hat) = chi^2_min + 1; (2.20)

compare with Eq. (2.8) for the general maximum-likelihood case. This equation, which defines a contour in theta-hat space, is often convenient for estimating errors in applications of least-squares techniques to nonlinear cases, where the second derivative [Eq. (2.19b)] may be a rapidly varying function of theta-hat. In general, contours at a standard deviation may be found by replacing the 1 in Eq. (2.20) by s^2. If the problem is highly nonlinear, all such contours are at best only approximations to desired exact confidence regions which would have some given probability of covering the true value of theta-hat. It may be that Eq. (2.20) will define a set of disjoint regions. In addition, iteration of Eq. (2.19a) may require sophisticated techniques [8] to reach convergence in a practical amount of computation. For example, in cases involving many variables in theta-hat, especially if the correlations are not small, simplex or other techniques which do not involve explicit calculation of derivatives are often to be preferred. Such techniques are designed to find their way through complicated nonlinear problems without diverging to infinite theta-hat (unless the minimum is actually at infinity).

Least-squares estimation requires that an error matrix V_2 be known (a matrix proportional to V_2 will suffice to find an estimator). For counting experiments it is therefore necessary to group the data in bins in order to associate a Poisson error with each bin. In this case theta-hat is the bin height and the error depends on the expectation value of the theory in each bin. theta-hat, as estimated by the best fit of the model. Thus the requirements of the Gauss-Markov theorem are not satisfied, since the errors are not fixed. Many experiments arrange the bins to contain enough expected events (say > 7 or 8) that the Gaussian approximation to the Poisson (Sec. 1.3.3) is accurate, in which case the expected error is the square root of the theoretical height and "chi^2" is approximately a true chi^2. If an approximate error is used, based on the actual observed height, theta-hat rather than the theoretical height, theta-hat, the Gauss-Markov conditions would be satisfied except that a bias favoring downward fluctuations would be

sigma_i^2. If this probability is larger than an agreed-upon value (0.001, 0.01, or 0.05 are common choices), the data are consistent with the assumptions, otherwise we may want to find improved assumptions. As for the converse, most people do not regard a model as being truly inconsistent unless the probability is as low as that corresponding to four or five standard deviations for a Gaussian (6 x 10^-3 or 6 x 10^-5; see Sec. 2.4.1). If the epsilon_i are not Gaussian, the method of least squares still gives an answer, but the goodness-of-fit test would have to be done using the correct distribution of the random variable which is still called "chi^2".

Finding the minimum of chi^2 in the linear case is straightforward:

-1/2 partial^2 chi^2 / partial theta_m = sum_i f_m(x_i) (y_i - sum_n a_nm f_n(x_i)) / sigma_i^2 (2.11) = sum_i y_i f_m(x_i) / sigma_i^2 - sum_n a_nm sum_i f_n(x_i) f_m(x_i) / sigma_i^2 (2.12)

With the definitions

theta-hat_m = sum_i y_i f_m(x_i) / sigma_i^2 (2.12)

and

(V_2^-1)_mn = sum_i f_n(x_i) f_m(x_i) / sigma_i^2, (2.13)

the k-element column vector of solutions theta-hat, for which partial^2 chi^2 / partial theta_m = 0 for all m, is given by

theta-hat = V_2 theta-hat (2.14)

More generally, the measured y_i's are not independent. Then the set of theta-hat's must be replaced by the N x N covariance matrix V_2. Then, if H is the N x k matrix with element H_m = f_m(x_i), the solution theta-hat is given by the solution to the normal equation

(H^T V_2^-1 H) theta-hat = H^T V_2^-1 y, (2.15a)

or, formally,

theta-hat = (H^T V_2^-1 H)^-1 H^T V_2^-1 y = D y, (2.15b)

where y is the N-element vector of measured y_i's. The normal equations may be solved by numerical methods much more computationally efficient than brute application of Eq. (2.15b). In particular, H^T V_2^-1 H is sometimes singular or nearly singular. In such cases there is at least one f_n which may be expressed as a linear combination of others (or nearly so) which evaluated at the data points. The best procedure is usually to drop such functions from the expansion (or set theta-hat_n = 0). See Press [6], Maindonald [9], or Basilevsky [10] for discussions.

In terms of the k x N matrix D, the standard covariance matrix for the theta-hat is estimated by

V_2 = D V_2 D^T (2.16)

If the measured y_i's are independent, V_2 is diagonal with ith element sigma_i^2 and V_2 is obtained from Eq. (2.13) above.

The expected covariance [see Eq. (1.8)] of theta-hat and theta-hat_m is estimated by

E[(theta-hat_m - theta-hat_m) (theta-hat - theta-hat)] = (V_2)_mm (2.17)

Even when the y_i's are independent (diagonal V_2), theta-hat and theta-hat_m may not be (non)diagonal V_2. For the model function y = sum a_nm f_n(x_i), the estimated variance of an interpolated or extrapolated value of y at a point x is

E[(y - y-hat)^2] = sigma^2(y) = sum_n (V_2)_nm f_n(x) f_m(x) (2.18)

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

This is because a fluctuation in the data which goes down from the true expectation value will be assigned a smaller error and therefore a greater weight than an equal fluctuation upward. For bins with few events, a procedure that converges to the above when N_i^{obs} is large and yields correct error estimates for all N_i^{obs} is to define

$$\hat{\sigma}_i^2 = \sum_{j=1}^2 \frac{2(N_i^{obs} - N_i^{exp}) + 2N_i^{exp} \ln(N_i^{obs}/N_i^{exp})}{N_i^{obs}} \quad (2.21)$$

