

Statistical Methods

Data Sets

Data set is typically a Sample, S

Normally a subset of a bigger population, P

Depend on some underlying distribution, p

Individual events are Independent and identically distributed, (i.i.d.)

P follows some "true" distribution

"Inter" or "estimate" props of P from Sample S.

Events S and p are drawn $P(\vec{x})$ with no dependence \rightarrow i.i.d

Histograms:

Collection of i.i.d events $X_i = \{x_1, x_2, \dots, x_N\}$

→ vis w/ hist predict bin edges

Cont events in each bin: $l_b \leq x < h_b$

Bin Contents:

$$N_b = \sum_i \begin{cases} 1 & \text{if } l_b \leq x_i < h_b \\ 0 & \text{otherwise} \end{cases} \quad \text{error: } \sigma_b = \sqrt{N_b} \quad \text{Poisson distribution}$$

$$\text{Can Scale Counts with a weight} \quad \left\{ \begin{array}{l} N_b = \sum_i \begin{cases} w_i & \text{if } l_b \leq x_i \leq h_b \\ 0 & \text{otherwise} \end{cases} \\ \sigma_b = \sqrt{\sum_{i \in b} w_i^2} \end{array} \right.$$

Look into Seaborn.

Histograms as probability distributions

Scaled to give prob density $\sum_b w_b p_b = 1$

Quantifying Data

Data Sample, S with N i.i.d events $\{x_1, x_2, x_3, \dots, x_N\}$

$$\text{mean: } \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \quad \text{Function mean: } \bar{f} = \frac{1}{N} \sum_{i=1}^N f(x_i)$$

If binned then average bin counts:

$$\bar{N}_b = \frac{1}{N} \sum_{j=1}^B n_j x_{c_j} \quad \begin{array}{l} \text{Central value} \\ \text{Counts in bin} \\ \text{total no. bins} \end{array} \quad \bar{f} = \frac{1}{N} \sum_{j=1}^B n_j f(x_{c_j})$$

\bar{x} is the "Sample mean" - estimation of 'population mean' - depends on 'true mean, N'

$$\hookrightarrow N = \lambda x$$

Measuring Spread

$$\text{Variance: } V(x) = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2 \quad (\text{Squared deviation from mean})$$

$$\text{Sample Variance} = \boxed{\bar{x}^2 - \bar{x}^2}$$

Sample Mean: \bar{x} 'An estimate'

True mean: N

$$\begin{aligned} \text{Proof: } V(x) &= \frac{1}{N} \sum_i^N (x_i - \bar{x})^2 \\ &= \bar{x}^2 - 2\bar{x}^2 + \bar{x}^2 \\ &= \frac{1}{N} \sum_i^N (x_i^2 - 2\bar{x}x_i + \bar{x}^2) \\ &= \bar{x}^2 - \bar{x}^2 \\ &= \frac{1}{N} \left(\sum_i^N x_i^2 - 2\bar{x} \sum_i^N x_i + \sum_i^N \bar{x}^2 \right) \end{aligned}$$

Function Variance $V(y) = \frac{1}{N} \sum_{i=1}^N (f(x_i) - \bar{f})^2 = \bar{f}^2 - \bar{f}^2$

Standard Deviation:

$$\sigma = \sqrt{V(x)} = \sqrt{\bar{x}^2 - \bar{x}^2}$$

Sample stand deviation = $\sqrt{\frac{1}{N} \sum (x_i - \bar{x})^2}$

Biased Estimator as used Sample data twice

For mean and the Std \rightarrow non independent of each other

these are sample estimates

True Variance / Std Deviation

$$V(x) = \langle x^2 \rangle - \langle x \rangle^2$$

$$\sigma = \sqrt{\langle x^2 \rangle - \langle x \rangle^2}$$

$$\sigma = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

Corrected Standard Deviation - Bessel Correction

$$S = \sqrt{\frac{1}{N-1} \sum_i (x_i - \bar{x})^2}$$

Higher Moments

n^{th} algebraic moment : $\alpha^n = \frac{1}{N} \sum_i x_i^n$

n^{th} central moment : $C^n = \frac{1}{N} \sum_i (x_i - \bar{x})^n$

1st Algebraic: mean

2nd Central: Variance

'Skew' (measures asymmetry)

Related to 3rd Central: $\gamma = \frac{1}{N \sigma^3} \sum (x_i - \bar{x})^3$

ensures Skew is dimensionless.

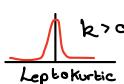


Curtosis (measures tails)

Related to 4th Central moment

$$k = \frac{1}{N \sigma^4} \left[\sum (x_i - \bar{x})^4 \right] - 3$$

ensures 0 for gaussian distribution



ensures dimensionality.

Covariance + Correlation

For datasets containing multiple variables $(x_i, y_i) = \{(x_1, y_1) \dots (x_n, y_n)\}$

Covariance: $\text{Cov}(x, y) = \frac{1}{N} \sum (x_i - \bar{x})(y_i - \bar{y})$

$$= \frac{1}{N} \sum_i (x_i y_i - \bar{x} \bar{y} - \bar{x} y_i - \bar{x} \bar{y})$$

$$= \bar{xy} - \bar{x} \bar{y}$$

For a data set of M variables

$M \times M$ Covariance (diagonal) matrix

$$\text{V}_{xy} = \text{Cov}(x, y) = \bar{xy} - \bar{x} \bar{y}$$

Covariance in words

+ve Covariance: If values of $x > \bar{x}$ then $y > \bar{y}$

-ve Covariance: Opposite.

Correlation (Dimensionless $-1 \leq \rho \leq 1$)

$$\rho(x, y) = \frac{\text{Cov}(x, y)}{\sigma_x \sigma_y} = \frac{\bar{xy} - \bar{x}\bar{y}}{\sigma_x \sigma_y}$$

For multi-dimensional:

$$V_{xy} = \rho_{xy} \sigma_x \sigma_y$$

ML algorithms

'attempt to fit some empirical function based on a set of features in data'

Sample data $\vec{x}_i = \{\vec{x}_1, \vec{x}_2, \vec{x}_3, \dots, \vec{x}_n\}$ $\xrightarrow{\text{typically } X}$
 \vec{x} is multidimensional - no. features

Overall (no. events, no. features)

Attempt to learn:

$$y_i = f(\vec{x}_i)$$

ML evaluation metric

- True Positive (TP) - False Positive (FP)
- True Negative (TN) - False Negative (FN)

All positive / Signal event: $P = TP + FN$

All negative / background events: $N = TN + FP$

All classified positive: $C_p = TP + FP$

All classified negative: $C_n = TN + FN$

True Positive rate: $TPR = \frac{TP}{TP + FN}$ } True Positive
} All positive

True Negative rate: $TNR = \frac{TN}{TN + FP}$

False Positive rate: $FPR = \frac{FP}{FP + TN}$

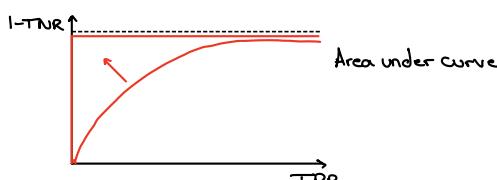
False Negative rate: $FNR = \frac{FN}{FN + TP}$

$$TPR + FNR = \frac{TP + FN}{TP + FN} = 1$$

Type-I Error - High False Positive Rate

Type II Error - High False Negative Rate

Receiver Operating Characteristic (ROC) Curve



Precision:

Fraction of all positively classified events that are correct

$$P = \frac{TP}{TP + FP}$$

Accuracy

\hookrightarrow Fraction Correctly Classified

$$\alpha = \frac{TP + TN}{P + N}$$

Error Rate

\hookrightarrow Fraction Incorrectly Classified

$$\epsilon = \frac{FP + FN}{P + N}$$

Purity:

$$\hookrightarrow P_p = \frac{TP}{TP + FP} \quad P_n = \frac{TN}{TN + FN}$$

Significance

$$\sigma = \frac{TP}{\sqrt{TP + FP}}$$

Signal to Noise (SNR)

$$SNR = \frac{TP}{FP}$$

F-Score

$$F = \frac{2TPR}{2TPR + FPR + FNR}$$

Avoiding Overfitting and Hyperparameter Tuning

↳ training Sample - fit the model

Validation Sample - optimize the hyperparameters

Test Sample - evaluate the performance

Cross-Validation

k -fold - Split data into k ($= 4$) folds: A, B, C, D

	Train	Test
Fold 1	BCD	A
Fold 2	ACD	B
Fold 3	ABD	C
Fold 4	ABC	D

Probabilistic Theory

Mathematical Probability

Kolmogorov Axioms (for X_i exclusive outcomes in Ω)

- 1) $P(X_i) \geq 0 \ \forall i$ (probability must be zero or positive)
- 2) $P(X_i \text{ or } X_j) = P(X_i) + P(X_j)$ (exclusive probabilities sum)
- 3) $\sum_i P(X_i) = 1$ (all probabilities must sum to unity)

Frequentist Probability

Given by relative frequency of particular outcome.

$$P(X) = \lim_{N \rightarrow \infty} \frac{n}{N}$$

N - total no.
n - observed outcome of X

True Probability - fixed probability but never attainable

Only applicable to repeatable experiments - keep relevant conditions stable
does not require any prior belief (purely experimental)

Frequentist vs Bayesian

↳ Both equally valid but cannot be mixed

Often depends on 'field'.

Should converge to the same inference in the 'asymptotic nirvana'. \rightarrow Sample size $\rightarrow \infty$

Bayesian Probability

Given by a degree of belief

$$P(X) = \text{degree of belief that } X \text{ happens}$$

Based on - Finetti's Coherent Bet:

$$P(X) = \text{largest amount you would bet}$$

must still obey Kolmogorov Axioms

Amount you would gain if win

Subjective as property of observer - 'amount willing to bet'.

depends on observer's knowledge so will change.

Properties of Probability

Addition - Non exclusive sets A and B in Ω

$$P(A \text{ or } B) = P(A) + P(B) - P(A \text{ and } B)$$

↓ Stop overcounting

Conditional Probability

$P(A|B)$ - probability of 'A given B'

Probability event known to be in B also in A

$$\begin{aligned} P(A \text{ and } B) &= P(A|B) P(B) \\ &= P(B|A) P(A) \end{aligned}$$

→ Note $P(A|B) \neq P(B|A)$

↳ Leads to Baye's Theorem

Independence:

A and B are independent if:

$$P(A|B) = P(A) \quad (\text{occurrence of } B \text{ has no input on } A)$$

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

If both A and B are independent:

$$P(A \text{ and } B) = P(A) \cdot P(B)$$

↳ i.e. Probability factorises
into two independent

Bayes Theorem (discrete)

$$P(A \text{ and } B) = P(A|B) P(B) = P(B|A) P(A)$$

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

for Bayes' theorem
often don't know $P(B)$

Law of Total Probability

Exclusive Sets $\{A_1, \dots, A_n\}$ and separate event B

$$P(B) = P(B \text{ and } A_1) \text{ or } (B \text{ and } A_2) \text{ or } \dots \text{ or } (B \text{ and } A_n)$$

$$\begin{aligned} &= \sum_i^N P(B \text{ and } A_i) \\ &= \sum_i^N P(B|A_i) P(A_i) \end{aligned}$$

Law of Total Probability and Bayes Theorem

$$P(A|B) = \frac{P(B|A) P(A)}{\sum_i^N P(B|A_i) P(A_i)}$$

Monty Hall Problem - Bayes Theorem

Stick:

$$P(\text{car in } A_1 \mid \text{host opens } B) = \frac{P(\text{open } B \mid \text{car in } A_1) P(\text{car in } A_1)}{P(\text{open } B)}$$

$$\begin{aligned} P(\text{open } B) &= \sum_i^N P(\text{open } B|A_i) P(A_i) \\ &= P(\text{open } B \mid \text{car in } A_1) P(\text{car in } A_1) \\ &\quad + P(\text{open } B \mid \text{car in } A_2) P(\text{car in } A_2) \\ &= \frac{1}{2} \times \frac{1}{3} + \frac{1}{2} \times \frac{1}{3} = \frac{1}{6} + \frac{1}{3} = \frac{1}{2} \end{aligned}$$

Switch:

$$P(\text{car in } C \mid \text{open } B) = \frac{P(\text{open } B \mid \text{car in } C) P(\text{car in } C)}{P(\text{open } B)} = \frac{\frac{1}{2} \times \frac{1}{3}}{\frac{1}{2}} = \frac{2}{3}$$

Union and Intersection

Union $A \cup B$ events in A or B or Both

$$A \cup B = \{x \mid x \in A \text{ or } x \in B\}$$

Intersection $A \cap B$ events in both A and B

$$A \cap B = \{x \mid x \in A \text{ and } x \in B\}$$

Formal Proof of Probability Theory

Prove if $P(A) = P_A$ then $P(\text{not } A) = 1 - P_A$

Declare all other outcomes in $\Omega - B$ $P(B) = P_B$

Using Kol III: $\sum_i P(X_i) = 1$; $P_A + P_B = 1$

$$P_B = 1 - P_A$$

Prove Probability in A or B but not both

Call $P(\text{not } A) = P(A^c)$ then $P(A \cup A^c) = 1$

$$P(A \cup B) = P(A \cup B) P(A \cup A^c)$$

$$P(A \cap (B \cap A^c)) = P(A) + P(B \cap A^c)$$

$$P(B) = P(B \cap A) + P(B \cap A^c)$$

$$P(B \cap A^c) = P(B) - P(B \cap A)$$

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

Random Variables

Example: rolling a dice

random event: roll

random variable: $\omega = \{1, 2, 3, 4, 5, 6\}$

Probabilities: $P(\omega) = \{P(\omega_1), P(\omega_2), \dots, P(\omega_6)\}$

Probability Mass Functions

↳ Discrete outcomes

Sum of probability over all possible values

$$\sum_i P(\omega_i) = 1 \quad \text{- Kolmogorov axiom}$$

Probability Density Functions

↳ Continuous outcomes - can take any real value

Example: heights of people

$$p(\omega) = \lim_{\Delta\omega \rightarrow 0} \frac{P(\omega - \Delta\omega/2 < \omega < \omega + \Delta\omega/2)}{\Delta\omega}$$

prob $\omega \rightarrow$ in a small interval $\Delta\omega$

Dimensions of density $[\omega^{-1}]$

Always normalised: 3rd Kolmogorov axiom

$$\int_{-\infty}^{\infty} p(\omega) d\omega = 1$$

For all ω ($\forall \omega$)

$$p(\omega) \geq 0$$

Probability that $a < \omega < b$:

$$P(a < \omega < b) = \int_a^b p(\omega) d\omega$$

Change of variables

If $\omega \sim f(\omega)$ (ω is distributed as $f(\omega)$) and $Y = h(\omega)$

then what is $g(Y)$? \rightarrow Prob distribution of Y

$$Y \sim g(Y)$$

Must conserve probability:

$$g(Y) dY = f(\omega) d\omega$$

If h is 'monotonic and invertible'

Use $\omega = h^{-1}(Y)$

$$g(Y) dY = f(h^{-1}(Y)) \left| \frac{d\omega}{dY} \right| dY$$

$$g(Y) = f(h^{-1}(Y)) \left| \frac{d\omega}{dY} \right| = f(\omega) \left| \frac{d\omega}{dY} \right| = \frac{f(\omega)}{|h'(\omega)|}$$

Jacobian of transform

'modulus of the transformation differential'

For disjoint regimes ($h(\omega)$ is not a 1-1 mapping)

$$g(Y) = \sum_{h^{-1}(Y)} \frac{f(\omega)}{|h'(\omega)|}$$

Extended to two indep. Variables X, Y

Distributed by joint density $f(x, y)$ - transformed by $U = U(x, y)$ $V = V(x, y)$

Joint density (U, V) :

$$G(U, V) = F(x, y) \left| J \left(\frac{x, y}{U, V} \right) \right| = F(x, y) \left| \begin{array}{cc} \frac{\partial x}{\partial u} & \frac{\partial x}{\partial v} \\ \frac{\partial y}{\partial u} & \frac{\partial y}{\partial v} \end{array} \right|$$

