

1 Statistical models

$$E, \{P_\theta\}_{\theta \in \Theta}$$

E is a sample space for X i.e. a set that contains all possible outcomes of X
 $\{P_\theta\}_{\theta \in \Theta}$ is a family of probability distributions on E .

Θ is a parameter set, i.e. a set consisting of some possible values of θ .

θ is the true parameter and unknown.

In a parametric model we assume that

$\Theta \subset \mathbb{R}^d$, for some $d \geq 1$.

1.1 Identifiability

$$\theta \neq \theta' \Rightarrow P_\theta \neq P_{\theta'}$$

$$P_\theta = P_{\theta'} \Rightarrow \theta = \theta'$$

A Model is well specified if:

$$\exists \theta \text{ s.t. } P = P_\theta$$

2 Estimators

A statistic is any measurable function calculated with the data $(X_n, \max(X_i), \dots)$.

An estimator $\hat{\theta}_n$ of θ is any statistic which does not depend on θ .

Estimators are random variables if they depend on the data (= realizations of random variables).

An estimator $\hat{\theta}_n$ is **weakly consistent**

if: $\lim_{n \rightarrow \infty} \hat{\theta}_n = \theta$ or $\hat{\theta}_n \xrightarrow{P} \mathbb{E}[g(X)]$.

If the convergence is almost surely it is **strongly consistent**.

Asymptotic normality of an estimator:

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{(d)} N(0, \sigma^2)$$

σ^2 is called the **Asymptotic Variance** of the estimator $\hat{\theta}_n$. In the case of the sample mean it is the same variance as as the single X_i .

If the estimator is a function of the sample mean the **Delta Method** is needed to compute the asymptotic variance. **Asymptotic Variance** \neq Variance of an estimator.

Bias of an estimator:

$$\text{Bias}(\hat{\theta}_n) = \mathbb{E}[\hat{\theta}_n] - \theta$$

Quadratic risk of an estimator

$$\begin{aligned} R(\hat{\theta}_n) &= \mathbb{E}[(\hat{\theta}_n - \theta)^2] \\ &= \text{Bias}^2 + \text{Variance} \end{aligned}$$

3 Slutsky theorem and CMT

3.1 CMT

Let g be continuous on a set B such that $\mathbb{P}(X \in B) = 1$ then

$$\begin{aligned} X_n &\xrightarrow{P} X \Rightarrow g(X_n) \xrightarrow{P} g(X) \\ X_n &\xrightarrow{as} X \Rightarrow g(X_n) \xrightarrow{as} g(X) \\ X_n &\xrightarrow{d} X \Rightarrow g(X_n) \xrightarrow{d} g(X) \end{aligned}$$

3.2 Slutsky's Theorem:

Let c is constant and X_n and Y_n are

seq. of random elements. If $X_n \xrightarrow{d} X$ and $Y_n \xrightarrow{P} c$ then

$$\begin{aligned} (1) X_n + Y_n &\xrightarrow{d} X + c \\ (2) X_n \cdot Y_n &\xrightarrow{d} X \cdot c \\ (3) \frac{X_n}{Y_n} &\xrightarrow{d} \frac{X}{c} \end{aligned}$$

4 Distances between distributions
4.1 Total variation distance
The total variation distance TV between the probability measures P and Q with a sample space E is defined as: $TV(P, Q) = \max_{A \subseteq E} |P(A) - Q(A)|$. Calculation with f and g :

$$TV(P, Q) =$$

$$\left\{ \begin{aligned} &\frac{1}{2} \sum_{x \in E} |f(x) - g(x)|, \text{ discr} \\ &\frac{1}{2} \int_{x \in E} |f(x) - g(x)| dx, \text{ cont} \end{aligned} \right.$$

Symmetry: $TV(P, Q) = TV(Q, P)$

Positive: $TV(P, Q) \geq 0$

Definite: $TV(P, Q) = 0 \iff P = Q$

Triangle inequality: $TV(P, V) \leq$

$$TV(P, Q) + TV(Q, V)$$

If the support of P and Q is disjoint:

$$TV(P, V) = 1$$

TV between continuous and discrete r.v.:

$$TV(P, V) = 1$$

4.2 KL divergence

The KL divergence (aka relative entropy) KL between two probability measures P and Q with the common sample space E and pmf/pdf functions f and g is defined as:

$$KL(P, Q) =$$

$$\left\{ \begin{aligned} &\sum_{x \in E} p(x) \ln \left(\frac{p(x)}{q(x)} \right), \text{ discr} \\ &\int_{x \in E} p(x) \ln \left(\frac{p(x)}{q(x)} \right) dx, \text{ cont} \end{aligned} \right.$$

The KL divergence is not a distance measure! Always sum over the support of P !