This assumes that N_i^{obs} is the outcome of a Poisson process, with Poisson parameter $\mu = N_i^{exp}$, in the i th bin. In bins where $N_i^{obs} = 0$, the second term is zero. For any N_i^{obs} , a standard-deviation error estimates are constructed as in Eq. (2.20) and subsequent discussion. If we drop the requirement that χ^2 converge to a true χ^2 for large numbers of events in each bin, then minimizing $\chi^2 = \sum_{i=1}^n (N_i^{obs} - N_i^{exp})^2 / (N_i^{obs})$ will give the same answer and errors, with slightly faster execution, as the above.

In the more general maximum likelihood case, the small-number distributions are well known and there are no corresponding requirements concerning large numbers or even of binning.

Example: straight-line fit
 For the case of a straight-line fit, the following estimates of σ_1 and σ_2 , independent measurements y_i , the following estimates of σ_1 and σ_2 ,

$$\hat{\sigma}_1 = (S_y S_{xx} - S_x S_{xy}) / D, \quad (2.22)$$

$$\hat{\sigma}_2 = (S_1 S_{xx} - S_x S_{y1}) / D, \quad (2.23)$$

where

$$D = S_1 S_{xx} - S_x S_{y1}$$

The covariance matrix of the fitted parameters is:

$$\begin{pmatrix} V_{11} & V_{12} \\ V_{21} & V_{22} \end{pmatrix} = \frac{1}{D} \begin{pmatrix} S_{xx} & -S_x \\ -S_x & S_1 \end{pmatrix} \quad (2.24)$$

The estimated variance of an interpolated or extrapolated value of y at point x is:

$$(\hat{y} - y_{true})^2_{est} = \frac{1}{S_1} + \frac{S_x^2}{D} \left(\frac{x - S_x}{S_1} \right)^2 \quad (2.25)$$

2.4 Errors and confidence intervals
2.4.1 Gaussian errors

If the data are such that the distribution of the estimator ($\hat{\theta}$) satisfies the central limit theorem discussed in Sec. 1.3.4, the Gaussian distribution is the basis of the error analysis. If there is more than one parameter being estimated, the multivariate Gaussian is used. We define a confidence interval as being an interval constructed from the data to have probability at least $1 - \alpha$ (α is called the confidence coefficient) of covering the true value of θ . For the univariate case with known σ ,

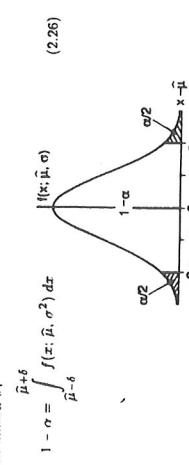
$$1 - \alpha = \int_{\hat{\mu} - \delta}^{\hat{\mu} + \delta} f(x; \hat{\mu}, \sigma) dx \quad (2.26)$$


Fig. 4. Illustration of a two standard-deviation confidence interval (unshaded) for a measurement of a single quantity with Gaussian errors. Integrated probabilities, defined by α , are as shown.

is the probability that the true value of μ will fall within $\pm \delta$ ($\delta > 0$) of the measured $\hat{\mu}$. This interval will cover μ in a fraction $1 - \alpha$ of all similar measurements. Fig. 4 shows a $\delta = 2\sigma$ confidence interval unshaded. The choice $\delta = \sqrt{\text{Var}(\hat{\mu})} = \sigma$ gives an interval called the standard error which has $1 - \alpha = 68.33\%$ if σ is known. Other frequently used choices for δ , in terms of α are:

α (%)	δ	α (%)	δ
31.73	1.28σ	20	1.28σ
4.55	2σ	10	1.64σ
0.27	3σ	5	1.96σ
6.3×10^{-3}	4σ	1	2.58σ
5.7×10^{-5}	5σ	0.1	3.29σ
2.0×10^{-7}	6σ	0.01	3.89σ

For other δ , find α as the ordinate of Fig. 1 on the $n = 1$ curve at $\chi^2 = (\delta/\sigma)^2$. We can set a one-sided (upper or lower) limit by excluding above $\hat{\mu} + \delta$ (or below $\hat{\mu} - \delta$), α 's for such limits are $1/2$ the values in the table above.

Note that we have increased confidence that the interval covers the true value as $1 - \alpha$ increases, or χ^2 increases. We must be careful to distinguish this case from the other major use of Fig. 1, evaluation of goodness-of-fit (Sec. 2.3.3). In that case we have increased confidence in the fit as χ^2 decreases. In an attempt to reduce possible confusion in this discussion, we will use the α notation (which corresponds to notation used in hypothesis testing [3]) when discussing confidence intervals and CL notation when discussing goodness-of-fit. Elsewhere in this Review, where the confusion between fit confidence level and interval (usually an upper or lower limit) confidence level does not arise, we follow the common practice of using "CL" to refer to the confidence level of the interval. This CL is understood to represent $1 - \alpha$.