Extended to multi-variable

For $\vec{X} = (X_1, X_2, \dots, X_n)$ Variables with transformation functions $\vec{Y} = (Y_1, Y_2, \dots, Y_m)$

$$g(\vec{Y}) = \left| J \left(\frac{\vec{X}}{\vec{Y}} \right) \right| f(\vec{X})$$

Cumulative Distribution

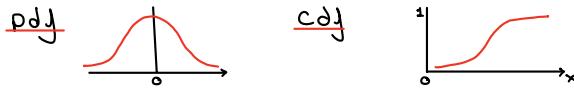
$$F(x) = \int_{-\infty}^x f(x') dx'$$

Note: $F(x_{\min}) = 0$

$$F(x_{\max}) = 1$$

pdj: $f(x)$ - Lower Case

Cdf: $F(x)$ - Upper Case



$$P(x < x') = \int_{x_{\min}}^{x'} f(x) dx = F(x')$$

$$P(x' < x < x'') = \int_{x'}^{x''} f(x) dx = F(x'') - F(x')$$

Percentage Point Functions F^{-1}

'The inverse of the C.d.f.'

maps us from random variable, x $[0, 1]$ \rightarrow random variable x

useful to generate samples from a distribution.

$$X = F^{-1}(p) \rightarrow \text{Uses } F^{-1} \text{ takes } p \text{ as an argument} \rightarrow \text{returns } x$$

Truncated Distribution

Standard distributions often $x \in [-\infty, \infty]$ or $x \in [0, \infty]$

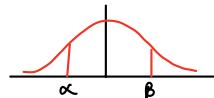
Often require: Restricted range $x \in [\alpha, \beta]$ \leftarrow Requires

Renormalisation

Truncated pdj

↳ uses truncated c.d.f to normalise

$$f(x) \rightarrow f'(x) = \frac{f(x)}{F(\beta) - F(\alpha)}$$



Truncated c.d.f

$$F(x) \rightarrow F'(x) = \frac{F(x) - F(\alpha)}{F(\beta) - F(\alpha)}$$

Truncated p.p.f.

$$\text{Def } G(x) = \frac{F(x) - \delta}{N} = p \quad S = F(x)$$

$$N = F(b) - F(a)$$

$$F(x) = Np + \delta$$

$$F^{-1}(F(x)) = x = F^{-1}(Np + \delta)$$

$$F^{-1}(p) = F^{-1}((F(b) - F(a))p + F(a))$$

Normal Distribution: notation

$$\text{P.d.f.} \sim N(\mu; \sigma^2)$$

$$\text{C.d.f.} \sim \Phi(\mu; \sigma^2)$$

$$\text{P.P.f.} \sim \Phi^{-1}(p; \sigma^2) \rightarrow \text{'erf'}$$

Joint p.d.f.

Prob dist of two random variables

Consider p.d.f. $x, y \sim f(x, y)$

$$\text{if } x, y \text{ are independent: } f(x, y) = g(x)h(y)$$

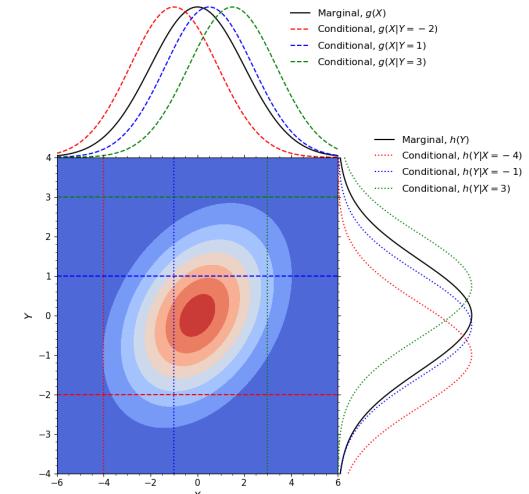
$g(x), h(y)$ - marginal distributions.

Working backwards:

Can 'integrate out' other variables

$$g(x) = \int f(x, y) dy \quad \text{"marginal in } x\text{"}$$

$$h(y) = \int f(x, y) dx \quad \text{"marginal in } y\text{"}$$



Conditional p.d.f.

↳ a 'slice through the joint density' 3D-Slice through condition in 1 dimension

$$\text{i) Prob of } x \text{ given } y: \quad g(x|y) = \frac{f(x, y)}{h(y)} = \frac{f(x, y)}{\int f(x, y) dx}$$

$$\text{ii) Prob of } y \text{ given } x: \quad h(y|x) = \frac{f(x, y)}{g(x)} = \frac{f(x, y)}{\int f(x, y) dy}$$

An example:

$$f(x, y) = \frac{e^{-x^2/2}}{\sqrt{2\pi}} \cdot e^{-y^2} \quad \begin{cases} x \in (-\infty, \infty) \\ y \in [0, \infty) \end{cases}$$

Unit Gaussian Decaying Exponential

Marginal in x

$$g(x) = \int_0^\infty f(x, y) dy = \left[-\frac{e^{-x^2/2}}{\sqrt{2\pi}} e^{-y^2} \right]_0^\infty = \frac{e^{-x^2/2}}{\sqrt{2\pi}} \left[-e^{-y^2} \right]_0^\infty = \frac{e^{-x^2/2}}{\sqrt{2\pi}}$$

Marginal in y

$$h(y) = \int_{-\infty}^\infty f(x, y) dx = \frac{e^{-y^2}}{\sqrt{2\pi}} \int_{-\infty}^\infty e^{-x^2/2} dx = e^{-y^2}$$

$\sqrt{2\pi}$ - See proof

$$f(x, y) = g(x) h(y) \quad \rightarrow \quad \text{Independent.}$$

$$g(x|y) = \frac{f(x, y)}{h(y)} = g(x)$$

$$h(x|y) = \frac{f(x, y)}{g(x)} = h(y)$$

In independent - no effect 'given'.

Integral Proof: $e^{-\alpha x^2}$

$$I = \int_{-\infty}^{\infty} e^{-\alpha x^2} dx \rightarrow I^2 = \int_{-\infty}^{\infty} e^{-\alpha x^2} dx \int_{-\infty}^{\infty} e^{-\alpha y^2} dy = \int_{-\infty}^{\infty} e^{-\alpha(x^2+y^2)} dy dx$$

Convert to polar: $x(r, \theta) = r \cos \theta \quad y(r, \theta) = r \sin \theta \quad r^2 = x^2 + y^2$

$$I^2 = \int_{-\infty}^{\infty} \int_{-\infty}^{\infty} e^{-\alpha(r^2)} dr dy$$

Jacobian: $(x, y) \rightarrow (r, \theta)$

$$dx dy = |J| dr d\theta$$

$$J = \begin{pmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{pmatrix}$$

$$= \begin{vmatrix} \frac{\partial x}{\partial r} & \frac{\partial x}{\partial \theta} \\ \frac{\partial y}{\partial r} & \frac{\partial y}{\partial \theta} \end{vmatrix} dr d\theta = \begin{vmatrix} \cos \theta & -r \sin \theta \\ \sin \theta & r \cos \theta \end{vmatrix} dr d\theta$$

$$= r(\cos^2 \theta + \sin^2 \theta) dr d\theta$$

$$= r dr d\theta$$

$$I^2 = \int_0^{2\pi} \int_0^{\infty} e^{-\alpha r^2} r dr d\theta \quad \text{using } u = r^2, \frac{du}{dr} = 2r \quad I^2 = \int_0^{2\pi} d\theta \int_0^{\infty} \frac{e^{-\alpha u}}{2} du = 2\pi \left[-\frac{e^{-\alpha u}}{2} \right]_0^{\infty} = \frac{\pi}{\alpha}$$

$$I = \sqrt{\frac{\pi}{\alpha}}$$

Bayes Theorem - Continuous

$$f(x|y) = g(x|y)h(y) = h(y|x)g(x)$$

$$g(x|y) = \frac{h(y|x)g(x)}{h(y)} \quad \text{holds for any r.v.s } x, y$$

for bayesian all parameters are random variables

The likelihood function

For a continuous random variable X and some data $\{x_1, x_2, \dots, x_n\}$

And a continuous range of hypothesis, Θ :

Assume definition $x \sim f(x)$ then

Likelihood is:

$$L(\Theta) = p(x|\Theta) = \prod_{i=1}^N f(x_i|\Theta)$$

The likelihood that we observe the data/measurement x_i given particular values of parameters Θ

Likelihood: $L(\Theta)$
Function of Θ

product of p.d.f's of assumed distribution over all observations

Not itself a p.d.f \rightarrow p.d.f is a function of (x) \rightarrow depends on data set.

A function of the parameters Θ (and the data) not the random variable x

Represents: agreement between our model and observed data.

Bayesian Inference

$$p(\theta|x) = \frac{p(x|\theta) p(\theta)}{p(x)}$$

$p(\theta|x)$ → Posterior distribution, probability distribution for parameter θ given the data we observed x

$$p(\theta|x) \propto p(x|\theta) p(\theta)$$

$p(x|\theta)$ → Likelihood Function, likelihood we observe data x given particular value, θ

$p(\theta)$ → The Prior, prior beliefs of θ (bayesian - prior belief required)

$p(x)$ → The Evidence, normalisation factor - ensures posterior is a p.d.f
- just a number

Statistical Inference

↳ deducing properties of a population or distribution from observed data

From posterior distribution

↳ Can determine: Inferred value (point estimate)

↳ Normally use 'maximum a posteriori' probability (MAP) i.e. mode

Uncertainty:

↳ Often High Density Interval (HDI)

(narrowest interval containing $(1-\alpha)\%$)

Properties of Distribution

Expectation

Random variable, x with p.d.f $f(x)$ ($x \sim f(x)$)

Function $g(x)$ operating on x

Expectation of $g(x)$ (continuous distribution)

$$E[g(x)] = \langle g(x) \rangle = \int g(x) f(x) dx$$

↳ Integral over all x

Expectation of $g(x)$ (discrete distribution)

$$E[g(x)] = \langle g(x) \rangle = \sum x g(x) P(x)$$

↳ Upper case (discrete)

Acts as a linear operator:

$$E[a g(x) + b h(x)] = a E[g(x)] + b E[h(x)]$$

Mean Density, N (of probability function)

↳ expectation value of x itself ($g(x) = x$)

$$E[x] = \langle x \rangle = N = \int x f(x) dx$$

Variance, σ^2

Expectation value of $(x-N)^2$

$$\sigma^2 = E[(x-N)^2]$$

σ , std deviation ↪ $= E[x^2 - 2N\bar{x} - N^2]$

$$= E[x^2] - 2N E[x] + E[N^2]$$

↳ $N^2 - \langle \text{constant} \rangle = \text{constant}$

$$= E[x^2] - 2N^2 + N^2 = E[x^2] - N^2$$

$$\sigma^2 = \int (x-N)^2 f(x) dx = \int x^2 f(x) dx - N^2$$

Moments

$$N_n = E[x^n] \rightarrow n^{\text{th}} \text{ algebraic moment}$$

$$\alpha_n = E[(x - N_1)^n] \rightarrow n^{\text{th}} \text{ Central moment.}$$

<u>Mean:</u> N_1	<u>Skew</u> $\gamma_1 = \sqrt{B_1} = N_3/N_2^{3/2}$
<u>Variance:</u> α_2	<u>Kurtosis</u> $\gamma_2 = \beta_2 - 3 = \frac{N_4}{N_2^2} - 3$

Covariance and Correlation

Between two random variables $X, Y - f(x, y)$

$$\begin{aligned} \text{Covariance: } V(x, y) &= E[(x - N_x)(y - N_y)] \\ &= \iint (x - N_x)(y - N_y) f(x, y) dx dy \\ &= \iint (xy - xN_y - yN_x - N_x N_y) f(x, y) dx dy \\ &= \iint xy f(x, y) dx dy - N_y \underbrace{\int x f(x, y) dx dy}_{N_x} - N_x \underbrace{\int y f(x, y) dx dy}_{N_y} + N_x N_y \\ &= \iint xy f(x, y) dx dy - N_x N_y \end{aligned}$$

$$V(x, y) = E[XY] - N_x N_y = E[XY] - E[X]E[Y]$$

$$\text{Correlation: } \rho(x, y) = \frac{V(x, y)}{\sigma_x \sigma_y} = \frac{E[XY]}{\sigma_x \sigma_y} - \frac{N_x N_y}{\sigma_x \sigma_y}$$

An example: $X \sim f(x) = Ne^{-x^2}$

$$1) \text{ Find } N: \int_{-\infty}^{\infty} Ne^{-x^2} dx = 1 \quad (\text{Normalised})$$

$$\text{Using } \int_{-\infty}^{\infty} e^{-x^2} dx = \sqrt{\pi} \rightarrow N = 1$$

$$X - f(x) = \frac{e^{-x^2}}{\sqrt{\pi}}$$

2) Find mean

$$N = \langle x \rangle = E[X] = \int_{-\infty}^{\infty} x N e^{-x^2} dx$$

Method 1: Odd function $\rightarrow 0$

$$\text{Method 2: } N = x^2 \frac{du}{dx} = 2x \quad du = 2x dx$$

$$\int x e^{-x^2} dx \rightarrow \int_{-\infty}^{\infty} \frac{e^{-u}}{2} du = \left[-\frac{e^{-u}}{2} \right]_{-\infty}^{\infty} = 0 \quad N = 0$$

3) Variance $E[(x - N)^2] \rightarrow E[x^2]$ as $N = 0$

$$E[X^2] = \int_{-\infty}^{\infty} \frac{x^2 e^{-x^2}}{\sqrt{\pi}} dx$$

$$\int N V' = UV - \int VU' \quad U = x \quad U' = 1 \quad V(x) = \frac{1}{2} \\ V' = x e^{-x^2} \quad V = -\frac{e^{-x^2}}{2x}$$

Proof: $\int_{-\infty}^{\infty} x^2 e^{-ax^2}$

$$N = x \rightarrow N = 1$$

$$V = \frac{e^{-ax^2}}{2a} \leftarrow V' = x e^{-ax^2}$$

$$= \int x^2 e^{-ax^2} dx = \left[-\frac{x e^{-ax^2}}{2} \right]_{-\infty}^{\infty} + \int_{-\infty}^{\infty} \frac{e^{-ax^2}}{2a} dx = \frac{1}{2a} \sqrt{\frac{\pi}{a}}$$

Standard Normal $f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$

- Shift in x by N $f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-N)^2}{2\sigma^2}}$ $x \rightarrow z = \frac{x-N}{\sigma}$

- Scale by $1/\sigma$

Characteristic Function

Fourier Transform of P.d.f. $f(x)$

For random variable $x \sim f(x)$ with C.d.f $F(x)$

$$\Psi(t) = E[e^{itx}] = \int_{-\infty}^{\infty} e^{itx} f(x) dx$$

Characteristic function

Completely defines $f(x)$

Can freely switch between $f(x) \leftrightarrow \Psi(t)$

Inverse Fourier transform

$$f(x) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \Psi(t) e^{-ixt} dt$$

'Characteristic function can completely define probability distribution'

Algebraic Moments (from Characteristic function)

- N_n - n^{th} moment around mean

Previously $N_n = E[x^n] = \int_{-\infty}^{\infty} x^n f(x) dx \rightarrow N_n$ can be obtained from n^{th} differential of $\Psi(t)$ at $t=0$

$$\Psi(t) = \int_{-\infty}^{\infty} e^{itx} f(x) dx$$

$\frac{d^n}{dt^n} \Psi(t) = \frac{d^n \Psi}{dt^n} = i^n \int_{-\infty}^{\infty} x^n e^{itx} f(x) dx \rightarrow$ 'Can generate moments by differentiating the Characteristic function'

$$\Psi_n(0) = i^n \int_{-\infty}^{\infty} x^n f(x) dx = i^n N_n = i^n E[x^n]$$

Advantage

Often much more useful to differentiate using Characteristic function than integrate P.d.f

Central Moment / Moment about mean

(Sub $V = x - N$)

Example: Normal Distribution

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

$$\Psi(t) = e^{wt} e^{-\frac{1}{2}\sigma^2 t^2} \rightarrow \text{proved below}$$

$$N = E[x] = i^{-1} \left[\frac{d \psi(t)}{dt} \right]_{t=0}$$

$$\frac{d \psi(t)}{dt} = iN e^{int} e^{-\frac{1}{2}\sigma^2 t^2} - e^{int} \frac{1}{2} \sigma^2 2t e^{-\frac{1}{2}\sigma^2 t^2}$$

$$= (iN - \sigma^2 t) \psi(t)$$

$$N = i^{-1} (iN - \sigma^2 t) \Big|_{t=0} = N \quad \underline{\text{Q.E.D}}$$

$$V(x) = E[(x-N)^2] = E[x^2] - N^2$$

$$\text{where } E[x^2] = i^{-2} \left[\frac{d^2 \psi(t)}{dt^2} \right]_{t=0}$$

$$\begin{aligned} \frac{d^2 \psi(t)}{dt^2} &= -\sigma^2 \psi(t) + (iN - \sigma^2 t) \frac{d \psi(t)}{dt} \\ &= -\sigma^2 \psi(t) + (iN - \sigma^2 t)^2 \psi(t) \end{aligned}$$

$$\text{Thus } E[x^2] = -1 (-\sigma^2 - N^2 - 2iN\sigma^2 t + \sigma^4 t^2) \psi(t) \Big|_{t=0}$$

$$\underline{E[x^2] = \sigma^2 - N^2}$$

$V(x) = E[x^2] - N^2 = \sigma^2$ \longrightarrow Characteristic function often easier to compute moments \longrightarrow Used to show N, σ can be undefined.