Asymmetric in general:

$$KL(P, Q) \neq KL(Q, P)$$

Nonnegative: $KL(P, Q) \geq 0$

Definite: if $P = Q$ then $KL(P, Q) = 0$

Does not satisfy triangle inequality in general: $KL(P, V) \not\leq KL(P, Q) + KL(Q, V)$

Estimator of KL divergence:

$$KL(P_\theta, P_\theta) = \mathbb{E}_{P_\theta} \left[\ln \left(\frac{p_\theta(X)}{p_\theta(X)} \right) \right]$$

$$\widehat{KL}(P_\theta, P_\theta) = \text{const} - \frac{1}{n} \sum_{i=1}^n \log(p_\theta(X_i))$$

5 LLN and CLT

Let $X_1, \dots, X_n \stackrel{iid}{\sim} P_\mu$, where $E(X_i) = \mu$ and $\text{Var}(X_i) = \sigma^2$ for all $i = 1, 2, \dots, n$ and $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$.

Law of large numbers:

$$\bar{X}_n \xrightarrow{P, a.s.} \mu$$

$$\frac{1}{n} \sum_{i=1}^n g(X_i) \xrightarrow{P, a.s.} \mathbb{E}[g(X)]$$

Central Limit Theorem for Mean:

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sqrt{\sigma^2}} \xrightarrow{(d)} N(0, 1)$$

$$\sqrt{n}(\bar{X}_n - \mu) \xrightarrow{(d)} N(0, \sigma^2)$$

Central Limit Theorem for Sums:

$$\sum_{i=1}^n X_i \xrightarrow{(d)} N(n\mu, \sqrt{n}\sigma^2)$$

Variance of the Mean:

$$\begin{aligned} \text{Var}(\bar{X}_n) &= \left(\frac{\sigma^2}{n} \right)^2 \text{Var}(X_1 + X_2, \dots, X_n) \\ &= \frac{\sigma^2}{n} \end{aligned}$$

Expectation of the mean:

$$\begin{aligned} E(\bar{X}_n) &= \frac{1}{n} E[X_1 + X_2, \dots, X_n] \\ &= \mu. \end{aligned}$$

6 Quantiles of a Distribution
Let α in $(0, 1)$. The quantile of order $1 - \alpha$ of a random variable X is the number q_α such that:

$$\mathbb{P}(X \leq q_\alpha) = q_\alpha = 1 - \alpha$$

$$\mathbb{P}(X \geq q_\alpha) = \alpha$$

$$F_X(q_\alpha) = 1 - \alpha$$

$$F_X^{-1}(1 - \alpha) = q_\alpha$$

If the distribution is **standard normal** $X \sim N(0, 1)$:

$$\mathbb{P}(|X| > q_\alpha) = \alpha$$

$$= 2\Phi(q_{\alpha/2})$$

Use **standardization** if a gaussian has unknown mean and variance $X \sim N(\mu, \sigma^2)$ to get the quantiles by using Z-tables (standard normal tables).

$$\mathbb{P}(X \leq t) = \mathbb{P}\left(Z \leq \frac{t - \mu}{\sigma}\right)$$

$$= \Phi\left(\frac{t - \mu}{\sigma}\right)$$

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

$$q_\alpha = \frac{t - \mu}{\sigma}$$

7 Confidence intervals

Confidence Intervals follow the form:

(statistic) \pm (critical value)(estimated standard deviation of statistic)

Let $(E, \{P_\theta\}_{\theta \in \Theta})$ be a statistical model based on observations X_1, \dots, X_n and assume $\Theta \subseteq \mathbb{R}$. Let $\alpha \in (0, 1)$.

Non asymptotic confidence interval of level $1 - \alpha$ for θ :

Any random interval \mathcal{I} , depending on the sample X_1, \dots, X_n but not at θ and such that:

$\mathbb{P}_\theta[\mathcal{I} \ni \theta] \geq 1 - \alpha, \forall \theta \in \Theta$
Confidence interval of **asymptotic level** $1 - \alpha$ for θ :

Any random interval \mathcal{I} whose boundaries do not depend on θ and such that: $\lim_{n \rightarrow \infty} \mathbb{P}_\theta[\mathcal{I} \ni \theta] \geq 1 - \alpha, \forall \theta \in \Theta$

7.1 Two-sided asymptotic CI

Let $X_1, \dots, X_n \stackrel{iid}{\sim} P_\theta$. A two-sided CI is a function depending on \bar{X} giving an upper and lower bound in which the estimated parameter lies $\mathcal{I} = [l(X), u(X)]$ with a certain probability $\mathbb{P}(\theta \in \mathcal{I}) \geq 1 - q_\alpha$ and conversely $\mathbb{P}(\theta \notin \mathcal{I}) \leq \alpha$.