If the variance σ^2 of the estimator is not known, but must be estimated from the data, then we need to incorporate the error in $\hat{\sigma}$ into our confidence interval using Student's t distribution. If we have N data points with which we estimate k parameters, the Gaussian approximation is adequate for $N - k \gg 1$. Otherwise replace δ by a factor $T\sigma$, T being defined by

$$1 - \alpha = \int_{-T}^T f(x; N - k) dx, \quad (2.27)$$

where f is defined in Eq. (1.34). T is tabulated in Ref. 2 and here:

$N - k$	1.84	6.31	12.71	13.97	63.66	235.78
1	1.32	2.92	4.30	4.53	9.92	19.21
2	1.20	2.35	3.18	3.31	5.84	9.22
3	1.14	2.13	2.78	2.87	4.60	6.62
4	1.11	2.01	2.57	2.65	4.03	5.51
5	1.05	1.81	2.23	2.28	3.17	3.96
10	1.03	1.72	2.09	2.13	2.85	3.42
20	1.00	1.64	1.96	2.00	2.58	3.00

For multivariate θ we must consider pairwise correlations. Assuming a multivariate Gaussian, Eq. (1.28a), and subsequent discussion the standard error ellipse for the pair $(\hat{\theta}_m, \hat{\theta}_k)$ may be drawn as in Fig. 5. The minimum χ^2 or maximum likelihood solution is at $(\hat{\theta}_m, \hat{\theta}_k)$. The standard errors σ_m and σ_k are defined as shown, where the ellipse is at a constant value of $\chi^2 = \chi^2_{min} + 1$ or $\ln L = \ln L_{max} - 1/2$. The angle of the major axis of the ellipse is given by

$$\tan 2\theta = \frac{2\hat{\theta}_m \sigma_m \sigma_k}{\sigma_m^2 - \sigma_k^2} \quad (2.28)$$

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

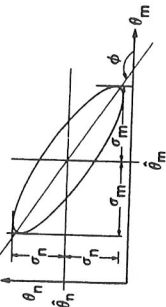


Fig. 5. Standard error ellipse for the estimators $\hat{\theta}_m$ and $\hat{\theta}_n$. In this case the correlation is negative.

For non-Gaussian or nonlinear cases, one may construct an analogous contour from the same χ^2 or $\ln L$ relations. Any other parameters $\hat{\theta}_i$, $i \neq m, n$, must be allowed freely to find their optimum values for every trial point.

For an unbiased procedure (e.g., least squares or maximum likelihood) being used to estimate k parameters θ_i , $i = 1, \dots, k$, the probability $1 - \alpha$ that the true values of all k lie within the standard deviation ellipsoid may be found from Fig. 1. Read the ordinate as α , the correct value of α occurs on the $n = k$ curve at $\chi^2 = \sigma^2$. For example, for $k = 2$, the probability that the true values of θ_1 and θ_2 simultaneously lie within the one-standard-deviation error ellipse ($\sigma = 1$), centered on $\hat{\theta}_1$ and $\hat{\theta}_2$, is 39%. This probability only assumes Gaussian errors, unbiased estimators, and that the model describing the data in terms of the θ_i is correct.

2.4.2 Gaussian errors—bounded physical region
 In certain statistical problems the true value of the parameter to be estimated, μ , is constrained to lie within a bounded physical region (e.g., the mass of a neutrino is bounded from below by 0). However, due to random measurement error, real measured values may or may not occur inside the physical region. For this case no completely satisfactory approach exists, but here we suggest a technique for obtaining limits within the physical region approximately at specified confidence levels. The "classical" statistical techniques of the previous section can still be used for confidence intervals at some α . However, such limits are useful mainly in the statistical sense where it is assumed that no bound exists. In bad cases, the limit may exclude the physical region entirely, or extend into it a small distance and create the false impression of a powerful limit close to the edge of the physical region.

We assume a measurement x , which represents one observation (or the result of combining multiple measurements as in Sec. 2.2) from a Gaussian of true (but unknown) mean μ and known, fixed, variance σ^2 . We estimate μ by $\hat{\mu} = x$ and attempt to construct a confidence interval for μ from the resultant Gaussian, as above. If $\hat{\mu}$ (Fig. 6), the result, while statistically perfectly correct as stated, is physically unsatisfactory.

If we assume μ is bounded from below by μ_{min} (the argument for μ bounded from above is similar), we may estimate a reasonable upper limit for μ at the $1 - \alpha$ (e.g., 90% or 95%) level by the following procedure: (1) renormalize the Gaussian probability distribution for x such that the integral of Eq. (1.24) with $\mu = \hat{\mu}$ over x from μ_{min} to infinity (i.e., over the physical region), unshaded in the figure below, is equal to $1 - \alpha$; (2) find the value μ_1 such that the integral over x of the renormalized distribution from μ_{min} to μ_1 is equal to the desired value of $1 - \alpha$; (3) set μ_1 to be the desired upper limit with confidence $1 - \alpha$. In fact, it can be shown that this is conservative, in the sense that the probability that this interval actually covers the true value of μ is $\geq 1 - \alpha$.

The "classical" approach as described above can be derived formally by the application of Bayes' theorem with the explicit assumption that all values of the parameter are equally probable. This means, for example, that limits on m^2 are different than limits on m . A recent treatment is given by James and Roos [11].

For $\mu - \mu_{min} \gg \sigma$, this technique, which may be applied for any measured x (physical or unphysical), converges smoothly to that of

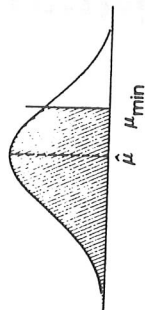


Fig. 6. An example of a bounded physical region with Gaussian errors. In this case the estimator $\hat{\mu}$ has fallen within the unphysical region due to random error.

the previous section since x is then effectively confined to the physical region.