Relation to moments

Expand exponential

$$\psi(t) = E[e^{itx}] = E \left[\sum_{n=0}^{\infty} \frac{(itx)^n}{n!} \right] = \sum_{n=0}^{\infty} \frac{(it)^n}{n!} E[x^n] \quad \text{Using } e^{itx} = \sum_{n=0}^{\infty} \frac{(itx)^n}{n!}$$

$$= \sum_{n=0}^{\infty} \frac{(it)^n}{n!} N_n \quad \longrightarrow \quad \boxed{\psi(t), \text{ characteristic function can be expanded as the sum of the moments which appear as coefficients of } \frac{(it)^n}{n!}}$$

Moment Generating Function (MGF) $E[e^{itx}]$

Characteristic Function \longrightarrow Always exists

Disadvantages: Only exists if finite for all t (real)

Properties of Characteristic Function

$$\boxed{\psi_x(t) = E[e^{itx}]}$$

1) Sum of many independent random variables

$$\psi_z(t) \text{ where: } z = \sum_i x_i, \psi_z(t) = E[e^{itz}] = E[e^{it \sum_i x_i}]$$

$$\psi_z = \int e^{itz} \underbrace{j(z)}_{\text{Independent r.v.'s}} dz = \int e^{it \sum_i x_i} \cdot \underbrace{j(x_1) j(x_2) \dots j(x_n)}_{\text{Thus Separable p.d.f. } j(x_1, \dots, x_n) = j(x_1) \dots j(x_n)} \cdot d(x_1) d(x_2) \dots d(x_n)$$

$$= \int_{x_1, \dots, x_n} e^{itx_1} j(x_1) dx_1 \cdot e^{itx_2} j(x_2) dx_2 \dots e^{itx_n} j(x_n) dx_n$$

$$= \prod_{i=1}^n \int e^{itx_i} j(x_i) dx_i = \boxed{\prod_{i=1}^n \psi_{x_i}(t)}$$

Charac func for Sum of random variables

Is product of individual Charac functions

2) Random Variable multiplied by a constant

$$\psi_{\text{acc}}(t) = E[e^{it\text{acc}}] = E[e^{it\text{acc}}] = \psi_{\text{acc}}(at)$$

Proof of $\psi(t)$ for normal

Starting With Std Normal:

$$f(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

even function

$$\psi(t) = E[e^{itx}] = E[\cos(tx) + i\sin(tx)] = E[\cos(tx)] + iE[\sin(tx)] \rightarrow \text{useful expansion}$$

$$= \underbrace{\int_{-\infty}^{\infty} \cos(tx) f(x) dx}_{\text{Difficult to solve: trick, investigate first derivative of } \psi(t)} + i \underbrace{\int_{-\infty}^{\infty} \sin(tx) f(x) dx}_{\int_{-\infty}^{\infty} \text{odd} = 0}$$

Investigating trick

$$\frac{d}{dt} \psi(t) = \frac{d}{dt} E[e^{itx}] = E\left[\frac{d}{dt}(e^{itx})\right]$$

$$= E[ixe^{itx}]$$

$$= iE[x\cos(tx)] - E[x\sin(tx)]$$

$$= i \underbrace{\int_{-\infty}^{\infty} x \cos(tx) f(x) dx}_{0 \text{ (odd)}} - \underbrace{\int_{-\infty}^{\infty} x \sin(tx) f(x) dx}_{x f(x) = x \left(\frac{1}{\sqrt{2\pi}} e^{-x^2/2}\right)} = -\frac{d}{dx} \left(\frac{1}{\sqrt{2\pi}} e^{-x^2/2}\right)$$

$$= \int \sin(tx) \frac{d}{dx} f(x) dx$$

By parts

$$= -\frac{d}{dx} f(x)$$

$$= [\sin(tx) f(x)]_{-\infty}^{\infty} - \int t \cos(tx) f(x) dx$$

$$= -t \int_{-\infty}^{\infty} \cos(tx) f(x) dx$$

$$\text{Using } \psi(t) = \int_{-\infty}^{\infty} \cos(tx) f(x) dx$$

$$\frac{d}{dt} \psi(t) = -t \psi(t)$$

Condition $\psi(0) = E[e^0] = 1$

$$\text{General sol: } \psi(t) = e^{-t^2/2}$$

For non-standard normal

$$X \rightarrow Z = \frac{X - \mu}{\sigma}$$

$$\psi(t) = E[e^{itX}] = E[e^{it(\mu + \sigma Z)}] = e^{i\mu t} E[e^{i\sigma t Z}]$$

$$= e^{i\mu t} \psi_z(\sigma t) = e^{i\mu t} e^{-\sigma^2 t^2/2}$$

Common Probability Densities

Probability distribution depends on r.v.s, $\vec{\omega}$ and parameters, Θ

$$P(\vec{x}; \Theta) \quad \text{eg normal } P(x; \mu, \sigma) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}$$

Binomial Distribution (fixed no. outcomes, two outcomes)

n trials, $P(\text{success}) = p$, $P(\text{fail}) = 1-p$

P that first k are Success: \times Combinations that k success in n trials:

$$P = p^k (1-p)^{n-k}$$

$$C(n, k) = \binom{n}{k} = \frac{n!}{k!(n-k)!}$$

Binomial p.d.f

$$P(k; p, n) = \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

Mean: $N = E[k] = np$

$$N = E[k] = \sum_{k=0}^n \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

$$= \sum_{k=0}^n \frac{n! (n-1)!}{(k-1)! (n-k)!} p p^{k-1} (1-p)^{(n-1)-(k-1)} \quad \text{Sub } j = k-1$$

$$m = n-1$$

$$= np \sum_{j=0}^m \frac{m!}{j!(m-j)!} p^j (1-p)^{m-j} = np$$

Variance:

$$V(k) = np(1-p)$$

$$V(k) = E[(k - N)^2] = E[k^2] - E[k]^2 = E[k^2] - (np)^2$$

$$E[k^2] = \sum_{k=0}^n k^2 \frac{n!}{k!(n-k)!} p^k (1-p)^{n-k}$$

$$= \sum_{k=1}^n k n \frac{(n-1)!}{(k-1)!(n-k)!} p p^{k-1} (1-p)^{(n-1)-(k-1)} \quad \text{Sub } j = k-1$$

$$m = n-1$$

$$= np \sum_{j=0}^m j+1 \frac{m!}{j!(m-j)!} p^j (1-p)^{m-j} = np \left[\frac{\sum_{j=0}^m j!}{m!} + 1 \right] \quad \text{mp} = (n-1)p$$

$$= np((n-1)p + 1) = np(np - p + 1) = (np)^2 - np^2 + np$$

$$V(k) = E[k^2] - E[k]^2 = (np)^2 - np(np-1) - (np)^2 =$$

Standard Deviation:

$$\sigma = \sqrt{np(1-p)}$$

Poisson Distribution

Binomial where $n \rightarrow \infty$ but $np = \lambda$ is fixed

Binomial becomes:

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

$$\text{As } n \rightarrow \infty: \frac{n!}{(n-k)!} = n(n-1)(n-2) \dots (n+k-1) \xrightarrow{n \rightarrow \infty} n^k$$

$$\left(1 - \frac{\lambda}{n}\right)^{n-k} \longrightarrow \left(1 - \frac{\lambda}{n}\right)^n \longrightarrow e^{-\lambda}$$

$$P(k; \lambda) = \frac{n^k \lambda^k}{k! n^k} e^{-\lambda} = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$P(k; \lambda) = \frac{\lambda^k e^{-\lambda}}{k!}$$

$$\text{Mean: } E[k] = \sum_{k=0}^{\infty} k \frac{\lambda^k e^{-\lambda}}{k!} = \lambda e^{-\lambda} \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!}$$

$$= \lambda e^{-\lambda} \sum_{j=0}^{\infty} \frac{\lambda^j}{j!} = \lambda e^{-\lambda} e^{\lambda} = \lambda$$

$$E[k] = \lambda$$

Variance: $E(k^2) - E(k)^2$

$$E[k^2] = \sum_{k=0}^{\infty} k^2 \frac{\lambda^k e^{-\lambda}}{k!} = \lambda e^{-\lambda} \sum_{k=1}^{\infty} k \frac{\lambda^{k-1}}{(k-1)!}$$

$$= \lambda e^{-\lambda} \left[\sum_{k=1}^{\infty} \frac{(k-1)}{(k-1)!} \frac{\lambda^{k-1}}{(k-1)!} + \sum_{k=1}^{\infty} \frac{\lambda^{k-1}}{(k-1)!} \right]$$

$$\lambda \sum_{k=2}^{\infty} \frac{\lambda^{k-1}}{(k-2)!} e^{\lambda}$$

$$= \lambda e^{-\lambda} (\lambda e^{\lambda} + e^{\lambda}) = \lambda(\lambda + 1) = \lambda^2 + \lambda$$

$$\text{Var}(k) = \lambda^2 + \lambda - \lambda^2 = \lambda$$

$$\text{Var}(k) = \lambda$$

$$\sigma = \sqrt{\lambda}$$

Note: Error for weights per event

If each entry has its own weight, w_i

Count in bins

$$n_B = \sum_{i \in B} w_i$$

Compute Variance

$$\text{Var}(n_B) = \sum_{i \in B} w_i^2$$

In each bin

Uncertainty

$$\sigma(n_B) = \sqrt{\sum_{i \in B} w_i^2}$$

Normal/Gaussian Distribution

Standard Normal $\mu=0, \sigma=1 \sim N(0,1)$

$$P(x) = \frac{1}{\sqrt{2\pi}} e^{-x^2/2}$$

Properties of Normal Distribution

1) Any linear combination of norm dist r.v.s is also normal dist.

2) Sample mean and Variance are independent

3) If X_i are Standard normally distributed:

→ p.d.f is constant on the hypersphere ($\sum_i^N x_i^2 = \text{constant}$) → Unique for Normal dist

Non Standard Normal

Shift by μ ; Scale by $1/\sigma$ $Z = \frac{x-\mu}{\sigma}$

Characteristic Function

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

$$e_x(t) = e^{i\mu t} e^{-\frac{1}{2}\sigma^2 t^2}$$

Useful Integral for Normal

$$\int_{-\infty}^{\infty} e^{-ax^2} dx = \sqrt{\frac{\pi}{a}}$$

Cumulative Distribution Function (c.d.f) for normal/Gaussian

$$\Phi(x) = \frac{1}{2} \left[1 + \text{erf} \left(\frac{x-\mu}{\sqrt{2}\sigma} \right) \right]$$

$$\text{erf}(z) = \frac{2}{\sqrt{\pi}} \int_0^z e^{-t^2} dt$$

Inverse C.d.f Ξ^{-1}

Z-score

$$\text{or } F(x) = \Phi \left(\frac{x-\mu}{\sigma} \right) \text{ where } \Phi(z) = \frac{1}{\sqrt{2\pi}} \int_{-\infty}^z e^{-t^2/2} dt$$

Mean

$$E[x] = \int_{-\infty}^{\infty} \frac{x}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx$$

(Sub $z = \frac{x-\mu}{\sigma}$)

$$= \int_{-\infty}^{\infty} \frac{z}{\sigma \sqrt{2\pi}} e^{-\frac{z^2}{2}} dz + \int \frac{\mu}{\sigma \sqrt{2\pi}} e^{-\frac{z^2}{2}} dz = \mu$$

$$\mu \frac{\sqrt{2\pi}}{\sqrt{2\pi}} = \mu$$

Variance

$$E[x^2] = \int_{-\infty}^{\infty} \frac{x^2}{\sigma \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x-\mu}{\sigma}\right)^2} dx$$

(Sub $z = \frac{x-\mu}{\sigma}$ $\rightarrow x^2 = \sigma^2 z^2 + 2\mu \sigma^2 z + \mu^2$)

$$= \int_{-\infty}^{\infty} \frac{z^2}{\sigma^2 \sqrt{2\pi}} e^{-\frac{z^2}{2}} dz + \int_0^{\infty} \frac{2\mu z}{\sigma^2 \sqrt{2\pi}} e^{-\frac{z^2}{2}} dz + \int_{\infty}^{\infty} \frac{\mu^2}{\sigma^2 \sqrt{2\pi}} e^{-\frac{z^2}{2}} dz$$

$$= \sigma^2 + \mu^2$$

$$V(x) = E[x^2] - \mu^2$$

$$V(x) = \sigma^2$$

Multi-Variate Normal

For Independent r.v.s \vec{x} then \vec{x} random variable with means $\vec{\mu}$

$$p(\vec{x}; \vec{\mu}, \vec{\sigma}) = \prod_i p(x_i; \mu_i, \sigma_i)$$

$$= \prod_i \frac{1}{\sigma_i \sqrt{2\pi}} e^{-\frac{1}{2} \left(\frac{x_i - \mu_i}{\sigma_i}\right)^2}$$

We can include dependent r.v.s with a Correlation term:

term. $\vec{\sigma}^2 = \begin{pmatrix} \sigma_1^2 & & & \\ & \sigma_2^2 & & \\ & & \ddots & \\ & & & \sigma_n^2 \end{pmatrix}$ $\underline{\sigma}^2 = \begin{pmatrix} \sigma_1^2 & \rho_{12}\sigma_1\sigma_2 & \dots & \rho_{1n}\sigma_1\sigma_n \\ \rho_{12}\sigma_1\sigma_2 & \sigma_2^2 & & \vdots \\ \vdots & & \ddots & \\ \rho_{1n}\sigma_1\sigma_n & \dots & & \sigma_n^2 \end{pmatrix}$

Term in exponential becomes:

$$(\vec{x} - \vec{\mu})^T \underline{\sigma}^{-1} (\vec{x} - \vec{\mu}) \quad \text{where } \underline{\sigma} = \sigma^2$$

Multivariate Normal P.d.f.

$$p(\vec{x}; \vec{\mu}, \underline{\sigma}) = \frac{1}{(2\pi)^{n/2} \sqrt{|\underline{\sigma}|}} \exp \left[-\frac{1}{2} (\vec{x} - \vec{\mu})^T \underline{\sigma}^{-1} (\vec{x} - \vec{\mu}) \right]$$

Properties of Multi-Variate Normal

1) Contours of Constant probability where:

$$(\vec{x} - \vec{\mu})^T \underline{\sigma}^{-1} (\vec{x} - \vec{\mu}) = \text{Const.}$$

2) Any Slice through a m.v.n distribution is also a m.v.n (Proof below)

\rightarrow if $p(\vec{x})$ is a m.v.n at dim n

$p(\vec{x}_{n+1} | \vec{x}_n)$ is also a m.v.n in dim n-1

With Covariance matrix $\underline{\sigma}_{n-1}$ obtained by removing ith row and column from $\underline{\sigma}^{-1}$ and inverting