Since the estimator is a r.v. depending on \bar{X} it has a variance $\text{Var}(\hat{\theta}_n)$ and a mean $\mathbb{E}[\hat{\theta}_n]$. Since the CLT is valid for every distribution standardizing the distributions and massaging the expression yields an an asymptotic CI:

$$\begin{aligned} \mathcal{I} &= [\hat{\theta}_n - \frac{q_{\alpha/2} \sqrt{\text{Var}(\bar{X}_n)}}{\sqrt{n}}, \\ &\quad \hat{\theta}_n + \frac{q_{\alpha/2} \sqrt{\text{Var}(\bar{X}_n)}}{\sqrt{n}}] \end{aligned}$$

This expression depends on the real variance $\text{Var}(X_i)$ of the r.v.s, the variance has to be estimated. Three possible methods: plugin (use sample mean or empirical variance), solve (solve quadratic inequality), conservative (use the theoretical maximum of the variance).

7.2 Sample Mean and Sample Variance

Let $X_1, \dots, X_n \stackrel{iid}{\sim} P_\mu$, where $E(X_i) = \mu$ and $\text{Var}(X_i) = \sigma^2$ for all $i = 1, 2, \dots, n$

Sample Mean:

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$$

Sample Variance:

$$\begin{aligned} S_n &= \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \\ &= \frac{1}{n} \sum_{i=1}^n X_i^2 - \bar{X}_n^2 \end{aligned}$$

Unbiased estimator of sample variance:

$$\begin{aligned} \hat{S}_n &= \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \\ &= \frac{1}{n-1} S_n \end{aligned}$$

7.3 Delta Method

To find the asymptotic CI if the estimator is a function of the mean. Goal is to find an expression that converges a function of the mean using the CLT.

Let Z_n be a sequence of r.v. $\sqrt{n}(Z_n - \theta) \xrightarrow{(d)} N(0, \sigma^2)$ and let $g: \mathbb{R} \rightarrow \mathbb{R}$ be continuously differentiable at θ , then:

$$\begin{aligned} \sqrt{n}(g(Z_n) - g(\theta)) &\xrightarrow{(d)} N(0, g'(\theta)^2 \sigma^2) \\ &= N(0, g'(\theta)^2 \sigma^2) \end{aligned}$$

Example: let $X_1, \dots, X_n \exp(\lambda)$ where $\lambda > 0$. Let $\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i$ denote the sample mean. By the CLT, we know

that $\sqrt{n}(\bar{X}_n - \frac{1}{\lambda}) \xrightarrow{(d)} N(0, \sigma^2)$ for some value of σ^2 that depends on λ . If we set $g: \mathbb{R} \rightarrow \mathbb{R}$ and $x \mapsto 1/x$, then by the Delta method:

$$\begin{aligned} \sqrt{n}(g(\bar{X}_n) - g(\frac{1}{\lambda})) &\xrightarrow{(d)} N(0, g'(E[X])^2 \text{Var}(X)) \\ &\xrightarrow{(d)} N(0, g'(\frac{1}{\lambda})^2 \frac{1}{\lambda^2}) \\ &\xrightarrow{(d)} N(0, \lambda^2) \end{aligned}$$

8 Asymptotic Hypothesis tests

Two hypotheses (Θ_0 disjoint set from Θ_1): $\begin{cases} H_0: \theta \in \Theta_0 \\ H_1: \theta \in \Theta_1 \end{cases}$. Goal is to reject H_0 using a test statistic.

A test ψ has level α if $\alpha_\psi(\theta) \leq \alpha, \forall \theta \in \Theta_0$. and **asymptotic level** α if $\lim_{n \rightarrow \infty} P_\theta(\psi = 1) \leq \alpha$.

A hypothesis-test has the form

$$\psi = 1\{T_n \geq c\}$$

for some test statistic T_n and threshold $c \in \mathbb{R}$. Threshold c is usually $q_{\alpha/2}$

Rejection region:

$$R_\psi = \{T_n > c\}$$

Symmetric about zero and acceptance Region interval:

$$\psi = 1\{|T_n| - c > 0\}.$$

Power of the test:

$$\pi_\psi = \inf_{\theta \in \Theta_1} (1 - \beta_\psi(\theta))$$

Where β_ψ is the probability of making a Type2 Error and \inf is the maximum.

Two-sided test:

$$\begin{aligned} H_1: \theta \neq \Theta_0 \\ 1\{|T_n| > q_{\alpha/2}\} \end{aligned}$$

One-sided tests:

$$\begin{aligned} H_1: \theta > \Theta_0 \\ 1\{T_n < -q_\alpha\} H_1 \quad : \theta < \Theta_0 \\ 1\{T_n > q_\alpha\} \end{aligned}$$

Type1 Error:

Test rejects null hypothesis $\psi = 1$ but it is actually true $H_0 = \text{TRUE}$ also known as the level of a test.

Type2 Error:

Test does not reject null hypothesis $\psi = 0$ but alternative hypothesis is true $H_1 = \text{TRUE}$

Example: Let $X_1, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Ber}(p^*)$. Question: is $p^* = 1/2$.