One should exercise caution for values of x which lie many standard deviations outside the physical region. It may be that the particular probability model (Gaussian with variance σ^2) may not be a correct description of the measurement process (e.g., the true variance may have unanticipated components and be $> \sigma^2$, or there may be a bias), in which case confidence levels of this sort will not be correct.

If $\hat{\mu} < \mu_{min}$, some authors prefer to use a fixed upper limit calculated for $\hat{\mu} = \mu_{min}$ or $\hat{\mu} = \mu_{min} + \sigma$, rather than allow the upper limit to decrease as $\hat{\mu}$ decreases. In any case, averaging of experiments requires that $\hat{\mu}$ and its variance be quoted, in addition to any upper limits, even if $\hat{\mu}$ is unphysical.

2.4.3 Poisson processes—upper limits
 Because the outcome of a Poisson process is an integral number of events, n_0 , it is usually not possible to set confidence intervals for the true Poisson parameter μ at a certain exact α . For large n_0 an approximate interval can be set using the Gaussian approximation, Sec. 1.3.3, and the techniques of Sec. 2.4.1.

For small n_0 we can define an upper limit N for μ as being that value of μ such that it would be at least $1 - \alpha$ (e.g., 90% or 95%) probable that a random observation of n would then lie above the observed n_0 . Thus

$$1 - \alpha = \sum_{n=n_0+1}^{\infty} f(n; N); \quad \alpha = \sum_{n=0}^{n_0} f(n; N). \quad (2.29)$$

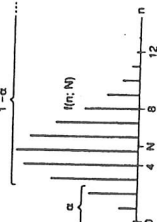


Fig. 7. Illustration of Eq. (2.29) Poisson probabilities for an assumed mean of N . With an observed count $n_0 = 2$, $N = 5.3$ as shown gives summed probability $1 - \alpha = 90\%$.

Fig. 7 illustrates the case with $n_0 = 2$ and $1 - \alpha = 90\%$, for which it may be shown that $N = 5.3$. For any given n_0 and desired α we can obtain N from the χ^2 Confidence Level figure because of a relation between the Poisson and the χ^2 : read the ordinate as α , find χ^2 on the curve for $n = 2(n_0 + 1)$, then $N = \chi^2/2$. Some useful values are:

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

Poisson upper limits N for n_0 observed events

n_0	$\alpha = 2.30$	$\alpha = 10\%$	$\alpha = 5\%$	$\alpha = 10\%$	$\alpha = 5\%$
0	2.30	3.00	6	10.53	11.84
1	3.89	4.74	7	11.77	13.15
2	5.32	6.30	8	13.00	14.44
3	6.68	7.75	9	14.21	15.71
4	7.99	9.15	10	15.41	16.96
5	9.27	10.51			

The meaning of these upper limits is that, for a given true μ , the probability is at least $1 - \alpha$ that one will observe n_0 which will result in N which is $\geq \mu$. The probability for that to occur may be higher than $1 - \alpha$; for example, if $\mu \leq 2.30$ a "90%" upper limit will actually exceed μ 100% of the time. Note from Eq. (2.28) that for $n_0 = 0$, $N = -\ln(1 - \alpha)$.

2.4.4 Poisson processes with background [12]

If we observe n_0 events in a Poisson process which has two components, signal and background, estimating a limit on the signal is more complicated. Let μ_S be the unknown mean (the Poisson parameter) for the signal and μ_B be the known mean for the sum of all backgrounds. Assume μ_B is known with negligible error; however we don't know n_B , the actual number of events resulting from the background. We do know that $n_B \leq n_0$. If $\mu_B + \mu_S$ is large, the Gaussian approximation to the Poisson distribution interval or limits is usually adequate, and one can define confidence intervals or limits as above, assuming $\bar{n}_B \approx \mu_B$ and therefore $\bar{\mu}_S = n_0 - \mu_B$ with variance equal to n_0 (larger than μ_S to allow for the error in \bar{n}_B).

Otherwise an upper limit can be defined by extension of the argument of the preceding section. Let N be the desired upper limit on μ_S with confidence coefficient α . Set N to be that value of μ_S such that any random repeat of the current experiment with $\mu_S = N$ and the same μ_B would observe more than n_0 events in total and would have $n_B \leq n_0$, all with probability $1 - \alpha$. For any assumed N and μ_B we can calculate this probability:

$$1 - \alpha = 1 - e^{-\mu_B} \sum_{n=0}^{n_0} \frac{(\mu_B + N)^n}{n!} e^{-\mu_B} \sum_{n=0}^{n_0} \frac{\mu_B^n}{n!} \quad (2.30)$$

We adjust N to obtain a desired α . For $\mu_B = 0$ this converges to (2.29). As in that case (see the last paragraph of Section 2.4.3) this gives a conservative upper limit in that for any given true μ_S we get a true probability $\geq 1 - \alpha$ that $N \geq \mu_S$, averaged over a large set of identically performed experiments. For $\alpha = 0.10$, Fig. 8 shows N as a function of n_0 and μ_B .