3) Any projection to a lower dimension is also a m.v.n

get $\underline{\sigma}_{n-1}$ by removing relevant row + column from $\underline{\sigma}$

Mathematical Proof 2)

Partition $\vec{x} = \begin{pmatrix} \vec{x}_1 \\ \vec{x}_2 \end{pmatrix}$ with dimensions $\begin{bmatrix} q_1 \times 1 \\ (N-q) \times 1 \end{bmatrix}$

Thus $\vec{N} = \begin{pmatrix} \vec{N}_1 \\ \vec{N}_2 \end{pmatrix}$ w/ $\begin{bmatrix} q_1 \times 1 \\ (N-q) \times 1 \end{bmatrix}$ and $\underline{\underline{\Sigma}} = \begin{pmatrix} \underline{\underline{\Sigma}}_{11} & \underline{\underline{\Sigma}}_{12} \\ \underline{\underline{\Sigma}}_{21} & \underline{\underline{\Sigma}}_{22} \end{pmatrix}$ w/ $\begin{bmatrix} q_1 \times q_1 & q_1 \times (N-q) \\ (N-q) \times q_1 & (N-q) \times (N-q) \end{bmatrix}$

the $p(\vec{x}_1 | \vec{x}_2) \sim \text{m.v.n } N(\vec{N}', \underline{\underline{\Sigma}}')$ with:

$$\vec{N}' = \vec{N}_1 + \underline{\underline{\Sigma}}_{12} \underline{\underline{\Sigma}}_{22}^{-1} (\vec{x}_2 - \vec{N}_2)$$

$$\underline{\underline{\Sigma}}' = \underline{\underline{\Sigma}}_{11} - \underline{\underline{\Sigma}}_{12} \underline{\underline{\Sigma}}_{22}^{-1} \underline{\underline{\Sigma}}_{11}$$

Example in 2D

r.v.s X, Y

$$\vec{N} = (N_x, N_y) \quad \underline{\underline{\Sigma}} = \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix}$$

$$|\underline{\underline{\Sigma}}| = \sigma_x^2 \sigma_y^2 - \rho^2 \sigma_x^2 \sigma_y^2 = \sigma_x^2 \sigma_y^2 (1 - \rho^2)$$

$$\sqrt{|\underline{\underline{\Sigma}}|} = \sigma_x \sigma_y \sqrt{1 - \rho^2}$$

$$\underline{\underline{\Sigma}}^{-1} = \frac{1}{\sigma_x^2 \sigma_y^2 (1 - \rho^2)} \begin{pmatrix} \sigma_y^2 & -\rho \sigma_x \sigma_y \\ -\rho \sigma_x \sigma_y & \sigma_x^2 \end{pmatrix}$$

$$\underline{\underline{\Sigma}}^{-1} \begin{pmatrix} X - N_x \\ Y - N_y \end{pmatrix} = \frac{1}{\sigma_x^2 \sigma_y^2 (1 - \rho^2)} \begin{pmatrix} \sigma_y^2 (X - N_x) & -\rho \sigma_x \sigma_y (X - N_x) \\ -\rho \sigma_x \sigma_y (Y - N_y) & \sigma_x^2 (Y - N_y) \end{pmatrix}$$

$$(X - N_x | \underline{\underline{\Sigma}}^{-1} \begin{pmatrix} V - N_x \\ Y - N_y \end{pmatrix}) = \frac{1}{\sigma_x^2 \sigma_y^2 (1 - \rho^2)} \left[\sigma_y^2 (X - N_x)^2 - \rho \sigma_x \sigma_y (Y - N_y)(X - N_x) - \rho \sigma_x \sigma_y (X - N_x)(Y - N_y) + \sigma_x^2 (Y - N_y)^2 \right]$$

$$= \frac{1}{1 - \rho^2} \left[\left(\frac{X - N_x}{\sigma_x} \right)^2 + \left(\frac{Y - N_y}{\sigma_y} \right)^2 - 2\rho \frac{(X - N_x)(Y - N_y)}{\sigma_x \sigma_y} \right]$$

$$\text{Thus: } p(X, Y) = \frac{1}{2\pi \sigma_x \sigma_y \sqrt{1 - \rho^2}} \exp \left[-\frac{1}{2(1 - \rho^2)} \left\{ \left(\frac{X - N_x}{\sigma_x} \right)^2 + \left(\frac{Y - N_y}{\sigma_y} \right)^2 - 2\rho \frac{(X - N_x)(Y - N_y)}{\sigma_x \sigma_y} \right\} \right]$$

Exponential Decay

$$p(x) = Ne^{-\lambda x}$$

$$\text{Normalise: } \int_0^\infty Ne^{-\lambda x} dx = 1$$

$$\frac{1}{N} = \left[-\frac{e^{-\lambda x}}{\lambda} \right]_0^\infty \quad \frac{1}{N} = \frac{1}{\lambda} \quad N = \lambda$$

P.d.f $p(x) = \lambda e^{-\lambda x}$

C.d.f $F(x) = \lambda \int_0^x e^{-\lambda x} dx = \lambda \left[-\frac{1}{\lambda} e^{-\lambda x} \right]_0^x$

$F(x) = 1 - e^{-\lambda x}$

$$U = \lambda x \quad V = e^{-\lambda x}$$

Mean: $N = E[x] = \int_0^\infty x \lambda e^{-\lambda x} dx$ $U' = \lambda$ $V' = -\frac{e^{-\lambda x}}{\lambda}$

$$= \left[-x e^{-\lambda x} \right]_0^\infty + \int_0^\infty e^{-\lambda x} dx = 0 + \left[-\frac{e^{-\lambda x}}{\lambda} \right]_0^\infty = \frac{1}{\lambda} = \tau$$

$N = E[x] = \frac{1}{\lambda} = \tau$

P.d.f alternative

$p(x) = \frac{e^{-x/\tau}}{\tau}$

Variance $E[x^2] - E[x]^2$

$$U = \lambda x^2 \quad V = e^{-\lambda x}$$

$$E[x^2] = \int x^2 \lambda e^{-\lambda x} dx \quad U' = 2\lambda x \quad V' = -\frac{e^{-\lambda x}}{\lambda}$$

$$= \left[-x^2 e^{-\lambda x} \right]_0^\infty + \int_0^\infty 2\lambda x e^{-\lambda x} dx$$

$$= 0 + \frac{2}{\lambda} E[x] = \frac{2}{\lambda^2}$$

$$\sqrt{V(x)} = \sqrt{E[x^2] - E[x]^2} = \sqrt{\frac{2}{\lambda^2} - \frac{1}{\lambda^2}} = \frac{1}{\lambda^2}$$

$\sqrt{V(x)} = \frac{1}{\lambda^2}$

$\sigma = \frac{1}{\lambda^2}$

Relationship to Poisson:

Poisson — Counting for n occurrences (in some interval)

Exp — length of time between occurrences.

Take Poisson w/ λt in time t

Expectation is λt with Poisson distribution

$p(x) = \frac{e^{-\lambda t} (\lambda t)^x}{x!}$ expectation λt

Probability of 0 occurrences in time t

$$P(0) = e^{-\lambda t}$$

Probability that event occurs within time t

$$P(T \leq t) = \underbrace{1 - e^{-\lambda t}}$$

C.d.f of exp.

Polynomial P.d.f's

Advantages: very flexible

Disadvantage: Go negative

↳ Kolmogorov's Axiom: $P(x) > 0$ for $\forall x$

Bernstein Basis (Basis of Polynomials)

↳ positive definite everywhere

defined for $0 \leq x \leq 1$ → can be easily transformed to any range $[a, b]$

$$x \rightarrow \frac{x-a}{x-b}$$

Bernstein Basis Polynomials

$$b_{i,n} = \binom{n}{i} x^i (1-x)^{n-i} \quad \text{for } i=0, \dots, n$$

degree $n \rightarrow$ degree $n+1$ polynomial.

$$b_{0,0} = 1$$

$$b_{0,1} = 1-x \quad b_{1,1} = x$$

$$b_{0,2} = (1-x)^2$$

$$b_{1,2} = 2x(1-x) \quad b_{2,2} = x^2$$

Change
given order
 n

Bernstein in generic range $[a, b]$:

$$b_{i,n}(x) = \frac{1}{(b-a)^n} \frac{n!}{i!(n-i)!} (x-a)^i (b-x)^{n-i}$$

$$\text{Derivative: } \frac{d b_{i,n}(x)}{dx} = \frac{n}{(b-a)} [b_{i+1,n-1}(x) - b_{i,n-1}(x)]$$

$$\text{Integral: } \int b_{i,n}(x) dx = \frac{(b-a)}{n+1} \sum_{j=i+1}^{n+1} B_{j,n+1}(x)$$

Def Integral: All basis funcs of same order have same
def integral over $[a, b]$

$$\int_a^b b_{i,n}(x) dx = \frac{b-a}{n+1}$$

The Bernstein Polynomial: $B_n = \sum_{i=0}^n C_i b_{i,n}$

$$P(x) = \sum_{i=0}^n C_i b_{i,n} \quad \int P(x) dx = 1$$

$$\frac{1}{N} = \int \sum_{i=0}^n C_i b_{i,n} dx = \sum C_i \int b_{i,n}(x) dx$$

$$= \sum_{i=0}^n C_i \frac{1}{n+1} = \frac{\sum C_i}{(n+1)}$$

$$P(x) = \frac{n+1}{\sum C_i} \sum C_i b_{i,n}(x)$$

χ^2 distribution

K Degrees of Freedom

↳ Distribution of Sum of Squares of k -independent normal variables

If X_i are K independent std. normal

$$Z = \sum_{i=1}^k X_i^2 \quad (\text{Sum of Squares})$$

↳ Distributed according to $\chi^2(k)$

χ^2 P.d.f

$$P(x; k) = \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2}$$

Mean: $E[\infty] = k$

$$\begin{aligned} E[\infty] &= \int_0^\infty x \frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2} dx \\ &= \frac{1}{2^{k/2} \Gamma(k/2)} \int_0^\infty x^{k/2} e^{-x/2} dx \quad U = \frac{x}{2} \quad dx = 2du \\ &= \frac{2^{k/2+1}}{2^{k/2} \Gamma(k/2)} \int_0^\infty u^{k/2} e^{-u} du \quad \text{gamma function: } \Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du \quad a = k/2 + 1 \\ &= \frac{2}{\Gamma(k/2)} \cdot \Gamma\left(\frac{k}{2} + 1\right) \quad \text{gamma identity: } \Gamma(a+1) = a\Gamma(a) \\ &= \frac{2}{\Gamma(k/2)} \cdot \frac{k}{2} \Gamma(k/2) \\ &= k \end{aligned}$$

Variance

$$V[\infty] = 2k$$

$$\sigma = \sqrt{2k}$$

$$E[\infty^2] - E[\infty]^2 = E[\infty^2] - k^2$$

$$\begin{aligned} E[\infty^2] &= \int_0^\infty x^2 \left(\frac{1}{2^{k/2} \Gamma(k/2)} x^{k/2-1} e^{-x/2} \right) dx \\ &= \frac{1}{2^{k/2} \Gamma(k/2)} \int_0^\infty x^{k/2+1} e^{-x/2} dx \quad U = \frac{x}{2} \quad dx = 2du \\ &= \frac{2}{2^{k/2} \Gamma(k/2)} \int_0^\infty (2u)^{k/2+1} e^{-u} du \\ &= \frac{4}{\Gamma(k/2)} \int_0^\infty u^{k/2+1} e^{-u} du \quad \text{gamma function: } \Gamma(a) = \int_0^\infty u^{a-1} e^{-u} du \quad a = k/2 + 2 \\ &= \frac{4}{\Gamma(k/2)} \Gamma(k/2 + 2) \quad \text{gamma identity: } \Gamma(a+1) = a\Gamma(a) \quad \Gamma(k/2 + 2) = (k/2 + 1) \Gamma(k/2 + 1) \\ &= \frac{4}{\Gamma(k/2)} (k/2 + 1) (k/2) \Gamma(k/2) \\ &= 4 (k/2 + 1) (k/2) = k^2 + 2k \end{aligned}$$

$$V[\infty] = (k^2 + 2k) - k^2 = k^2 + 2k - k^2 = 2k$$

C.d.f.: allows us to read off fraction contained within given Standard deviation of multivariate normal distribution

Note: Fraction contained within given Z Score for different degrees of freedom (k)

DOF (k)	Within 1σ
k = 1	0.683
k = 2	0.393
k = 3	0.199

Convolutions

Imagine 2 independent r.v.s X, Y then $Z = X + Y$

$$Z = X + Y$$

$$f(x, y) = g(x)h(y)$$

C.d.f. $F(z)$

$$\begin{aligned} F(z) &= P(X+Y \leq z) = \int_{-\infty}^{\infty} g(x) dx \int_{-\infty}^{z-x} h(y) dy \\ &= \int_{-\infty}^{\infty} h(y) dy \int_{-\infty}^{z-y} g(x) dx \end{aligned}$$

P.d.f. $f(z)$

$$\begin{aligned} f(z) &= \frac{dF(z)}{dz} = \int_{-\infty}^{\infty} g(x)h(z-x)dx \\ &= \int_{-\infty}^{\infty} h(y)g(z-y)dy \end{aligned} \quad \left. \right\} \text{known as Convolution: } f(z) = g(x) \otimes h(y)$$

If $Z = XY$ (mellin Convolution)

$$\begin{aligned} f(z) &= \int_{-\infty}^{\infty} \frac{g(x)h(z/x)}{|x|} dx \\ &= \int_{-\infty}^{\infty} \frac{g(z/y)h(y)}{|y|} dy \end{aligned}$$

Generating Samples from distributions and inverse C.d.f.

Simulation knowing Inverse C.d.f.

C.d.f.: Random Variable \rightarrow Cumulative probability $[0, 1]$ $F(x) = \int_{-\infty}^x f(x)dx$

Inverse C.d.f.: $[0, 1] \rightarrow$ Random Variable x $F^{-1}(p) = x$

↳ Commonly 'Percentage point function' (p.p.f.)

↳ $F^{-1}(p)$ or $\Phi^{-1}(p)$

1) Generate Uniform random numbers

2) Sub into p.p.f to transform into relevant dist.

Accept Reject Method (Brute force)

1) Compute mass value p.d.f. (f_{mass})

↳ Not necessarily normalised

2) Generate random uniform normal in x

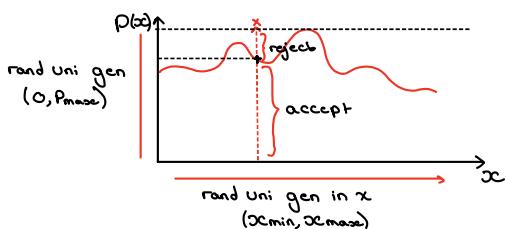
↳ along random variable x

3) Generate random uniform normal in y

↳ between

4) If $y < p(x)$ then accept

5) If $y > p(x)$ then reject.