$H_0: p^* = 1/2; H_1: p^* \neq 1/2$

If asymptotic level α then we need to standardize the estimated parameter $\hat{p} = \bar{X}_n$ first.

$$\begin{aligned} T_n &= \sqrt{n} \frac{[\bar{X}_n - 0.5]}{\sqrt{0.5(1-0.5)}} \\ \psi_n &= 1\{T_n > q_{\alpha/2}\} \end{aligned}$$

where $q_{\alpha/2}$ denotes the $q_{\alpha/2}$ quantile of a standard Gaussian, and α is determined by the required level of ψ . Note the absolute value in T_n for this two sided test.

Pivot:

Let T_n be a function of the random samples X_1, \dots, X_n, θ . Let $g(T_n)$ be a random variable whose distribution is the same for all θ . Then, g is called a pivotal quantity or a pivot.

Example: let X be a random variable with mean μ and variance σ^2 . Let X_1, \dots, X_n be iid samples of X . Then,

$$g_n \triangleq \frac{\bar{X}_n - \mu}{\sigma/\sqrt{n}}$$

is a pivot with $\theta = [\mu, \sigma^2]^T$ being the parameter vector (not the same set of parameters that we use to define a statistical model).

8.1 P-Value

The (asymptotic) p-value of a test ψ_α is the smallest (asymptotic) level α at which ψ_α rejects H_0 . It is random since it depends on the sample. It can also be interpreted as the probability that the test-statistic T_n is realized given the null hypothesis.

If $pvalue \leq \alpha$, H_0 is rejected by ψ_α at the (asymptotic) level α .

The smaller the p-value, the more confidently one can reject H_0 .

Left-tailed p-values:

$$\begin{aligned} pvalue &= \mathbb{P}(X \leq x|H_0) \\ &= \mathbb{P}(Z < T_n, \theta_0(\bar{X}_n)) \\ &= \Phi(T_n, \theta_0(\bar{X}_n)) \\ Z &\sim \mathcal{N}(0, 1) \end{aligned}$$

Right-tailed p-values:

$$pvalue = \mathbb{P}(X \geq x|H_0)$$

Two-sided p-values: If asymptotic, create normalized T_n using parameters from H_0 . Then use T_n to get to probabilities.

$$\begin{aligned} pvalue &= 2\min[\mathbb{P}(X \leq x|H_0), \mathbb{P}(X \geq x|H_0)] \\ \mathbb{P}(|Z| > |T_n, \theta_0(\bar{X}_n)|) &= 2(1 - \Phi(T_n)) \\ Z &\sim \mathcal{N}(0, 1) \end{aligned}$$

8.2 Comparisons of two proportions

Let $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p_x)$ and $Y_1, \dots, Y_n \stackrel{iid}{\sim} \text{Bern}(p_y)$ and be X

independent of Y . $\hat{p}_x = 1/n \sum_{i=1}^n X_i$ and $\hat{p}_y = 1/n \sum_{i=1}^n Y_i$

$$H_0: p_x = p_y; H_1: p_x \neq p_y$$

To get the asymptotic Variance use multivariate Delta-method. Consider $\hat{p}_x - \hat{p}_y = g(\hat{p}_x, \hat{p}_y); g(x, y) = x - y$, then

$$\begin{aligned} \sqrt{n}(g(\hat{p}_x, \hat{p}_y) - g(p_x - p_y)) &\xrightarrow{(d)} N(0, \nabla g(p_x - p_y)^T \Sigma \nabla g(p_x - p_y)) \\ &\Rightarrow N(0, p_x(1 - p_x) + p_y(1 - p_y)) \end{aligned}$$

9 Non-asymptotic Hypothesis tests
9.1 Chi squared

The χ_d^2 distribution with d degrees of freedom is given by the distribution of $Z_1^2 + Z_2^2 + \dots + Z_d^2$, where $Z_1, \dots, Z_d \stackrel{iid}{\sim} \mathcal{N}(0, 1)$

If $V \sim \chi_k^2$:

$$\begin{aligned} \mathbb{E}[\mathbb{E}[Z_1^2]] + \mathbb{E}[Z_2^2] + \dots + \mathbb{E}[Z_d^2] &= d \\ \text{Var}(V) &= \text{Var}(Z_1^2) + \text{Var}(Z_2^2) + \dots + \text{Var}(Z_d^2) = 2d \end{aligned}$$

Cochran's Theorem:

If $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$, then sample mean \bar{X}_n and the sample variance S_n are independent. The sum of squares of n variables follows a chi squared distribution with $(n-1)$ degrees of freedom:

$$\frac{nS_n}{\sigma^2} \sim \chi_{n-1}^2$$

If formula for unbiased sample variance is used:

$$\frac{(n-1)S_n}{\sigma^2} \sim \chi_{n-1}^2$$

9.2 Student's T Test

Non-asymptotic hypothesis test for small samples (works on large samples too), data must be gaussian.

Student's T distribution with d degrees of freedom: $t_d := \frac{Z}{\sqrt{V/n}}$ where $Z \sim \mathcal{N}(0, 1)$ and $V \sim \chi_k^2$ are independent.