Averaging of experiments and other comparisons require that n_0 and μ_B be quoted and the technique used for upper limit extraction be given. If $\mu_B \gg n_0$ the experimenter should question the probability of observing n_B as that n_0 . If this is very small the background, μ_B , may not have been calculated properly and the upper limit for μ_S obtained under these assumptions may be too low. For example, in Fig. 8, the dashed portions of the curves lie in the region where n_0 is expected to exceed the observed value 90% of the time (or more), even in the complete absence of signal. In these regions one should be cautious about accepting the results of the measurement. As in the Gaussian case (2.4.2), whenever $n_0 < \mu_B$ some experimenters may prefer to use N calculated as if $n_0 \approx \mu_B$ rather than the smaller value obtained from the observed n_0 .

2.5 Propagation of errors

Suppose we have a set of N random variables y_i which may be direct measurements or derived estimators δ_i , and we have a covariance matrix $V(y)$ for these. We can make a transformation to a different

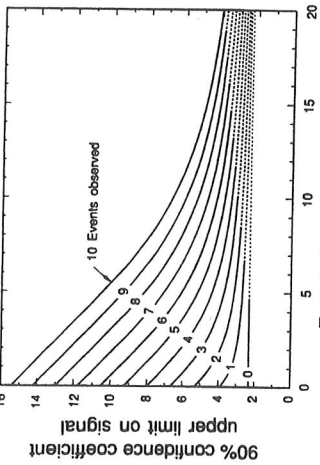


Fig. 8. 90% confidence coefficient upper limit on the number of signal events as a function of the expected number of background events. For example, if the expected background is 8 events and 5 events are observed, then the signal is 4.0 (approximately) or less with 90% confidence. Dashed portions indicate regions where it is to be expected that the number observed would exceed the number actually observed $\geq 90\%$ of the time, even in the complete absence of signal.

set of variables $f_j \equiv f_j(y)$, $j = 1, \dots, M$ ($M \leq N$) and obtain best estimates for the f_j from

$$\hat{f}_j \approx f_j(\hat{y}) + \sum_{k,n} \frac{\partial f_j}{\partial y_k} \left[\frac{\partial^2 f_j}{\partial y_k \partial y_n} \right]^{-1} V_{nm}(\hat{y}) \quad (2.31)$$

with covariance matrix

$$V_j(\hat{f}) \approx \sum_{n,m} \frac{\partial f_j}{\partial y_n} \left[\frac{\partial f_j}{\partial y_m} \right]^{-1} V_{nm}(\hat{y}) \quad (2.32)$$

For a single-valued function, f of a single measurement y with variance σ^2 (i.e., $M = 1, N = 1$), this becomes

$$\hat{f} \approx f(\hat{y}) + \frac{1}{2} \sigma^2 f''(\hat{y}) \quad (2.33)$$

where the primes denote differentiation with respect to y , evaluated at \hat{y} . These approximations are based on a Taylor expansion of f about the true value of y . If f is approximately linear in y over a range of roughly $\pm \sigma(y)$, the approximation is good and the second-order terms in (2.31) and (2.33) can be neglected. This is what is usually done. However, if linearity is badly violated (e.g., $f \propto 1/y$ and \hat{y} is no more than a few σ from zero), it should be recognized that propagation of errors will give very approximate results. In such cases $f \approx f(\hat{y})$ may be a biased estimator for f even if \hat{y} is unbiased for y , and the second-order terms in (2.31) and (2.33) will help to reduce that bias.

3. MONTE CARLO TECHNIQUES

Monte Carlo techniques are used to simulate on a computer random calculations are based upon pseudorandom numbers, a reproducible sequence of numbers generated on the open interval (0,1) in such a way that they satisfy various statistical tests for a uniform distribution, with independent numbers. (Caution: some commercial random number generators fill the closed interval [0,1]. The occurrence of 0 or 1 can sometimes cause problems for the algorithms below). No such numbers are truly uniform and independent. Many commercial random number generators sacrifice randomness in favor of speed. It

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

is not rare that unforeseen correlations will introduce non-negligible errors in the results. A useful test for this is to recalculate the same results with a different algorithm for the pseudorandom numbers. To improve the performance of an existing generator one may use the *Bray-Durham algorithm* (see Ref. 8 for discussion): (a) Initialize by generating and storing N (e.g., $N = 97$) random numbers in an array u , using the available generator. Generate a new random number u and save it. (b) On the next call, use this u as an address $j = 1 + (\text{integer part of } Nu)$ to select u_j as the random number to be returned. Also save this u_j as u for the next call. Replace u_j in the array with a new random number using the available generator. On the next call, go to (b).

A second problem sometimes encountered in computations requiring long sequences of random numbers is that all pseudorandom number generators will eventually begin over and repeat the same sequence. One may choose algorithms which minimize the number used. One may also use two or three different generators in different parts of the program.

Monte Carlo simulations of complex processes break them down into a sequence of steps. At each step a particular outcome is chosen from a set of possibilities according to a certain p.d.f. To do this we must transform our uniform random numbers into random numbers sampled from different distributions on different ranges.

Two techniques are in wide use to do this. We will discuss only single variable cases; multiple variable cases use straightforward extensions of these techniques. We assume we are in possession of a random number u chosen from a uniform distribution on (0,1).