Convergence

Convergence in Distribution

For $N \rightarrow \infty$, distribution of sample tends to $p(x)$
 ↳ Sample Size ↳ probability density (true)

For (x_1, x_2, \dots, x_n) with empirical c.d.f.s (F_1, F_2, \dots, F_n)

$$\lim_{n \rightarrow \infty} F_n(x) = F(x) \rightarrow \text{Requirement for 'Convergence in distribution'}$$

Convergence in Probability

Implies Convergence in distribution (stronger requirement)

↳ places a requirement on distance x_n from x (bounding limit)

'Probability for 'outliers' becomes smaller as Sample Size Increases

For (x_1, x_2, \dots, x_n) converges in probability to the random variable x

↳ If for all $\epsilon > 0$:

$$\lim_{n \rightarrow \infty} P(|x_n - x| > \epsilon) = 0$$

Law of Large Numbers

Often distribution not a 'priori' known — obtained using experiment

Frequentist Probability

Expectation of frequency tends to probability

h - Sample frequency/mean

$$h = \frac{1}{n} \sum_{i=1}^n x_{ci}$$

↳ 'outcome of i^{th} trial' → If 'match'
 ↳ no trials

$\sum x_{ci}$ — no. of times event occurred

in n samples

$\frac{1}{n} \sum x_{ci}$ — frequency/proportion of trials in which event occurred

$$E[h] = \hat{h} = \frac{1}{n} \sum E[x_{ci}] = \frac{1}{n} np = p$$

Although $E[x_{ci}] = np$ ↳ Whole distribution
 $E[x_{ci}] = p$ ↳ Outcome of i^{th} trial

Law of Large No.: As $n \rightarrow \infty$ then $h \rightarrow p$

$$E[h] = p \quad (\text{consistent}).$$

Central Limit Theorem

For any x_i generated from any distribution (mean, N_i and variance, σ_i^2)

Sum of random variables:

$$S = \sum_i x_i \quad \text{has mean: } N = \sum_i N_i$$

$$\text{variance } \sigma^2 = \sum_i \sigma_i^2$$

Central Limit theorem: as $N \rightarrow \infty$, distribution of $S \rightarrow$ Normal/Gaussian $G(N, \sigma^2)$

Any distribution with defined mean and variance with large $N \rightarrow$ tends to Normal Dist

Proof of CLT

Assume N independent random variables x_i from any distribution with mean N_i and variance σ_i^2

Standardise them:

$$Z_i = \frac{x_i - N_i}{\sigma_i}$$

Wish to find probability distribution: $Y = \frac{\sum_i Z_i}{\sqrt{N}}$ as $N \rightarrow \infty$

Moments

$$E[Z_i] = \frac{E[x_i] - E[N_i]}{\sigma_i} = 0$$

$$\begin{aligned} E[Z_i^2] &= E\left[\frac{(x_i - N_i)^2}{\sigma_i^2}\right] = \frac{1}{\sigma_i^2} E[x_i^2 - 2x_i N_i + N_i^2] \\ &= \frac{1}{\sigma_i^2} \left\{ E[x_i^2] - 2N_i E[x_i] + N_i^2 \right\} \\ &= \frac{1}{\sigma_i^2} \left\{ \frac{E[x_i^2]}{\sigma_i^2 + N_i^2} - N_i^2 \right\} = \frac{1}{\sigma_i^2} \left\{ \sigma_i^2 + N_i^2 - N_i^2 \right\} \\ &= 1 \end{aligned}$$

Characteristic Functions

$$\psi_y = \prod_i \psi_{z_i}(t/\sqrt{N})$$

$$\begin{aligned} \text{Studying } \psi_{z_i} &= E[e^{it z_i}] = E\left[\sum_{n=0}^{\infty} \frac{(it)^n}{n!} z_i^n\right] \xrightarrow{\text{Expansion around moments}} \\ &= \sum_{n=0}^{\infty} \frac{(it)^n}{n!} E[z_i^n] = 1 + it E[z_i] + \frac{(it)^2}{2} E[z_i^2] + \frac{(it)^3}{3!} E[z_i^3] + \dots \\ &= 1 - \frac{t^2}{2} + \frac{(it)^3}{3!} E[z_i^3] + \dots \end{aligned}$$

$$\psi_y(t) = \prod_{i=0}^N \psi_{z_i}\left(\frac{t}{\sqrt{N}}\right) = \left(1 - \frac{t^2}{2N} + \frac{(it)^3}{3! N^{3/2}} E[z_i^3] + \dots\right)^N$$

Change: $\frac{N}{i=0}$ and $\frac{1}{\sqrt{N}}$

$$\begin{aligned} \ln \psi_y(t) &= N \ln \left(1 - \frac{t^2}{2N} + \frac{(it)^3}{3! N^{3/2}} E[z_i^3] + \dots \right) \xrightarrow{\text{Vanishes as } N \rightarrow \infty} \\ &= N \ln \left(1 - \frac{t^2}{2N} \right) = N \left(-\frac{t^2}{2N} - \frac{t^4}{4N^2} + \dots \right) = -\frac{t^2}{2} \end{aligned}$$

Recalling Properties of Characteristic Function

$$\psi_x(t) = E[e^{itx}]$$

$$\psi_y(t) = \psi_{\frac{\sum_i Z_i}{\sqrt{N}}}(t) = E[e^{it \frac{\sum_i Z_i}{\sqrt{N}}}]$$

$$= \psi_{Z_i}\left(\frac{t}{\sqrt{N}}\right)$$

$$= \prod_i \psi_{z_i}\left(\frac{t}{\sqrt{N}}\right)$$

↳ Char func of sum of r.v's
= Product of indiv characteristic functions

$$\ln Y_y(t) = -\frac{t^2}{2}$$

$$Y_y(t) = e^{-t^2/2}$$

→ C.f for std normal distribution

Propagation of Errors

Often give Central values and Uncertainties
 ↳ Convention: Uncertainty at 1 std

$x \pm \sigma_x$ → Require σ_y on $y(x)$

Linear Function

Say $Z = ax + b$

$$\begin{aligned} V(Z) &= E[Z^2] - E[Z]^2 \\ &= E[(ax+b)^2] - E[ax+b]^2 \\ &= E[a^2 x^2 + abx + b^2] - (aE[x] + b)^2 \\ &= a^2 E[x^2] + 2ab E[x] + b^2 - (a^2 E[x]^2 + 2ab E[x] + b^2) \\ &= a^2 \{E[x^2] - E[x]^2\} = a^2 V(x) \end{aligned}$$

$$V(Z) = a^2 V(x)$$

$\sigma_z = a \sigma_x$ → Shift of b : does nothing to Spread

Scale of a : Increases Spread

Generic Functions $y(x)$ - Using Taylor Expansion

Taylor Expansion:

$$y(x) = y(x_0) + (x - x_0) \left. \frac{dy}{dx} \right|_{x=x_0} + \mathcal{O}(x^2) + \dots$$

$$\begin{aligned} V(y(x)) &= E[y(x)^2] - E[y(x)]^2 \\ &= E \left[y(x_0)^2 + (x - x_0) \left. \frac{dy}{dx} \right|_{x=x_0}^2 + 2 y(x_0) (x - x_0) \left. \frac{dy}{dx} \right|_{x=x_0} \right] \\ &\quad - E[y(x_0)]^2 - E \left[(x - x_0) \left. \frac{dy}{dx} \right|_{x=x_0} \right]^2 \\ &= y(x_0)^2 + \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right)^2 E[x^2 - 2x x_0 + x_0^2] - y(x_0)^2 - \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right)^2 (E[x]^2 + E[x_0]^2 - 2E[x_0]E[x]) \\ &= \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right) \{E[x^2] - E[2x x_0] + E[x_0^2] - E[x]^2 - E[x_0]^2 - 2E[x_0]E[x]\} \\ &= \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right) \{E[x^2] - 2x_0 E[x] + x_0^2 - E[x]^2 - x_0^2 - 2x_0 E[x]\} \\ &= \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right) \{E[x^2] - E[x]^2\} = \left(\left. \frac{dy}{dx} \right|_{x=x_0} \right) V(x) \end{aligned}$$

$$\sigma_y \approx \left| \left. \frac{dy}{dx} \right| \right| \sigma_x$$

→ Only valid for small errors

Error Propagation - Two variables $y(x, y)$

$$V(y) = \left(\left. \frac{dy}{dx} \right| \right)^2 V(x) + 2 \left(\left. \frac{dy}{dx} \right| \right) \left(\left. \frac{dy}{dy} \right| \right) \text{Cov}(x, y) + \left(\left. \frac{dy}{dy} \right| \right)^2 V(y)$$

$$\sigma_y^2 = \left(\left. \frac{dy}{dx} \right| \right)^2 \sigma_x^2 + \left(\left. \frac{dy}{dy} \right| \right)^2 \sigma_y^2 + 2 \left(\left. \frac{dy}{dx} \right| \right) \left(\left. \frac{dy}{dy} \right| \right) \rho \sigma_x \sigma_y$$

Error propagation - Many variables

$$\text{Cov}(f_k, f_j) = \sum_i \sum_j \left(\frac{\partial f_k}{\partial x_i} \right) \left(\frac{\partial f_j}{\partial x_i} \right) \text{Cov}(x_i, x_j)$$

In terms of Jacobian Matrix

$$V_j = \overline{J} \vee \times \overline{J}^T$$

Jacobian Matrix

$$\begin{bmatrix} \frac{\partial f_1}{\partial x_1} & \cdots & \frac{\partial f_1}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1} & \cdots & \frac{\partial f_m}{\partial x_n} \end{bmatrix}$$

$f: \mathbb{R}^n \rightarrow \mathbb{R}^m$ where $\frac{\partial f_i}{\partial x_j}$ - partial derivative of i^{th} output function wrt j^{th} input variable

Averaging and Combining Measurements

$$N = \frac{1}{N} \sum_i N_i$$

$$\text{Var}(N) = \frac{1}{N^2} \sum v_i = \frac{\sigma^2}{N} \quad \text{Standard deviation falls like }$$

$$\sigma \rightarrow \frac{\sigma}{\sqrt{2}}$$

Combining measurements with different uncertainties

'Weighted' average according to precision

$$N = \frac{\sum_i (N_i / \sigma_i^2)}{\sum \sigma_i^2}$$

$$V(N) = \frac{1}{\sum_i (1/\sigma_i^2)}$$

Estimates

Measurement (or estimate) of a parameter based on a limited no. of observations.

Produced by 'estimator'

Note: $\hat{\theta} \neq \theta_0$

Consistency

↳ Estimator Consistent - estimate produced tends to true value as data size increases

Difference should get smaller as $N \rightarrow \infty$

$$\lim_{N \rightarrow \infty} \hat{\Theta} = \Theta_0 \rightarrow \text{Variance} \downarrow \text{as } N \uparrow$$

Bias

Deviation of the expectation of the estimate, $E[\hat{\theta}]$, from true value θ

$$b_N(\hat{\theta}) = E[\hat{\theta}] - \theta_0 = E[\hat{\theta} - \theta_0]$$

Estimator unbiased: for any Sample Size, N. $b_n(\hat{y}) = 0$

Efficiency

Estimator produces narrow spread of possible values, minimal variance

Bias - Variance Trade Off

Often trade off between Bias and Efficiency

↳ mathematical description of accuracy vs precision

Estimates of mean, variance and standard deviation

Mean:

Arithmetic mean \rightarrow Consistent and unbiased \rightarrow True mean of
of Sample as an estimate Population

$$\hat{\mu} = E[\bar{x}] = \bar{x} = \frac{1}{N} \sum_{i=1}^N x_i$$

Central Limit Theorem tells us variance of the estimate :

(ie Variance of the estimate - Not the estimate of Variance)

$$\hat{V}(\bar{x}) = \frac{\sigma^2}{N}$$

Variance of parent distribution
Size of Sample

$$\sigma_{\bar{x}} = \frac{\sigma}{\sqrt{N}}$$

Variance:

1) If we know the true mean, μ :

$$\hat{V}(\bar{x}) = \frac{1}{N} \sum_i (x_i - \mu)^2$$

$$E[\hat{V}(\bar{x})] = \frac{N E[(x - \mu)^2]}{N} = V(\bar{x})$$

unbiased and consistent estimate

2) Usually do not know μ true, thus estimate $\hat{\mu}$:

$$\hat{V} = \frac{1}{N} \sum (x_i - \bar{x})^2 = \frac{1}{N} \sum (x_i^2 - \bar{x}^2)$$

Take expectation:

$$\begin{aligned} E[\hat{V}] &= \frac{N E[x^2 - \bar{x}^2]}{N} = E[x^2] - E[\bar{x}^2] \\ &= \underbrace{\{E[x^2] - E[x]^2\}}_{V(x)} - \underbrace{\{E[\bar{x}^2] - E[\bar{x}]^2\}}_{V(\bar{x})} \\ &= V(x) - \frac{V(\bar{x})}{N} \end{aligned}$$

Variance on mean est.

$$E[\hat{V}] = \frac{N-1}{N} V(x) \neq V(x)$$

Thus biased As $N \rightarrow \frac{N-1}{N} \rightarrow 0$

Shows origin of Bessel Function

$$\hat{V} = S^2 = \frac{1}{N-1} \sum (x_i - \bar{x})^2$$

Exactly compensates for bias
Consistent and unbiased

Standard error on Variance

$$V(\hat{V}) = \frac{2\sigma^4}{N-1}$$

Standard Deviation on Variance

$$\begin{aligned} \hat{\sigma} &= \sqrt{\hat{V}} = S \\ &= \frac{\sigma}{\sqrt{2(N-1)}} \end{aligned}$$

Summary of Sample estimates

Given i.i.d Sample $X_i = \{x_1, x_2, \dots, x_N\}$

$$\text{Mean: } \hat{N} = \frac{1}{N} \sum x_i$$

$$\text{Std dev: } \hat{\sigma} = \sqrt{\frac{1}{N-1} \sum_i (x_i - \bar{x})^2}$$

Error on mean

$$\hat{\sigma}_N = \frac{\hat{\sigma}}{\sqrt{N}}$$

Error on standard

$$\hat{\sigma}_\sigma = \frac{\hat{\sigma}}{\sqrt{2(N-1)}}$$

The Likelihood function

↳ likelihood of observing the dataset we have given a particular set of parameters

For N independent observations of random variable $X = \{x_1, \dots, x_N\}$

$$\text{p.d.f: } p(X|\theta) = \prod_{i=1}^N f(x_i|\theta)$$

Replace variable X with particular set of observations x

$$L(\theta) = L(x, \theta) = p(x|\theta)$$

Maximum Likelihood Method

Estimate of θ obtained by finding value $\hat{\theta}$ which minimises $L(\theta)$

$$\left. \frac{\partial L}{\partial \theta} \right|_{\theta=\hat{\theta}} = 0$$

Log-Likelihood

$$\text{Log}(L(\theta)) = \sum_{i=1}^N \ln p(x_i|\theta)$$

$$\left. \frac{\partial \ln L(\theta)}{\partial \theta} \right|_{\theta=\hat{\theta}} = \frac{\partial}{\partial \theta} \sum_{i=1}^N \ln p(x_i|\theta) = 0$$

→ often requires numerical minimisation.