Student's T test (one sample + two-sided):

Let $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$ and suppose we want to test $H_0: \mu = \mu_0 = 0$ vs. $H_1: \mu \neq 0$.

9.6 Goodness of Fit Discrete Distributions

Let X_1, \dots, X_n be iid samples from a categorical distribution. Test $H_0: p = p^0$ against $H_1: p \neq p^0$. Example: against the uniform distribution $p^0 = (1/K, \dots, 1/K)^T$.

Test statistic under H_0 :

$$T_n = n \sum_{k=1}^K \frac{(\hat{p}_k - p_k^0)^2}{p_k^0} \xrightarrow[n \rightarrow \infty]{} \chi_{K-1}^2$$

Test at level α :

$$\psi_\alpha = \# \{T_n > q_\alpha(\chi_{K-1}^2)\}$$

9.7 Kolmogorov-Smirnov test

$$T_n = \sqrt{n} \max_i \left(\max_{\delta \in [-1, 0]} | \frac{i+\delta}{n} - F^0(X_{(i)}) | \right)$$

9.8 Kolmogorov-Lilliefors test

$$T_n = \sup_{t \in \mathbb{R}} | F_n(t) - \Phi_{\mu, \sigma^2}(t) |$$

9.9 QQ plots

test Heavier tails: below > above the diagonal.

Lighter tails: above > below the diagonal.

Right-skewed: above > below > above the diagonal.

Left-skewed: below > above > below the diagonal.

10 Maximum likelihood estimation

Let $\{E, \{P_\theta\}_{\theta \in \Theta}\}$ be a statistical model associated with a sample of i.i.d. random variables X_1, X_2, \dots, X_n . Assume that there exists $\theta^* \in \Theta$ such that $X_i \sim P_{\theta^*}$.

The likelihood of the model is the product of the n samples of the pdf/pmf:

$$L_n(X_1, X_2, \dots, X_n, \theta) =$$

$$\begin{cases} \prod_{i=1}^n p_\theta(X_i) & \text{if } E \text{ is discrete} \\ \prod_{i=1}^n f_\theta(X_i) & \text{if } E \text{ is continuous} \end{cases}$$

The maximum likelihood estimator is the (unique) θ that minimizes the KL divergence of the parameter space. (The minimizer of the KL divergence is unique due to it being strictly convex in the space of distributions once it is fixed.)

$$\begin{aligned} \hat{\theta}_n^{\text{MLE}} &= \arg\min_{\theta \in \Theta} \widehat{\text{KL}}_n(P_{\theta^*}, P_\theta) \\ &= \arg\max_{\theta \in \Theta} \sum_{i=1}^n \ln p_\theta(X_i) \\ &= \arg\max_{\theta \in \Theta} \ln \left(\prod_{i=1}^n p_\theta(X_i) \right) \end{aligned}$$

Since taking derivatives of products is hard but easy for sums and $\exp()$ is very common in pdfs we usually take the log of the likelihood function before maximizing it.

$$\ell((X_1, X_2, \dots, X_n, \theta)) = \ln(L_n(X_1, X_2, \dots, X_n, \theta)) = \sum_{i=1}^n \ln(L_i(X_i, \theta))$$

Cookbook: set up the likelihood function, take log of likelihood function. Take the partial derivative of the log-likelihood function wrt. the parameter(s). Set the partial derivative(s) to zero and solve for the parameter. If an indicator function on the pdf/pmf does not depend on the parameter, it can be ignored. If it depends on the parameter it can't be ignored because there is a discontinuity in the log-likelihood function. The maximum/minimum of the X_i is then the maximum likelihood estimator.

10.1 Fisher information

The Fisher information is the covariance matrix of the gradient of the loglikelihood function. It is equal to the negative expectation of the Hessian of the loglikelihood function and captures the negative of the expected curvature of the loglikelihood function.

Let $\theta \in \Theta \subset \mathbb{R}^d$ and let $(E, \{P_\theta\}_{\theta \in \Theta})$ be a statistical model. Let $f_\theta(x)$ be the pdf of the distribution P_θ . Then, the Fisher information of the statistical model is:

$$\begin{aligned} \mathcal{I}(\theta) &= \text{Cov}(\nabla \ell(\theta)) = \mathbb{E}[\nabla \ell(\theta) \nabla \ell(\theta)^T] \\ &= \mathbb{E}[\nabla \ell(\theta) \mathbb{E}[\nabla \ell(\theta)]^T] = -\mathbb{E}[\text{H}(\ell(\theta))] \end{aligned}$$

Where $\ell(\theta) = \ln f_\theta(X)$. If $\nabla \ell(\theta) \in \mathbb{R}^d$ it is a $d \times d$ matrix. The definition when the distribution has a pmf $p_\theta(x)$ is also the same, with the expectation taken with respect to the pmf.