3.1. Inverse transform method

If the desired probability density function is $f(x)$ on the range $-\infty < x < \infty$, its cumulative distribution function (expressing the probability that $x \leq \alpha$) is given by Eq. (1.1). If α is chosen with probability density $f(\alpha)$, then the integrated probability up to point α , $F(\alpha)$, is itself a random variable which will occur with uniform probability density on [0,1], ignoring the endpoints, we can then find a unique x distributed as $f(x)$ for $F(x)$ continuous, for a given u if we set

$$u = F(x), \quad (3.1)$$

provided we can find an inverse of F , defined by

$$x = F^{-1}(u), \quad (3.2)$$

as is illustrated in Fig. 9

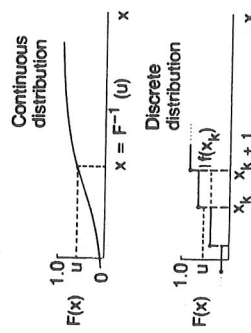


Fig. 9. Use of a random number u chosen from a uniform distribution (0,1) to find a random number x from a distribution with cumulative distribution function $F(x)$.

For a discrete distribution, $F(x)$ will have a discontinuous jump of size $f(x_k)$ at each allowed x_k , $k = 1, 2, \dots$. Choose u from a uniform distribution on (0,1) as before. Find x_k such that

$$F(x_{k-1}) < u \leq F(x_k) \equiv \text{Prob}(x \leq x_k) = \sum_{j=1}^k f(x_j); \quad (3.3)$$

then x_k is the value we seek (note: $F(x_0) \equiv 0$).

3.2 Acceptance-rejection method (Von Neumann)
Very commonly an analytical form for $F(x)$ is unknown or too complex to work with, so that obtaining an inverse as in Eq. (3.2) is impractical. We suppose that for any given value of x the probability density function $f(x)$ can be computed and further that enough is known about $f(x)$ that we can enclose it entirely inside a shape which is C times an easily generated distribution $h(x)$ as illustrated in Fig. 10.

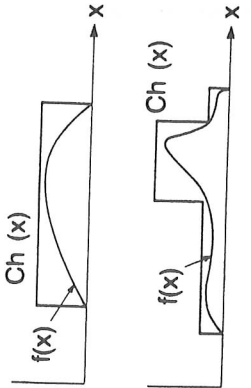


Fig. 10. Illustration of the acceptance-rejection method. Random points are chosen inside the upper bounding figure, and rejected if the ordinate exceeds $f(x)$. Lower figure illustrates importance sampling.

Frequently $h(x)$ is uniform or is a normalized sum of uniform distributions. Note that both $f(x)$ and $h(x)$ must be normalized to unit area and therefore the proportionality constant $C > 1$. To generate $f(x)$, first generate a candidate x according to $h(x)$. Calculate $f(x)$ and the height of the envelope $Ch(x)$, generate u and test if $uCh(x) \leq f(x)$. If so, accept x ; if not, reject x and try again. If we regard x and $uCh(x)$ as the abscissa and ordinate of a point in a two-dimensional plot, these points will populate the entire area $Ch(x)$ in a smooth manner; then we accept those which fall under $f(x)$. The efficiency is the ratio of areas, which must equal $1/C$; therefore we must keep C as close as possible to 1.0. Therefore we try to choose $Ch(x)$ to be as close to $f(x)$ as convenience dictates, as in the lower part of Fig. 10. This practice is called importance sampling, because we generate more trial values of x in the region where $f(x)$ is most important.

3.3 Algorithms

Many algorithms for generating common distributions are given by Rubinstein (1981) [13], Devroye (1986) [14], Press (1986) [8], Walk (1987) [15], and Everett (1983) [16]; a few of these are reproduced here. For many distributions alternative algorithms exist, varying in complexity, speed, and accuracy. For time-critical applications, these algorithms may be coded in-line to remove the significant overhead often encountered in making function calls. Variables named "u" are assumed to be independent and uniform on (0,1).

3.3.1 Sine and cosine of random angle

Generate u_1 and u_2 . Then $v_1 = 2u_1 - 1$ is uniform on $(-1,1)$, and $v_2 = u_2$ is uniform on (0,1). Calculate $r^2 = v_1^2 + v_2^2$. If $r^2 > 1$, start over. Otherwise, the sine (S) and cosine (C) of a random angle are given by

$$S = 2v_1v_2/r^2 \quad \text{and} \quad C = (v_1^2 - v_2^2)/r^2$$

3.3.2 Gaussian distribution

If u_1 and u_2 are uniform on (0,1), then

$$z_1 = \sin 2\pi u_1 \sqrt{-2 \ln u_2} \quad \text{and} \quad z_2 = \cos 2\pi u_1 \sqrt{-2 \ln u_2}$$

are independent and Gaussian distributed with mean 0 and $\sigma = 1$. There are many faster variants of this basic algorithm. For example, construct $v_1 = 2u_1 - 1$ and $v_2 = 2u_2 - 1$, which are uniform on $(-1,1)$.

PROBABILITY, STATISTICS, AND MONTE CARLO (Cont'd)

(calculator) $r^2 = v_1^2 + v_2^2$, and if $r^2 > 1$ start over. If $r^2 < 1$, it is uniform on $(0, 1)$. Then

$$z_1 = r_1 \sqrt{\frac{-2 \ln r^2}{r^2}} \quad \text{and} \quad z_2 = r_2 \sqrt{\frac{-2 \ln r^2}{r^2}}$$

are independent numbers chosen from a normal distribution with mean 0 and variance 1. $z_1^2 = \mu + \sigma z_1$ distributes with mean μ and variance σ^2 .