Example: Normal Dist

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

$$\ln p(x; \mu, \sigma) = -\frac{(x-\mu)^2}{2\sigma^2} + C$$

The Score

$$S(\theta) = \frac{\partial \ln L(\theta)}{\partial \theta}$$
 partial derivative of likelihood function wrt θ

Maximum Likelihood Estimation: Where $S(\theta) = 0$

Expectation value of Score

Likelihood function gives weight: expected Score function over all possible observations x given θ .

$$E[S(\theta)] = \int \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 0$$

total p.d.f is normalised $\int L(\theta) dx = \int \prod_{i=1}^N p(x_i; \theta) dx = 1$

$$\int L(\theta) dx$$

$$\underbrace{\frac{\partial}{\partial \theta} \int L(\theta) dx}_{1} = \int \frac{\partial L(\theta)}{\partial \theta} dx = \int \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx$$

$$= \frac{\partial}{\partial \theta} \cdot 1 = 0$$

$$\text{Overall: } \int \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 0 \quad \therefore E[S(\theta)] = 0 \quad \text{QED.}$$

Information (variance of score)

$$I(\theta) = V(S(\theta)) = E[S(\theta)^2] - \cancel{E[S(\theta)]^2} = E[S(\theta)^2]$$

$$= E \left[\left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 \right]$$

Using knowledge

$$E[S(\theta)] = E \left[\frac{\partial \ln L(\theta)}{\partial \theta} \right] = \int \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 0$$

$$\frac{\partial}{\partial \theta} E[S(\theta)] = \frac{\partial}{\partial \theta} \int \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx \quad \text{Def: } \frac{\partial \ln L(\theta)}{\partial \theta} \cdot L(\theta)$$

$$= \int \frac{\partial^2 \ln L(\theta)}{\partial \theta^2} L(\theta) dx + \int \frac{\partial \ln L(\theta)}{\partial \theta} \frac{\partial L(\theta)}{\partial \theta} dx$$

$$= \int \frac{\partial^2 \ln L(\theta)}{\partial \theta^2} L(\theta) dx + \int \left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 L(\theta) dx = 0$$

$$\int \frac{\partial^2 \ln L(\theta)}{\partial \theta^2} L(\theta) dx = - \int \left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 L(\theta) dx$$

$$= E \left[\frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \right] = - E \left[\left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 \right] = - I(\theta)$$

$$I(\theta) = E \left[\left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 \right] = - E \left[\frac{\partial^2 \ln L(\theta)}{\partial \theta^2} \right]$$

The Minimum Variance Bound

$$E[\hat{\theta}] = \theta = \int \hat{\theta} L(\theta) dx = \theta$$

Def: Unbiased $\hat{\theta}$

Differentiate w.r.t θ (provided integral limits do not depend on θ)

$$\int \hat{\theta} \frac{\partial L(\theta)}{\partial \theta} dx = 1$$

$$\frac{\partial L(\theta)}{\partial \theta} = \frac{\partial \ln L(\theta)}{\partial \theta} \cdot L(\theta)$$

$$\int \hat{\theta} \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 1$$

Subtract $\theta \cdot E[S(\theta)] = \int \theta \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 0$

$$\int (\hat{\theta} - \theta) \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx = 1 \quad \text{Exploit Cauchy Schwartz Inequality:}$$

$$\int u^2 dx \int v^2 dx \geq \left(\int uv dx \right)^2$$

$$\text{With } u = (\hat{\theta} - \theta) \sqrt{L(\theta)} \quad v = \frac{\partial \ln L(\theta)}{\partial \theta} \sqrt{L(\theta)}$$

$$\left(\int uv dx \right)^2 = \left(\int (\hat{\theta} - \theta) \frac{\partial \ln L(\theta)}{\partial \theta} L(\theta) dx \right)^2 \leq \int (\hat{\theta} - \theta)^2 L(\theta) dx \int \left(\frac{\partial \ln L(\theta)}{\partial \theta} \right)^2 L(\theta) dx$$

$$1 \leq \underbrace{\int (\theta - \hat{\theta})^2 L(\theta) dx}_{V(\hat{\theta})} \leq \int \underbrace{\left(\frac{\partial \ln L}{\partial \theta} \right)^2 L(\theta) dx}_{E \left[\left(\frac{\partial \ln L}{\partial \theta} \right)^2 \right]}$$

$$V(\hat{\theta}) \geq \left(E \left[\left(\frac{\partial \ln L}{\partial \theta} \right)^2 \right] \right)^{-1}$$

$$\geq - \left(E \left[\frac{\partial^2 \ln L}{\partial \theta^2} \right] \right)^{-1}$$

Shows minimum Variance bound from likelihood

Allows Calc of Variance

Can Show as $n \rightarrow \infty$:

Maximum Likelihood Estimate = minimum Variance bound

i.e. ML efficient as $n \rightarrow \infty$

Many Parameter Likelihood

Covariance of the estimates:

$$\text{Cov}(\theta_i, \theta_j) = -E \left[\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \right] = - \left(\frac{\partial^2 \ln L}{\partial \theta_i \partial \theta_j} \Big|_{\theta=\hat{\theta}} \right)^{-1}$$

For $N \rightarrow \infty$:

ML estimates: Normally distributed estimate for $\hat{\theta}$

Unbiased.

Distributed about the true value

Variance \approx Minimum Variance bound

Useful Properties of MLE

$$E \left[\frac{\partial^2 \ln L}{\partial \theta^2} \right] = \frac{\partial^2 \ln L}{\partial \theta^2} \Big|_{\theta=\hat{\theta}}$$

$$\text{At minimum } V(\hat{\theta}) = -E \left[\frac{\partial^2 \ln L}{\partial \theta^2} \right]^{-1}$$

Variance bound

Large Sample \rightarrow ML is best estimator

(unbiased, consistent, efficient)

Small N \rightarrow Variable

Profiled Log-Likelihood

Scan log likelihood in dimension of interest.

At each point; fix parameter to that value

Minimize $-\ln L$ with respect to all other parameters

Uncertainties profiled log-likelihoods

$$\text{N} \sigma \text{ error} \quad -\Delta \ln L = \frac{N^2}{2}$$

Extended Maximum Likelihood

Previously: assumed total probability density is normalised

i.e: for N i.i.d observations x_i with $x_i \sim f(x_i; \theta)$

$$\int f(x; \theta) dx = 1 \quad L(\theta) = \prod_{i=1}^N f(x_i; \theta)$$

Often: Total no. events in dataset is also a random variable

Wish to determine Overall Scale/yield as well as Shape parameters

N - Poisson distributed random variable with expectation ν

$$N \sim \text{Pois}(\nu; \nu)$$

Likelihood

$$L(\nu, \bar{\theta}) = \frac{\nu^N e^{-\nu}}{N!} \prod_{i=1}^N f(x_i; \theta)$$

$$\ln L(\nu, \bar{\theta}) = N \ln \nu - \nu - \ln(N!) + \sum_{i=1}^N \ln f(x_i; \bar{\theta})$$

Binned Maximum Likelihood

Used for very large datasets \rightarrow more computationally efficient
More numerically stable

Dataset split into N bins

Relevant p.d.f. is now Poisson with expectation ν_b for each bin.

Observed data in each bin n_b :

$$p(x_b) = \frac{\nu_b^{n_b} e^{-\nu_b}}{n_b!}$$

n_b - no. observed events

ν - Poisson expectation

Likelihood:

$$L = \prod_{b=1}^N \frac{\nu_b^{n_b} e^{-\nu_b}}{n_b!}$$

$$\nu_b = n_b \int_{b_l}^{b_h} f(x; \theta) dx$$

$$-\ln L = -\ln \prod_{b=1}^N \frac{\nu_b^{n_b} e^{-\nu_b}}{n_b!}$$

$$= -\sum n_b \ln \nu_b - \nu_b - \ln(n_b!) \quad \text{Stirling approx} \\ = -\sum [n_b \ln \nu_b - \nu_b - n_b \ln(n_b) + n_b] \quad \ln(n!) = n \ln n - n$$

$$-\ln L(\theta) = -\sum_{b=1}^N [n_b \ln (F(b_h, \theta) - F(b_l, \theta)) - (F(b_h, \theta) - F(b_l, \theta))] + C$$

b_h - upper edge bin boundary

b_l - lower edge bin boundary

Least Squares Method

$$\chi^2 = \sum_{i=1}^N (y_i^{\text{obs}} - y_i^{\text{pred}})^2 = \sum_{i=1}^N [y_i - f(x_i; \theta)]^2$$

Assumes all y_i are measured to same precision

If y_i have associated errors σ_i :

$$\chi^2 = \sum_{i=1}^N \frac{(y_i - f(x_i; \theta))^2}{\sigma_i^2}$$

Number of Degrees of Freedom

$$n_{\text{d.o.f.}} = n_{\text{bins}} - n_{\text{params}}$$

Binned Chi Squares Fts

For sufficiently large N in each bin:

y_i - sum of events in bin

$$\text{Y}_{\text{error}} = \sqrt{y_i}$$

As Poisson distributed

Relationship to log Likelihood (if Gaussian Distributed)

$y_i, y_i \pm \sigma_i$ (gaussian distributed)

$$L(\theta) = \prod_{i=1}^N \text{Gauss}(y_i - f(x_i, \theta), \sigma_i)$$

$$= \prod_{i=1}^N \frac{1}{\sigma_i \sqrt{2\pi}} \exp \left[-\frac{(y_i - f(x_i, \theta))^2}{2\sigma_i^2} \right]$$

$$\ln L(\theta) = \sum_{i=1}^n \ln \left(\frac{1}{\sigma_i \sqrt{2\pi}} \right) + \sum_{i=1}^n -\frac{1}{2} \frac{(y_i - f(\sigma_i, \theta))^2}{\sigma_i^2} = \sum_{i=1}^n -\frac{1}{2} \frac{(y_i - f(\sigma_i, \theta))^2}{\sigma_i^2} + C$$

Notice $\chi^2 = -2 \ln L + C$

$$\Delta \chi^2 = -2 \Delta \ln L$$

↳ Significant: If uncertainties are normal

$-2 \Delta \ln L$ is distributed as χ^2 with $\text{df} = \text{pars}$

which under CLT true for $N \rightarrow \infty \rightarrow$ Thus true for all $N \rightarrow \infty$

Wilkes Theorem

As $N \rightarrow \infty$:

$$-2 \Delta \ln(\theta) \xrightarrow{m \text{ parameters}} \chi^2 \xrightarrow{m \text{ degrees of freedom}}$$

Allows drawing contours on profile likelihood plots

↳ Z Score - fraction contained within χ^2 distribution

Log-Likelihood ratio (LLR) $\Delta \ln L = \ln L_1 - \ln L_0 = \ln \left(\frac{L_1(\hat{\theta})}{L_0(\hat{\theta})} \right) \rightarrow$ Test Statistic

↳ a function of estimated parameters

$$T = j(\hat{\theta})$$

Example: LLR

Note drawing Contours

Z Score	$2 \Delta \ln L$ 1D	$2 \Delta \ln L$ 2D	Frac
1	1	2.29575	0.6827
2	4	6.18007	0.9545
3	9	11.8292	0.9973

Method of Moments

↳ Often difficult to implement Maximum Likelihood and Least Squares Methods

Recall Law of Large Numbers

↳ $\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i \rightarrow \mu$ as $N \rightarrow \infty$ 'Sample estimates \rightarrow true moments for large N '

For α^{th} moment about zero:

$$\hat{\mu}_\alpha = \frac{1}{N} \sum_{i=1}^N x_i^\alpha \rightarrow \mu_\alpha \text{ as } N \rightarrow \infty$$

For α^{th} moment about the mean:

$$\hat{\mu}_\alpha = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^\alpha \rightarrow \mu_\alpha \text{ as } N \rightarrow \infty$$

Data: random variable $X_i = \{x_1, x_2, \dots, x_N\}$

Model: $f(x; \theta)$

True mean: $\mu = g(\theta) = \int_{-\infty}^{\infty} x f(x; \theta) dx \rightarrow$ if θ is the mean; $g(\theta) \rightarrow \mu = \theta$

Estimate $\hat{\mu}$: using law of large numbers

$$\bar{x} = \hat{\mu} = \frac{1}{N} \sum x_i \rightarrow \mu \text{ as } N \rightarrow \infty$$

$\rightarrow g(\theta)$ as $N \rightarrow \infty$

Estimate $\hat{\theta}$: Invert $g(\theta)$ - Solve for θ

$$\hat{\theta} = g^{-1}(\bar{x}) \text{ (Method of Moments estimate of } \theta)$$

For m params $\vec{\Theta} = \{\Theta_1, \dots, \Theta_m\}$

↳ m unknowns

Cannot Solve $N = g(\vec{\Theta})$

Exploiting Higher Powers

$m^{\text{th}} \text{ moment: } \hat{N}_m = \frac{1}{N} \sum_i x_i^m \rightarrow N_m \text{ as } N \rightarrow \infty$

have m simultaneous equations:

$$\left. \begin{array}{l} N_1 = g_1(\Theta_1, \Theta_2, \dots, \Theta_m) \\ N_2 = g_2(\vec{\Theta}) \\ \vdots \\ N_m = g_m(\vec{\Theta}) \end{array} \right\}$$

Gives m functions $g(\vec{\Theta})$

$\vec{\Theta}$ is expressed as a function $\vec{h}(\vec{N})$ of moments ie $\vec{\Theta} = h(\vec{N})$

Plug in Sample estimates of moments, \vec{N} → gives estimates $\vec{\Theta}$

Uncertainties from MLE estimates

Estimate of Covariance of moment - with Bessel Correction

$$\widehat{\text{Cov}}(\hat{N}_m, \hat{N}_p) = \frac{1}{n(n-1)} \sum_{i=1}^n (x_i^m - \bar{x}_i) (x_i^p - \bar{x}_i)$$

Propagate to Covariance on Parameters (Θ)

↳ Using functions of Θ in terms of N

$$\text{Cov}(\Theta_i, \Theta_j) = \sum_{k,l} \frac{\partial \Theta_i}{\partial N_k} \frac{\partial \Theta_j}{\partial N_l} \text{Cov}(N_k, N_l)$$

1) Example with Normal Distribution $\sim N(\mu, \sigma^2)$

1st moment $N_1 = \mu$ (mean of normal dist)

2nd moment: $N_2 = \mu^2 + \sigma^2$

↳ (St. dev of normal distribution)

Solve for $\vec{\Theta} = (\mu, \sigma)$

$$\mu = N_1$$

$$\sigma^2 = N_2 - \mu^2$$

Plugging in Sample Moments

$$\hat{N} = \hat{N}_1 = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma} = \sqrt{N_2 - \hat{N}_1^2} = \sqrt{\frac{1}{N} \sum_{i=1}^N x_i^2 - \left(\frac{1}{N} \sum_{i=1}^N x_i \right)^2} = \sqrt{\frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2}$$

2) Particle physics example

$x \sim f(x; \alpha) = N(1 + \alpha \cos x)$ → Normalise by finding N

$$\frac{1}{N} = \int_{-\pi}^{\pi} 1 + \alpha \cos x \, dx = 2\pi \rightarrow f(x; \alpha) = \frac{1 + \alpha \cos x}{2\pi}$$

$$1^{\text{st}} \text{ Moment: } N_1 = \int_{-\pi}^{\pi} \frac{x(1 + \alpha \cos x)}{2\pi} \, dx = 0 \rightarrow \therefore$$

$$2^{\text{nd}} \text{ Moment: } N_2 = \int_{-\pi}^{\pi} \frac{x^2(1 + \alpha \cos x)}{2\pi} \, dx = \frac{1}{2\pi} \int_{-\pi}^{\pi} x^2 \, dx + \frac{1}{2\pi} \int_{-\pi}^{\pi} \alpha x^2 \cos x \, dx$$

$$\begin{aligned}
 &= \frac{1}{2\pi} \left[\frac{x^3}{3} \right]_{-\pi}^{\pi} + \frac{\alpha}{2\pi} \left\{ \left[x^2 \sin(x) \right]_{-\pi}^{\pi} - \int_{-\pi}^{\pi} 2x \sin x \, dx \right\} \\
 &= \frac{1}{2\pi} \left[\frac{x^3}{3} \right]_{-\pi}^{\pi} + \frac{\alpha}{2\pi} \left\{ \cancel{\left[x^2 \sin(x) \right]_{-\pi}^{\pi}} + \left[2x \cos x \right]_{-\pi}^{\pi} - \cancel{\left[2 \sin x \right]_{-\pi}^{\pi}} \right\} \\
 &= \frac{1}{\pi} \left[\frac{\pi^3}{3} \right] + \frac{\alpha}{2\pi} [4\pi] = \frac{\pi^2}{3} - 2\alpha
 \end{aligned}$$

Solve for α in terms of N_2

$$\hat{\alpha} = \frac{1}{2} \left(\frac{\pi^2}{3} - N_2 \right) = \frac{1}{2} \left(\frac{\pi^2}{3} - \bar{x}^2 \right)$$

Goodness-of-Fit Tests

↳ determining how well parameters values and uncertainties match dataset

Test Statistic for Goodness of Fit

↳ Quantify agreement between data and model:

'How well does the data, x_0 agree with the null hypothesis, H_0 '

↳ Hypothesis-model with fitted parameters

For some Test Statistic $\sim T$ distributed according to $P(T)$

↳ Probability I got this fit or worse:

$$P_{\text{got}} = \int_{T_0}^{\infty} P(T|H_0) \, dT \quad T_0 = T \text{ evaluated at data given}$$

Chi-Squared, χ^2 Test Statistic

χ^2 value - great test statistic

↳ χ^2 with k degrees of freedom:

$$\mathbb{E}[\chi^2] = k \quad \therefore \frac{\chi^2}{\text{d.o.f.}} \approx 1$$

For a χ^2 test: $\text{d.o.f.} = \text{N}_{\text{obs}} - \text{N}_{\text{pars}}$

↳ $k = \text{N}_{\text{obs}} - \text{N}_{\text{pars}}$

The Chi-Squared Value

$$\chi^2 = \sum_{i=1}^n \frac{(O_i - E_i)^2}{E_i}$$

O_i - Observed value for i^{th} data point

E_i - Expected value for i^{th} data point

Chi-Squared Distribution

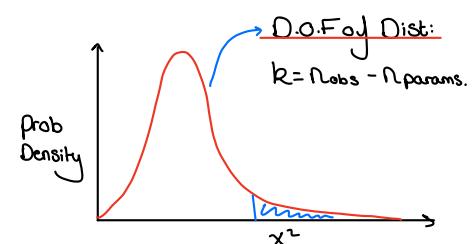
Expected outcome of the sum of squared 'independent standard normal variable'

Under null hypothesis

As DoF increases: Distribution resembles normal

χ^2 Test/Probability: $1 - F(\chi^2)$ where $F(\chi^2)$ - C.d.f. of χ^2 dist with k degrees of freedom

↳ Probability that a function that does define the data has as large or larger χ^2 value than the one you find. - Percentage change of better distribution



Small p-value: model doesn't describe data well (large discrepancies - large χ^2 value)

Larger than expected by random chance: Small chance of obtaining such a poor fit (or worse)
if model was true but data independent standard normally distributed

Large p-value: model describes data too well - overfitting - too many free parameters (very small discrepancy)

Very high chance χ^2 value for true model would be worse (larger)
ie Overfitting

Step 1: Define Hypothesis, H_0 (model with parameters)

Step 2: Calculate χ^2 value given x_i and model: χ^2

Step 3: Calculate C.d.f of χ^2 distribution with $k = n_{\text{obs}} - n_{\text{params}}$ D.O.F
With χ^2 -value, $F(\chi^2)$

Step 4: Percentage change of better distribution: $1 - F(\chi^2)$

Note: Not necessarily very powerful.