Let $(\mathbb{R}, \{P_\theta\}_{\theta \in \mathbb{R}})$ denote a continuous statistical model. Let $f_\theta(x)$ denote the pdf (probability density function) of the continuous distribution P_θ . Assume that $f_\theta(x)$ is twice-differentiable as a function of the parameter θ .

Formula for the calculation of Fisher Information of X :

$$\mathcal{I}(\theta) = \int_{-\infty}^{\infty} \frac{(\frac{\partial f_\theta(x)}{\partial \theta})^2}{f_\theta(x)} dx$$

Models with one parameter (ie. Bernoulli):

$$\mathcal{I}(\theta) = \text{Var}(\ell'(\theta))$$

$$\mathcal{I}(\theta) = -\mathbb{E}(\ell''(\theta))$$

Models with multiple parameters (ie. Gaussians):

$$\mathcal{I}(\theta) = -\mathbb{E}[\text{H}(\ell(\theta))]$$

Cookbook:

Better to use 2nd derivative.

- Find loglikelihood
- Take second derivative (=Hessian if multivariate)
- Message second derivative or Hessian (isolate functions of X_i to use with $-\mathbb{E}(\ell''(\theta))$ or $-\mathbb{E}[\text{H}(\ell(\theta))]$.
- Find the expectation of the functions of X_i and substitute them back into the Hessian or the second derivative. Be extra careful to substitute the right power back. $\mathbb{E}[X_i] \neq \mathbb{E}[X_i^2]$.
- Don't forget the minus sign!

10.2 Asymptotic normality of the maximum likelihood estimator

Under certain conditions the MLE is asymptotically normal and consistent. This applies even if the MLE is not the sample average.

Let the true parameter $\theta^* \in \Theta$. Necessary assumptions:

- The parameter is identifiable
- For all $\theta \in \Theta$, the support P_θ does not depend on θ (e.g. like in $\text{Unif}(0, \theta)$);
- θ^* is not on the boundary of Θ ;
- Fisher information $\mathcal{I}(\theta)$ is invertible in the neighborhood of θ^*

- A few more technical conditions

The asymptotic variance of the MLE is the inverse of the Fisher information.

$$\sqrt{n}(\hat{\theta}_n^{\text{MLE}} - \theta^*) \xrightarrow[n \rightarrow \infty]{} N_d(0, \mathcal{I}(\theta^*)^{-1})$$

11 Method of Moments

Let $X_1, \dots, X_n \stackrel{iid}{\sim} P_{\theta^*}$ associated with model $(\mathbb{E}, \{P_\theta\}_{\theta \in \Theta})$, with $\mathbb{E} \subseteq \mathbb{R}$ and $\Theta \subseteq \mathbb{R}$, for some $d \geq 1$. Population moments:

$$m_k(\theta) = \mathbb{E}_\theta[X_1^k], 1 \leq k \leq d$$

Empirical moments:

$$\widehat{m}_k(\theta) = \overline{X}_k^* = \frac{1}{n} \sum_{i=1}^n X_i^k$$

Convergence of empirical moments:

$$\widehat{m}_k \xrightarrow[n \rightarrow \infty]{P_{\theta^*}, a.s.} m_k$$

$$(\widehat{m}_1, \dots, \widehat{m}_d) \xrightarrow[n \rightarrow \infty]{P_{\theta^*}, a.s.} (m_1, \dots, m_d)$$

MOM Estimator M is a map from the parameters of a model to the moments of its distribution. This map is invertible, (ie. it results into a system of equations that can be solved for the true parameter vector θ^*). Find the moments (as many as parameters), set up system of equations, solve for parameters, use empirical moments to estimate.

$$\psi: \Theta \rightarrow \mathbb{R}^d$$

$$\theta \mapsto (m_1(\theta), m_2(\theta), \dots, m_d(\theta))$$

$M^{-1}(m_1(\theta^*), m_2(\theta^*), \dots, m_d(\theta^*))$ The MOM estimator uses the empirical moments:

$$M^{-1}\left(\frac{1}{n} \sum_{i=1}^n X_i, \frac{1}{n} \sum_{i=1}^n X_i^2, \dots, \frac{1}{n} \sum_{i=1}^n X_i^d\right)$$

Assuming M^{-1} is continuously differentiable at $M(0)$, the asymptotic variance of the MOM estimator is:

$$\sqrt{n}(\hat{\theta}_n^{\text{MOM}} - \theta) \xrightarrow[n \rightarrow \infty]{} N(0, \Gamma)$$

$$\begin{aligned} \text{where, } \Gamma(\theta) &= \left[\frac{\partial M^{-1}}{\partial \theta}(M(\theta)) \right]^T \Sigma(\theta) \left[\frac{\partial M^{-1}}{\partial \theta}(M(\theta)) \right] \\ \Gamma(\theta) &= \nabla_\theta (M^{-1})^T \Sigma_\theta (M^{-1}) \end{aligned}$$

Σ_θ is the covariance matrix of the random vector of the moments $(X_1^1, X_1^2, \dots, X_1^d)$.