For a multivariate Gaussian it often is simplest to find a transformation matrix H as described at the end of Sec. 1.3.4 and return $\bar{x} = H\bar{z} + \bar{\mu}$. For $n = 2$ it is convenient to choose H such that $x_1 = z_1\sigma_1 + \mu_1$ and $x_2 = V_{12} z_1/\sigma_2^2 + z_2 [\sigma_2^2 - V_{12}^2/\sigma_1^2]^{1/2} + \mu_2$, where $\sigma_i^2 = V_{ii}$.

3.3.3 $\chi^2(n)$ distribution

For n even, generate $n/2$ uniform numbers u_i ; then

$$v = -2 \ln \left(\prod_{i=1}^{n/2} u_i \right) \quad \text{is} \quad \chi^2(n).$$

For n odd, generate $(n-1)/2$ uniform numbers u_i and one Gaussian z as in 3.3.2; then

$$v = -2 \ln \left(\prod_{i=1}^{(n-1)/2} u_i \right) + z^2 \quad \text{is} \quad \chi^2(n).$$

For $n \geq 30$ the much faster Gaussian approximation for the χ^2 may be preferable; generate z as in 3.3.2 and use $v = [z + \sqrt{2n-1}]^2/2$; if $z < -\sqrt{2n-1}$ reject and start over.

3.3.4 Binomial distribution

If $p \leq 1/2$ in Eq. (1.20), iterate until a successful choice is made: begin with $k = 1$; compute $P_k = q^n$ [for $k \neq 1$, use $P_k \equiv P_{k-1} p / (k - p)$]; $P_k(1.20)$ and store P_k into B ; generate u . If $u \leq B$ accept $r_k = k$ and stop; otherwise increment k by 1 and compute next P_k and add to B ; generate a new u and repeat. If we arrive at $k = n + 1$, stop and accept $r_{n+1} = n$. If $p > 1/2$ it will be more efficient to generate r from $J(r; n, q)$, i.e., with p and q interchanged, and then set $r_k = n - r$.

3.3.5 Poisson distribution

Iterate until a successful choice is made: Begin with $k = 1$ and set $A = 1$ to start. Generate u . Replace A with uA ; if now $A < \exp(-\mu)$, where μ is the Poisson parameter, accept $n_k = k - 1$ and stop. (Otherwise increment k by 1, generate a new u and repeat, always starting with the value of A left from the previous try. For large $\mu (\geq 10)$, it may be satisfactory (and much faster) to approximate the Poisson distribution by a Gaussian distribution [Sec. 1.3.4] and generate z from $J(z; 0, 1)$; then accept $r = \max(0, \lfloor \mu + z\sqrt{\mu} - 0.5 \rfloor)$ where $\lfloor \rfloor$ signifies the greatest integer \leq the expression.

3.3.6 Student's t distribution

For $n > 0$ degrees of freedom (n not necessarily integer), generate z from a Gaussian with mean 0 and $\sigma^2 = 1$ according to the method of 3.3.2. Next generate y , an independent gamma random variate with $k = n/2$ degrees of freedom. Then $z = z\sqrt{2n}/\sqrt{y}$ is distributed as a t with n degrees of freedom.

For the special case $n = 1$, the Breit-Wigner distribution, generate u_1 and u_2 ; set $v_1 = 2u_1 - 1$ and $v_2 = 2u_2 - 1$. If $v_1^2 + v_2^2 \leq 1$ accept $z = v_1/v_2$ as a Breit-Wigner distribution with unit area, center at 0.0, and FWHM 2.0. Otherwise start over. For center M_0 and FWHM Γ , use $W = z\Gamma/2 + M_0$.

Revised April 1992.

1. H. Gränér, *Mathematical Methods of Statistics*, Princeton Univ. Press, New Jersey (1968).
2. M. Abramowitz and I. Stegun, eds., *Handbook of Mathematical Functions* (Dover, New York, 1972).
3. W.T. Eadie, D. Drijard, F.E. James, M. Roos, and B. Sadulek, *Statistical Methods in Experimental Physics* (North Holland, Amsterdam and London, 1971).
4. L. Lyons, *Statistics for Nuclear and Particle Physicists* (Cambridge University Press, New York, 1986).
5. S.L. Meyer, *Data Analysis for Scientists and Engineers* (John Wiley and Sons, Inc., New York, 1975).
6. A.G. Prodesen, O. Skjeggstad, and H. Tøffe, *Probability and Statistics in Particle Physics* (Universitetsforlaget, Oslo, Norway, 1979).
7. R.A. Fisher, *Statistical Methods for Research Workers*, 8th edition, Edinburgh and London (1941).
8. W.H. Press et al., *Numerical Recipes* (Cambridge University Press, New York, 1986).
9. W.H. Maindonald et al., *Statistical Computation* (John Wiley and Sons, Inc., New York, 1984).
10. A. Basilevsky et al., *Applied Matrix Algebra in the Statistical Sciences* (North Holland, New York, 1983).
11. F. James and M. Roos, *Phys. Rev. D* **44**, 299 (1991).
12. O. Helene, *Nucl. Instr. and Meth.* **212**, 319 (1983).
13. R.Y. Rubinstein, *Simulation and the Monte Carlo Method* (John Wiley and Sons, Inc., New York, 1981).
14. L. Devroye, *Non-Uniform Random Variate Generation* (Springer-Verlag, New York, 1986).
15. Ch. Waack, *Random Number Generation*, University of Stockholm Physics Department Report 1987-10-20 (Vers. 3.0).
16. C.J. Everett and E.D. Cashwell, *A Third Monte Carlo Sampler*, Los Alamos report LA-9721-MS (1983).