Only assesses compatibility between data and hypothesis

Not good at comparing two.

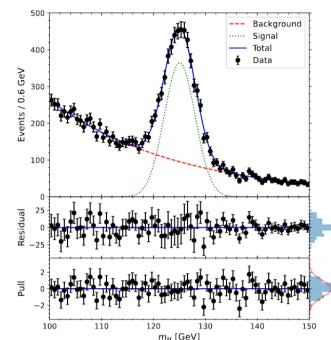
Residuals and Pull Distribution

↳ 'pull': residual normalised by the error

$$\text{pull} = \frac{n_{\text{obs}} - n_{\text{pred}}}{\sigma_{\text{obs}}}$$

Error on pull-unity

For a good fit: Expect pulls to be 'standard normally distributed'



Kolmogorov-Smirnov, KS Test

Two Sample KS Test

Compares two empirical C.d.f's of two datasets to assess whether they are likely drawn from the same underlying distribution

KS Score: Maximum deviation between two distributions

C.d.f (times $\sqrt{\text{Sample Size}}$)

$$P_{\text{KS}} = \sqrt{N} \cdot \left| \max_{\infty} (F(x_1, \infty) - F(x_2, \infty)) \right|$$

Often Empirical C.d.f:

$$F(x) = \frac{\text{No. elements in sample} \leq x}{n} = \frac{1}{n} \sum_{i=1}^n \mathbb{1}_{(-\infty, x]} (x_i)$$

One Sample KS Test

Compares dataset to specified theoretical distribution

Difference between empirical C.d.f and C.d.f of hypothesized theoretical distribution

$$P_{\text{KS}} = \sqrt{N} \cdot \left| \max_{\infty} (F_n(x_1, \infty) - F(x, \infty)) \right|$$

Hypothesised CDF
Empirical CDF

Confidence Intervals

Classical Frequentist Confidence Interval

We believe: "True value of the parameter θ_0 lies within some interval, $\theta_{\text{a}} < \theta_0 < \theta_{\text{b}}$ with some probability, B "

↳ True value of parameter, θ_0 is considered fixed

Confidence Intervals: true value will fall within interval $\beta\%$ of time

Typically quote: 10% $\beta = 68.3\%$

20% $\beta = 95.4\%$

Z Score: (Std: $1\sigma \rightarrow Z \text{ score}: 1$)

Any normal distribution can be translated to the standard normal

$$Z = \frac{X - \mu}{\sigma}$$

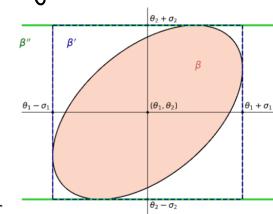
Under Coverage: too small a interval

Over Confident.

Over Coverage: too large an interval

'Conservative'

Confidence Levels 2D



	$Z = 1$	$Z = 2$	$Z = 3$
β	0.393	0.865	0.989
β' for $\rho = 0.00$	0.466	0.911	0.995
β' for $\rho = 0.20$	0.471	0.912	0.995
β' for $\rho = 0.50$	0.498	0.917	0.995
β' for $\rho = 0.80$	0.561	0.929	0.995
β' for $\rho = 0.90$	0.596	0.936	0.996
β' for $\rho = 0.95$	0.622	0.941	0.996
β' for $\rho = 1.00$	0.683	0.954	0.997
β''	0.683	0.954	0.997

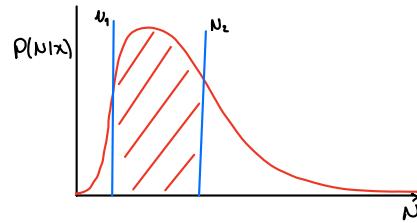
Bayesian Credible Interval

defined on posterior p.d.f

$$\beta = \int_{\Theta_1}^{\Theta_2} p(\theta | x) d\theta$$

→ Infinitely many by shifting Θ_1, Θ_2 around

Convention: quote narrowest (central) interval



To quote upper limit: Set Θ_2 to Θ_{\min}

Neyman-Pearson (Classical) Intervals

Frequentist Parameter, Θ - fixed value

Confidence Interval

$$P(N \in [N_1, N_2]) = \beta$$

Interval is built from the p.d.f of the observation, X - will vary over experiments

Difference from Bayesian

Classical $P(x|N)$ - p.d.f of observation given fixed parameter N

Bayesian $P(N|x)$ - fixed observation and varying parameter, N

For different values, N

Create Confidence belt: (for central intervals - equal probability either side)

Compute values x_1, x_2 for each N :

$$P(X < x_1 | N) = P(X > x_2 | N) = \frac{1-\beta}{2}$$

Constraints in fits

When fitting to a combination of p.d.f's

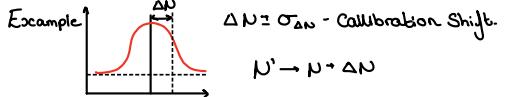
$$f_1 P_1 + f_2 P_2 + f_3 P_3$$

Constraints: $f_1, f_2, f_3 > 0$

$$f_1 + f_2 + f_3 = 1$$

Step1: $j_1 p_1 + j_2 p_2 + (1-j_1 - j_2) p_3$

Step2 (External Constraint)



Likelihood Model: $L = \prod_i p(x_i; N = \Delta N, \sigma) \times p(N + \Delta N; \Delta N \sigma_{\Delta N})$

Constraint term
p.d.f of $N + \Delta N$

Estimates and Intervals Near Physical Boundaries

Bayesian Approach Advantages:

Allows removal of physical boundaries using prior.

ie $N > 0$: $p(N) = 0 \forall N > 0$

Frequentist approach:

Frequentist typically allows parameter value and interval which goes into unphysical region

↳ possible to measure value with empty Confidence Interval → Undercoverage

Suppose physical restriction $N > 0$ is required

Confidence belt cuts off at $N > 0$

Likelihood abruptly stops at physical boundary

↳ Cannot calculate: gradient and double differential near minimum.

Values 'stack up' at boundary

Flip-Flopping (incorrect)

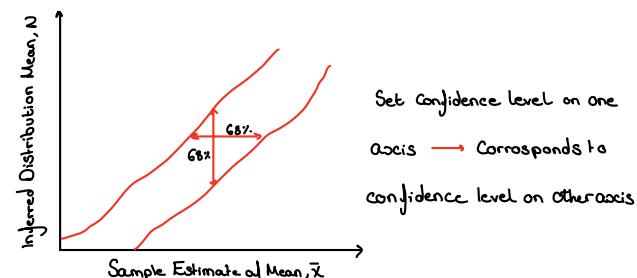
1. Quote central value and uncertainty above threshold
 2. Quote upper limits below some threshold up to the boundary
 3. Quote upper limit as boundary for any value beyond boundary
- } Leads to discontinuities in the confidence belt

Feldman-Cousins Method:

↳ Purely frequentist approach:

Creates interval directly from coverage - determined by iterations

↳ Producing simulation samples and compute were $\beta\%$ lie.



Procedure: For each x , start with Maximum Likelihood Estimate of N , \hat{N}

With physical bound applied

For some fixed (different) N :

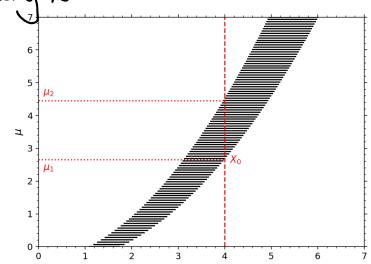
$$\text{Compute likelihood ratio: } R = \frac{p(x|N)}{p(x|\hat{N})}$$

Compares likelihood of observing x under two different assumptions

Arrange values of x from higher to lower R

Add values to interval under desired probability Content reached

} Can do these in 1 to stop over computing
R for unnecessary x



Feldman Cousins - Monte Carlo Simulation

1. Initial Fit and Maximum Likelihood

Run a maximum likelihood fit for a dataset to find best estimate \hat{N} of parameter N

Calculate likelihood value, to serve as reference ($-2 \ln L(\hat{N})$)

2. Choosing Test values of N

To find Confidence Interval - choose some 'test' values of N around \hat{N}

For each test value N_0

↳ Wish to determine confidence level that corresponds to it

Iterate over values until you find a range that meets desired confidence level

3. Define Test Statistic

$$T = -2 \Delta \ln L(N_0)$$

↳ Compares likelihood at N_0 against maximum

$$T = -2 \ln(L) = -2(\ln p(x|N) - \ln p(x|\bar{N}))$$
$$= -2(\ln L(N) - \ln L(\bar{N})) = -2 \Delta \ln(L)$$

4. Generate Pseudo Experiments

Simulate datasets generated under assumption $N=N_0$

For each toy dataset:

↳ repeat max likelihood fitting process to find best fit N_i for that toy

Calculate likelihood at both \hat{N}_0 and test value N_0

Use these to calculate Test Statistic: $T_i = -2 \Delta \ln(N_i)$ for each toy

5. Determine Confidence Level (p-value)

CL at N_0 → fraction of toy experiments where $T_i >$ observed T

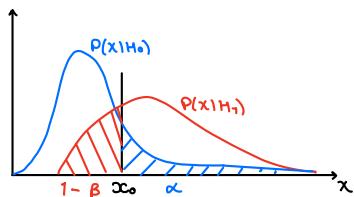
↳ indicates how often simulated data (under N_0) produces worse results than actual data.

Hypothesis Testing

Goodness of fit: test a hypothesis H_0 against all other hypotheses

Wish to: Test between null hypothesis, H_0 and alternative hypothesis, H_1

Use $P(x)$ of H_0 and H_1 to ascertain how much of each distribution falls above/below particular value.



$$\alpha = \int_{x_0}^{\infty} p(x|H_0) dx$$

$$\beta = \int_{x_0}^{\infty} p(x|H_1) dx$$

2 Hypothesis: H_0, H_1

4 outcomes:

Accept	True
H_0	H_0 (Correct)
H_1	H_0 (Incorrect - Type I error)
H_0	H_1 (Correct)
H_1	H_1 (Incorrect - Type II error)

Two types of error:

Type I: "loss" reject H_0 when true (probability: α)

Type II: "Contamination" accept H_0 when false (probability: $1 - \beta$)

α : probability of incorrectly rejecting H_0

$1 - \beta$: probability of incorrectly accepting H_0

$$\alpha = p(x > x_0 | H_0) \quad \beta = p(x > x_0 | H_1)$$

↳ Typically low α (low Type I error rate) favoured over larger $1 - \beta$ (high Type II error rate).

Power of Test: β

Probability of correctly rejecting null hypothesis H_0 (when it is false)

↳ Higher power - test more trustworthy rejection

β : Probability of correctly rejecting H_0

Assume hypothesis depends on some parameter, θ

↳ either θ - continuous

or θ_0, θ_1 - distinct two values.

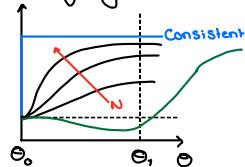
Power ρ

$$\rho(\theta_1) = \beta \quad \text{in general} \quad \rho(\theta) = \beta(\theta)$$

$$\rho(\theta_0) = \beta(\theta_0) = \alpha$$

Power Curve

↳ function of hypothesis:



Want test to be:

i) Consistent: $\lim_{N \rightarrow \infty} P(x > x_0)$ (Sample Size N increases, Prob correctly rejecting $H_0 = 1$)

ii) Unbiased: $\rho(\theta_1 | H_0) > \rho(\theta_0 | H_0)$ (Higher probability rejecting H_0 under θ_1 than θ_0)

The Neyman-Pearson Lemma

↳ Most powerful test: maximises power

Critical Region, W_α :

Region in x which is rejected $W_\alpha \in [x_0, \infty]$

W.r.t r.v $x = (x_1, x_2, \dots, x_n)$ w.p.d.f $j(x|\theta)$

$$\text{Type I error: } \alpha = \int_{x_0}^{\infty} j(\vec{x}|\theta_0) dx$$

$$\text{Type II error: } 1 - \beta = 1 - \int_{x_0}^{\infty} j(\vec{x}|\theta_1) dx$$

Find W_α (given fixed α) that maximises β (i.e. Maximal power by choosing optimal W)

$$\beta = \int_{x_0}^{\infty} \frac{j(\vec{x}|\theta_1)}{j(\vec{x}|\theta_0)} j(\vec{x}|\theta_0) dx = E_{W_\alpha} \left[\frac{j(\vec{x}|\theta_1)}{j(\vec{x}|\theta_0)} \Big|_{\theta=\theta_0} \right]$$

↳ maximal where: $\frac{j(\vec{x}|\theta_1)}{j(\vec{x}|\theta_0)}$ likelihood ratio maximal.
With constraint α .

Best most powerful test:

$$T = -2 \ln \left(\frac{L(x|H_1)}{L(x|H_0)} \right)$$

$$\begin{aligned} j(\vec{x}|\theta_1) &= j(x_1|\theta_1) j(x_2|\theta_1) \dots j(x_n|\theta_1) \\ &= T \prod_i j(x_i|\theta_1) = L(\theta_1) \end{aligned}$$

↳ T will be χ^2 distributed ($k=1$)

Limit Setting

Useful when insufficient precision to make a measurement.

↳ i.e.: Given didn't see evidence for H_1 , what is the confidence level of rejecting H_1 ?

↳ i.e. For a given confidence level, which H_1 's can I reject.

Signal Searches

Consider model of Hypotheses:

$$p(x; \bar{\theta}) = j \cdot S(x; \bar{\theta}_s) + (1-j) \cdot b(x; \bar{\theta}_b)$$

↳ Set of H dependant on j - each of which dependent on $\bar{\theta}_s$

H_0 : Null hypothesis (background only) Subst $p(x; \bar{\theta})$ where $j=0$.

H_1 : Continuous set of models with different values j .

↳ j - best fit value of j for the given data under alternative hypothesis, H_1

maximises likelihood function for observed data

'value j maximises favourability of rejecting H_0 given data.'

Changing into 'T Space' $x \rightarrow T$

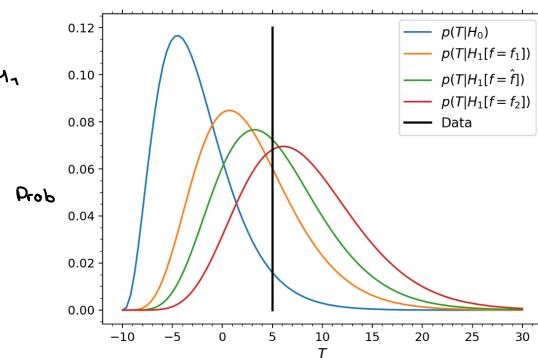
Once determined j by maximising likelihood of H_1

Calculate T (test statistic):

$$T = -2 \ln \left(\frac{L(j)}{L(j=0)} \right)$$

Also Convert $x \rightarrow T$

$$p(x|H) \rightarrow p(x|T)$$



In T space: $\alpha = \int_{T_0}^{\infty} p(T|H_0) dT$

$1 - \beta = \int_{-\infty}^{T_0} p(T|H_1) dT$

Observed Limit

For given observed data:

With confidence level β can limit values j

i.e. probability of incorrectly accepting H_0 :

$$1 - \beta = \int_{-\infty}^{T_0} p(T|H_1, j_{\text{lim}}) dT$$

'Values of j > j_{lim} excluded at $\beta \times$ Confidence level'

Expected Limits/Bands

Show the expected limits on j if H_0 is actually true.