12 Bayesian Statistics

Bayesian inference conceptually amounts to weighting the likelihood $L_n(\theta)$ by a prior knowledge we might have on θ . Given a statistical model we technically model our parameter θ as if it were a random variable. We therefore define the prior distribution (PDF):

$$\pi(\theta)$$

Let X_1, \dots, X_n . We note $L_n(X_1, \dots, X_n | \theta)$ the joint probability distribution of X_1, \dots, X_n conditioned on θ where $\theta \sim \pi$. This is exactly the likelihood from the frequentist approach.

The posterior distribution verifies:

$$\forall \theta \in \Theta, \pi(\theta | X_1, \dots, X_n) \propto$$

$$\pi(\theta) L_n(X_1, \dots, X_n | \theta)$$

The constant is the normalization factor to ensure the result is a proper distribution, and does not depend on θ :

$$\pi(\theta | X_1, \dots, X_n) = \frac{\pi(\theta) L_n(X_1, \dots, X_n | \theta)}{\int \pi(\theta) L_n(X_1, \dots, X_n | \theta) d\theta}$$

We can often use a prior that is not a proper probability distribution (whose integral diverges), and still get a proper posterior. For example, the improper prior $\pi(\theta) = 1$ on Θ gives the likelihood as a posterior.

12.2 Jeffreys Prior

$$\pi_J(\theta) \propto \sqrt{\det I(\theta)}$$

where $I(\theta)$ is the Fisher information. This prior is **invariant by reparameterization**, which means that if we have $\eta = \phi(\theta)$, then the same prior gives us a probability distribution for η verifying:

$$\pi_J(\eta) \propto \sqrt{\det I(\eta)}$$

The change of parameter follows the following formula:

$$\pi_J(\eta) = \det(\nabla \phi^{-1}(\eta)) \pi_J(\phi^{-1}(\eta))$$

12.3 Bayesian confidence region

Let $\alpha \in (0, 1)$. A Bayesian confidence region with level α^* is a random subset $\mathcal{R} \subset \Theta$ depending on X_1, \dots, X_n (and the prior π) such that:

$$P[\theta \in \mathcal{R} | X_1, \dots, X_n] \geq 1 - \alpha$$

Bayesian confidence region and confidence interval are **distinct notions**. The Bayesian framework can be used to estimate the true underlying parameter. In that case, it is used to build a new class of estimators, based on the posterior distribution.

12.4 Bayes estimator posterior mean:

$$\hat{\theta}_{(\pi)} = \int_{\Theta} \theta \pi(\theta | X_1, \dots, X_n) d\theta$$

Maximum a posteriori estimator (MAP):

$$\hat{\theta}_{(\pi)}^{\text{MAP}} = \arg\max_{\theta \in \Theta} \pi(\theta | X_1, \dots, X_n)$$

The MAP is equivalent to the MLE, if the prior is uniform.

13 OLS

Given two random variables X and Y , how can we predict the values of Y given X ?

Let us consider $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{iid}{\sim} \mathbb{P}$ where \mathbb{P} is an unknown joint distribution. \mathbb{P} can be described entirely by:

$$g(X) = \int f(X, y) dy$$

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

where f is the joint PDF, g the marginal density of X and h the conditional density. What we are interested in is $h(Y | X)$.

Regression function: For a partial description, we can consider instead the conditional expectation of Y given $X = x$:

$$x \mapsto f(x) = \mathbb{E}[Y | X = x] = \int y h(y | x) dy$$

We can also consider different descriptions of the distribution, like the median, quantiles or the variance.

Linear regression: trying to fit any function to $\mathbb{E}[Y | X = x]$ is a nonparametric problem; therefore, we restrict the problem to the tractable one of linear function:

$$f: x \mapsto a + bx$$

Linear regression: let X, Y be two random variables with two moments such as $\mathbb{V}[X] > 0$. The theoretical linear regression of Y on X is the line $a^* + b^*x$ where

$$(a^*, b^*) = \arg\min_{(a, b) \in \mathbb{R}^2} \mathbb{E}[(Y - a - bX)^2]$$

Which gives:

$$b^* = \frac{\text{Cov}(X, Y)}{\mathbb{V}[X]}, \quad a^* = \mathbb{E}[Y] - b^* \mathbb{E}[X]$$

Noise: we model the noise of Y around the regression line by a random variable $\varepsilon = Y - a^* - b^*X$, such as:

$$\mathbb{E}[\varepsilon] = 0, \quad \text{Cov}(X, \varepsilon) = 0$$

We have to estimate a^* and b^* from the data. We have n random pairs $(X_1, Y_1), \dots, (X_n, Y_n) \stackrel{iid}{\sim} (X, Y)$ such as:

$$Y_i = a^* + b^* X_i + \varepsilon_i$$

The **Least Squares Estimator (LSE)** of (a^*, b^*) is the minimizer of the squared sum:

$$(\hat{a}_n, \hat{b}_n) = \arg\min_{(a, b) \in \mathbb{R}^2} \sum_{i=1}^n (Y_i - a - bX_i)^2$$