ELECTROMAGNETIC RELATIONS

Quantity	Gaussian CGS	SI
Charge:	$2.987\,924\,58 \times 10^9$ esu	$= 1\text{ C} = 1\text{ A s}$
Electron charge e :	$4.803\,206\,8 \times 10^{-10}$ esu	$= 1.602\,177\,33 \times 10^{-19}\text{ C}$
Potential:	$(1/299\,792\,458)$ statvolt (ergs/esu)	$= 1\text{ V} = 1\text{ J C}^{-1}$
Magnetic field:	10^4 gauss = 10^4 dyne/esu	$= 1\text{ T} = 1\text{ N A}^{-1}\text{m}^{-1}$
Lorentz force:	$\mathbf{F} = q(\mathbf{E} + \frac{v}{c} \times \mathbf{B})$	$\mathbf{F} = q(\mathbf{E} + \mathbf{v} \times \mathbf{B})$
Maxwell equations:	$\nabla \cdot \mathbf{D} = 4\pi\rho$ $\nabla \times \mathbf{H} - \frac{1}{c} \frac{\partial \mathbf{D}}{\partial t} = \mathbf{J}$ $\nabla \cdot \mathbf{B} = 0$ $\nabla \times \mathbf{E} + \frac{1}{c} \frac{\partial \mathbf{B}}{\partial t} = 0$	$\nabla \cdot \mathbf{D} = \rho$ $\nabla \times \mathbf{H} - \frac{\partial \mathbf{D}}{\partial t} = \mathbf{J}$ $\nabla \cdot \mathbf{B} = 0$ $\nabla \times \mathbf{E} + \frac{\partial \mathbf{B}}{\partial t} = 0$
Materials:	$\mathbf{D} = \epsilon \mathbf{E}$, $\mathbf{H} = \mathbf{B}/\mu$	$\mathbf{D} = \epsilon \mathbf{E}$, $\mathbf{H} = \mathbf{B}/\mu$
Permittivity of free space:	1	$\epsilon_0 = 8.854\,187 \dots \times 10^{-12}\text{ F m}^{-1}$
Permeability of free space:	1	$\mu_0 = 4\pi \times 10^{-7}\text{ N A}^{-2}$
Fields from potentials:	$\mathbf{E} = -\nabla V - \frac{1}{c} \frac{\partial \mathbf{A}}{\partial t}$ $\mathbf{B} = \nabla \times \mathbf{A}$	$\mathbf{E} = -\nabla V - \frac{\partial \mathbf{A}}{\partial t}$ $\mathbf{B} = \nabla \times \mathbf{A}$
Static potentials: (coulomb gauge)	$V = \sum_{\text{charges}} \frac{q_i}{r_i} = \int \frac{\rho(\mathbf{r}')}{ \mathbf{r} - \mathbf{r}' } d^3x'$ $\mathbf{A} = \frac{1}{c} \sum_{\text{currents}} \frac{I_i}{r_i} = \frac{1}{c} \int \frac{\mathbf{J}(\mathbf{r}')}{ \mathbf{r} - \mathbf{r}' } d^3x'$	$V = \frac{1}{4\pi\epsilon_0} \sum_{\text{charges}} \frac{q_i}{r_i} = \frac{1}{4\pi\epsilon_0} \int \frac{\rho(\mathbf{r}')}{ \mathbf{r} - \mathbf{r}' } d^3x'$ $\mathbf{A} = \frac{\mu_0}{4\pi} \sum_{\text{currents}} \frac{I_i}{r_i} = \frac{\mu_0}{4\pi} \int \frac{\mathbf{J}(\mathbf{r}')}{ \mathbf{r} - \mathbf{r}' } d^3x'$
Relativistic transformations: (\mathbf{v} is the velocity of the primed frame as seen in the unprimed frame)	$\mathbf{E}'_{\parallel} = \mathbf{E}_{\parallel}$ $\mathbf{E}'_{\perp} = \gamma(\mathbf{E}_{\perp} + \mathbf{v} \times \mathbf{B})$ $\mathbf{B}'_{\parallel} = \mathbf{B}_{\parallel}$ $\mathbf{B}'_{\perp} = \gamma(\mathbf{B}_{\perp} - \frac{v}{c^2} \mathbf{v} \times \mathbf{E})$	$\mathbf{E}'_{\parallel} = \mathbf{E}_{\parallel}$ $\mathbf{E}'_{\perp} = \gamma(\mathbf{E}_{\perp} + \mathbf{v} \times \mathbf{B})$ $\mathbf{B}'_{\parallel} = \mathbf{B}_{\parallel}$ $\mathbf{B}'_{\perp} = \gamma(\mathbf{B}_{\perp} - \frac{1}{c^2} \mathbf{v} \times \mathbf{E})$
	$\frac{1}{4\pi\epsilon_0} = c^2 \times 10^{-7}\text{ N A}^{-2} = 8.987\,55 \dots \times 10^9\text{ F m}^{-1}$; $\frac{\mu_0}{4\pi} = 10^{-7}\text{ N A}^{-2}$;	$c = \frac{1}{\sqrt{\mu_0\epsilon_0}} = 2.987\,924\,58 \times 10^8\text{ m s}^{-1}$