Simulate many pseudo-datasets under H_0

↳ Calculating j_{lim} under pseudo-data set

Determine j_{lim} median - Expected Limit

Determine j_{lim} spread quantile to some confidence level - Expected Band

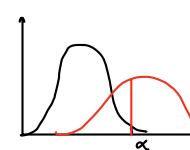
Note: One Sided vs Two Sided tests

One Sided:

eg Null hypothesis \rightarrow no signal

Alternate \rightarrow presence Signal

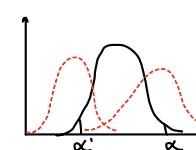
$$P_{\text{val}} = \alpha$$



Two-Sided test

$$P_{\text{val}} = \alpha + \alpha'$$

Convert Z-Score accordingly



Example:

For given observed data

↳ Calculate j

Simulate Set Pseudo Experiments H_1 ($j = j$)

↳ Calculate T for each pseudo-experiment.

↳ Distribution T under H_1 ($j = j$) \rightarrow Gives $T_{\text{exp}}(H_1, j = j)$

Simulate Set Pseudo Experiments H_0 ($j = 0$)

↳ Repeat. \rightarrow Gives $T_{\text{exp}}(H_0, j = 0)$

Compare T_0 with $T_{exp}(H_0)$ and $T_{exp}(H_1)$

↳ T_0 closer to $T_{exp}(H_0)$ - data more consistent with H_0

↳ T_0 closer to $T_{exp}(H_1)$ - data more consistent with H_1

CLs Method

Issues arise when: data inconsistent with both H_0 and H_1

or poor sensitivity to separation

CLs Method: Down weighting 'p-value' value by exclusion

P_i - probability of attaining a value as extreme or more extreme

P_b - Background only

$$P_b = \int_{T_0}^{\infty} p(T|H_0) dT \text{ or } \int_{-\infty}^{T_0} p(T|H_0) dT$$

P_{sb} - Signal and Background

$$P_{sb} = \int_{T_0}^{\infty} p(T|H_1) dT \text{ or } \int_{-\infty}^{T_0} p(T|H_1) dT$$

Adjusted p-value:

$$P_s = \frac{P_{sb}}{1 - P_b} \quad \text{Sometimes written } P_b \sim CL_s \quad P_{sb} \sim CL_{sb}$$

Reduces risk of false rejecting Signal hypothesis due to an overlap between H_0 and H_1 distributions. → Intentionally Conservative

Resampling Methods:

↳ Excellent for estimation methods and cross checking

Jackknife Resampling

$$X_i = \{x_1, x_2, \dots, x_N\}$$

$\hat{\Theta}_N$ - Biased estimate of Θ (on N observations)

Can write out bias as a Taylor expansion in N :

$$E[\hat{\Theta}_N - \Theta] = \underbrace{\frac{\alpha_1}{N}}_{\text{bias ab}} + \frac{\alpha_2}{N^2} + \frac{\alpha_3}{N^3}$$

$O(N)$

$$\text{Defn: } \hat{\Theta}' = N\hat{\Theta}_N - (N-1)\hat{\Theta}_{N-1}$$

$$E[\hat{\Theta}'] = N\left(\Theta + \frac{\alpha_1}{N} + \frac{\alpha_2}{N^2} + \dots\right) - (N-1)\left(\Theta + \frac{\alpha_1}{N-1} + \frac{\alpha_2}{(N-1)^2} + \frac{\alpha_3}{(N-1)^3} + \dots\right)$$

$$= N\Theta + \alpha_1 + \frac{\alpha_2}{N} + \frac{\alpha_3}{N^2} - (N-1)\Theta - \alpha_1 - \frac{\alpha_2}{N-1} - \frac{\alpha_3}{(N-1)^2} + \dots$$

$$= \Theta + \frac{\alpha_2}{N(N-1)} + \dots = \Theta + O\left(\frac{1}{N^2}\right)$$

↳ $\hat{\Theta}'$ is biased to $O\left(\frac{1}{N^2}\right)$

Produce: N jackknife samples omitted one point at a time

$$X_{(1)} = \{x_2, x_3, \dots, x_N\}$$

$$X_{(2)} = \{x_1, x_3, \dots, x_N\}$$

$$X_{(N)} = \{x_1, x_2, \dots, x_{N-1}\}$$

$$\hat{\theta}_1$$

$$\text{Note: } \hat{\theta}_{(i)} = f(x_1, \dots, x_{i-1}, x_{i+1}, x_N)$$

Create jackknife
estimates using:

$$\hat{\theta} = f(x)$$

$$\hat{\theta}_{(1)}$$

$$\vdots$$

$$\hat{\theta}_{(N)}$$

Average of Jackknife: $\hat{\theta}_{(1)} = \frac{1}{N} \sum_{i=1}^N \hat{\theta}_{(i)}$

from which estimate bias: $\hat{b}(\hat{\theta}) = (N-1)(\hat{\theta}_{(1)} - \hat{\theta})$

Bias Corrected $\hat{\theta}_c = \hat{\theta} - \hat{b}(\hat{\theta}) = \hat{\theta} - (N-1)(\hat{\theta}_{(1)} - \hat{\theta}) = N\hat{\theta} - (N-1)\hat{\theta}_{(1)}$

Often See Jackknife pseudo values:

$$\tilde{\theta}_{(1)} = N\hat{\theta} - (N-1)\hat{\theta}_{(1)}$$

↳ average of these gives bias Corrected est

$$\hat{\theta}_c = \frac{1}{N} \tilde{\theta}_{(1)} = \frac{1}{N} \sum (N\hat{\theta} - (N-1)\hat{\theta}_{(1)})$$

$$= N\hat{\theta} - \frac{(N-1)}{N} \sum \tilde{\theta}_{(1)} = N\hat{\theta} - (N-1)\hat{\theta}_{(1)}$$

$$= N\hat{\theta} - \frac{(N-1)}{N} \sum \hat{\theta}_{(1)} = N\hat{\theta} - (N-1)\hat{\theta}_{(1)}$$

Jackknife variance:

↳ "pseudo values"- $\tilde{\theta}_{(1)}$ are approx i.i.d

Sample variance, \hat{S}^2 , given est of variance on the estimate

Using CLT (std error on mean)

$$\hat{V}(\hat{\theta}_c) = \frac{\hat{S}^2}{N} = \frac{1}{N} \frac{1}{N-1} \sum (\hat{\theta}_{(1)} - \hat{\theta}_c)^2$$

$$= \frac{1}{N} \cdot \frac{1}{N-1} \sum (N\hat{\theta} - (N-1)\hat{\theta}_{(1)} - N\hat{\theta}_c - (N-1)\hat{\theta}_{(1)})^2$$

$$= \frac{N-1}{N} \sum_{i=1}^N \left[\hat{\theta}_{(1)} - \frac{1}{N} \sum \hat{\theta}_{(1)} \right]^2 = \frac{N-1}{N} \sum_{i=1}^N \left[\hat{\theta}_{(1)} - \hat{\theta}_{(1)} \right]^2$$

variance on jk average

w/ $\frac{N-1}{N}$ correction

Bootstrap (non-parametric) resampling

Sample with replacement:

↳ Same Sample Size as original

Randomly Choose N points - with replacement

Commonly end up with duplicate/missing events in sample.

Bootstrapping Intervals

↳ Can produce Confidence Intervals from quantiles of bootstrap distribution

Large Samples Needed for reasonable accuracy at high/low CL

↳ Typically accurate to $1/N$

Empirical CDF : \tilde{F} with Inverse: $\hat{\Phi} = \tilde{F}^{-1}$

$$[\hat{\theta}_L, \hat{\theta}_U] = [\hat{\Phi}(z_{\alpha/2}), \hat{\Phi}(1-\alpha/2)] = [\hat{\Phi}(z_{\alpha}), \hat{\Phi}(z_{1-\alpha})]$$

α - Confidence level of interval

Bias-Corrected-Accelerated (BCa) Intervals

Corrects for bias and Skewness in Bootstrap distribution.

$$\frac{\hat{\theta} - \theta}{\hat{\sigma}} = N(0, 1)$$

Standard Normal

Bias Correction:

$$[\theta_0, \theta_n] = [\tilde{\Phi}(\tilde{z}_0 + z_{\alpha}), \tilde{\Phi}(\tilde{z}_0 + z_{1-\alpha})]$$

↳ \tilde{z}_0 is a bias parameter - estimated from bootstrapped Sample

$$\rightarrow \frac{\hat{\theta} - \theta}{\hat{\sigma}} = N(-z_0, 1)$$

Additional acceleration

$$[\theta_0, \theta_n] = \left[\tilde{\Phi} \left(\tilde{z}_0 + \frac{\tilde{z}_0 + z_{\alpha}}{1 - \hat{\alpha}(\tilde{z}_0 + z_{\alpha})} \right), \tilde{\Phi} \left(\tilde{z}_0 + \frac{\tilde{z}_0 + z_{1-\alpha}}{1 - \hat{\alpha}(\tilde{z}_0 + z_{1-\alpha})} \right) \right]$$

↳ $\hat{\alpha}$ is accelerated using a jackknife of Sample

$$\rightarrow \frac{\hat{\theta} - \theta}{\hat{\sigma}} = N(-z_0(1 + \alpha\hat{\theta}), (1 + \alpha\hat{\theta})^2)$$

Parametric Bootstrap

Sampling from the p.d.f of the model directly (rather than via replacement)

'Bootstrapping based on model itself'.

Used for estimation - and check estimation procedures.

'Pull' of parameters

↳ Should find over an ensemble:

$$\frac{\hat{\theta} - \theta}{\hat{\sigma}} \sim N(0, 1)$$

→ often plot this distribution
called 'pull' distribution

Aidan Plot
(condensed)
version

If instead find:

$$N(z_0, \alpha) \rightarrow \text{biased by } z_0$$

Under/overcover by α or $1/\alpha$

Computing Systematic

Bootstrapping procedure:

Change an assumption in procedure and/or modelling

Run parametric bootstrap:

↳ generate under change of assumption
fit back with nominal procedure

Obtain Covariance matrix using Covariance estimate for BS sample.

Nuisance Parameters

Statistics in Machine learning

Measurement error (observation error)

↳ difference between measured (observed) value and true value

$$\hat{\theta} - \theta$$

Different from bias:

$$b(\hat{\theta}) = E[\hat{\theta} - \theta]$$

Two potential Sources:

↳ Statistical Fluctuations (random) - not problematic

Systematic Biases - problematic

Forward Modelling

↳ used in Supervised machine learning algorithms.

Predictions made based on input features

↳ wish to minimise difference between prediction and some reference target.

Training Sample (size N):

$$\vec{x}_i = \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_N\}$$

Target

$$y_i = \{y_1, y_2, \dots, y_N\}$$

Aim: build function:

$$y_i = f(\vec{x}_i, \vec{\theta}) \quad \vec{\theta}: \text{weights, } w_j \text{ and biases, } b_j$$

Done by minimising:

$$\chi^2 = \sum_i^N (y_i - \hat{y}_i)^2 = \sum_i^N (f(\vec{x}_i) - y_i)^2$$

Overtraining avoided by:

Statistically independent adjoint samples

Optimisation:

Given a function: $f(\vec{x})$

↳ find values \vec{x}_0 which either minimise/maximise $f(\vec{x})$

Optimisation Algorithms

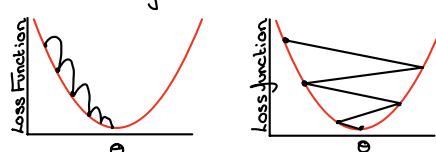
↳ Scipy.optimize

lmminuit - migrad algorithm

Gradient Descent

free parameter: learning rate

↳ Step Size of descent.



Local minima/Saddle points

↳ Training gets stuck over epochs

Parameters in optimisations

Scaling parameters

↳ optimisers are designed for / like parameters $\Theta(1)$

Scale parameters accordingly

Parameter limits

↳ help limit region in which optimiser is to search.

Regularisation

↳ terms added to cost function in order to steer or change the optimisation outcome.

Help learn trends not statistical fluctuations

Loss function becomes:

$$L = \sum_i^N V(f(x_i), y_i) + \lambda R(f)$$

V - empirical loss function

R - regularisation term

λ - input hyperparameter

- tune importance of regularisation term

Non-Parametric Density Estimation

Histograms as densities

↳ Approximate density based on independent 'control' sample

↳ Created by: Adjoint dataset

MC Simulation

For a binned fit: provides a template

Unbinned fit: approximation

↳ Any x value in bin will evaluate to same fit.

Note: Require large independent sample.

↳ avoid/smooth statistical fluctuations.

Kernel Density Estimation

↳ provides smooth continuous distribution (histogram step func \rightarrow Smoothing kernel).

For dataset: $X = \{x_1, x_2, \dots, x_N\}$

Approximate the shape:
$$f(x) = \frac{1}{N} \sum_i^N K\left(\frac{x-x_i}{h}\right)$$

Choice of bandwidth parameter:

Larger Bandwidth

↳ Smoother distribution

More bias

Smaller Bandwidth

↳ More fluctuating distribution

Greater variance.

Decision:

Low Dimensions: by eye

Higher Dimensions: Scott's, Silverman Algorithms

Types of Kernel Shape

↳ Less important than band width.

Common options: Gaussian

Exponential

Top hat

Linear.

KDEs at boundaries

All kernels typically fall off at end.

Often deal with truncated distributions

↳ Mirror distribution to other side.

Expectation Maximisation

→ Set of observed random variables: \vec{x}

Set of hidden random variables: \vec{y}

Normal Likelihood:

$$L(\vec{\theta} | \vec{x}, \vec{y}) = \prod_i p(x_i, y_i | \vec{\theta})$$

Issue: don't observe \vec{y} , so unsure of $P(x, y)$

Question: Is it possible to estimate, $\vec{\theta}$ if we don't observe \vec{y} ?

Marginalise out \vec{y} :

$$L(\vec{\theta} | \vec{x}) = \prod_i p(x_i | \vec{\theta}) = \prod_i \int p_i(x_i, y | \vec{\theta}) dy$$

$$= \prod_i \int p(x_i | y, \vec{\theta}) p(y | \vec{\theta}) dy \quad \begin{matrix} \text{→ Maximise to give } \vec{\theta} \text{ estimate} \\ \text{Still unknown} \end{matrix}$$

Using Bayes Theorem:

$$p(y | x, \vec{\theta}) = \frac{p(x | y, \vec{\theta}) p(y | \vec{\theta})}{p(x | \vec{\theta})}$$

$$\rightarrow p(x | \vec{\theta}) = \frac{p(x | y, \vec{\theta}) p(y | \vec{\theta})}{p(y | x, \vec{\theta})} = \frac{p(x | y, \vec{\theta}) p(y, \vec{\theta})}{q_y(y)} \cdot \frac{q_y(y)}{p(y | x, \vec{\theta})}$$

$$\rightarrow \ln(p(x | \vec{\theta})) = \ln\left(\frac{p(x | y, \vec{\theta}) p(y, \vec{\theta})}{q_y(y)}\right) + \ln\left(\frac{q_y(y)}{p(y | x, \vec{\theta})}\right)$$

Expectation over $q_y(y)$:

$$E[\ln p(x | \vec{\theta})] = \int q_y(y) \ln\left(\frac{p(x | y, \vec{\theta}) p(y, \vec{\theta})}{q_y(y)}\right) dy + \int q_y(y) \ln\left(\frac{q_y(y)}{p(y | x, \vec{\theta})}\right) dy$$

G(\vec{\theta}) - Lower Bound Functional Kullback-Liebler Divergence