The estimators are given by:

$$\hat{b}_n = \frac{\overline{XY} - \bar{X}\bar{Y}}{\overline{X^2} - \bar{X}^2}, \quad \hat{a}_n = \bar{Y} - \hat{b}_n \bar{X}$$

The **Multivariate Regression** is given by:

$$Y_i = \sum_{j=1}^p X_i^{(j)} \beta_j^* + \varepsilon_i = \underbrace{X_i^\top}_{1 \times p} \underbrace{\beta^*}_{p \times 1} + \varepsilon_i$$

We can assume that the $X_i^{(1)}$ are 1 for the intercept.

- If $\beta^* = (a^*, b^* \top)^\top$, $\beta_1^* = a^*$ is the intercept.
- the ε_i is the noise, satisfying $\text{Cov}(X_i, \varepsilon_i) = 0$

The **Multivariate Least Squares Estimator (LSE)** of β^* is the minimizer of the sum of square errors:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^p} \sum_{i=1}^n (Y_i - X_i^\top \beta)^2$$

Matrix form: we can rewrite these expressions. Let $Y = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$, and $\varepsilon = (\varepsilon_1, \dots, \varepsilon_n)^\top$. Let

$$X = \begin{pmatrix} X_1^\top \\ \vdots \\ X_n^\top \end{pmatrix} \in \mathbb{R}^{n \times p}$$

X is called the **design matrix**. The regression is given by:

$$Y = X\beta^* + \varepsilon$$

and the LSE is given by:

$$\hat{\beta} = \arg\min_{\beta \in \mathbb{R}^p} \|Y - X\beta\|_2^2$$

Let us suppose $n \geq p$ and $\text{rank}(X) = p$. If we write:

$$F(\beta) = \|Y - X\beta\|_2^2 = (Y - X\beta)^\top (Y - X\beta)$$

Then:

$$\nabla F(\beta) = 2X^\top (Y - X\beta)$$

Least squares estimator: setting $\nabla F(\beta) = 0$ gives us the expression of $\hat{\beta}$:

$$\hat{\beta} = (X^\top X)^{-1} X^\top Y$$

Geometric interpretation: $X\hat{\beta}$ is the orthogonal projection of Y onto the subspace spanned by the columns of X :

$$X\hat{\beta} = PY$$

where $P = X(X^\top X)^{-1} X^\top$ is the expression of the projector.

Statistic inference: let us suppose that:

* The design matrix X is deterministic and $\text{rank}(X) = p$. * The model is **homoscedastic**: $\varepsilon_1, \dots, \varepsilon_n$ are i.i.d. * and the same finite variance. * The noise is Gaussian: $\varepsilon \sim N_n(0, \sigma^2 I_n)$.

We therefore have:

$$Y \sim N_n(X\beta^*, \sigma^2 I_n)$$

Properties of the LSE:

$$\hat{\beta} \sim N_p(\beta^*, \sigma^2 (X^\top X)^{-1})$$

The quadratic risk of $\hat{\beta}$ is given by:

$$\mathbb{E}[\|\hat{\beta} - \beta^*\|_2^2] = \sigma^2 \text{Tr}((X^\top X)^{-1})$$

The prediction error is given by:

$$\mathbb{E}[\|Y - X\hat{\beta}\|_2^2] = \sigma^2(n - p)$$

The unbiased estimator of σ^2 is:

$$\hat{\sigma}^2 = \frac{1}{n - p} \|Y - X\hat{\beta}\|_2^2 = \frac{1}{n - p} \sum_{i=1}^n \varepsilon_i^2$$

By **Cochran's Theorem**:

$$(n - p) \frac{\hat{\sigma}^2}{\sigma^2} \sim \chi_{n-p}^2, \quad \hat{\beta} \perp \hat{\sigma}^2$$

Significance test: let us test $H_0: \beta_j = 0$ against $H_1: \beta_j \neq 0$. Let us call

$$\gamma_j = ((X^\top X)^{-1})_{jj} > 0$$

then:

$$\frac{\hat{\beta}_j - \beta_j}{\sqrt{\hat{\sigma}^2 \gamma_j}} \sim t_{n-p}$$

We can define the test statistic for our test:

$$T_n^{(j)} = \frac{\hat{\beta}_j}{\sqrt{\hat{\sigma}^2 \gamma_j}}$$

The test with non-asymptotic level α is given by:

$$\psi_\alpha^{(j)} = 1[\|T_n^{(j)}\| > q_{\alpha/2}(t_{n-p})]$$

Bonferroni's test: if we want to test the significance level of multiple tests at the same time, we cannot use the same level α for each of them. We must use a stricter test for each of them. Let us consider $S \subseteq \{1, \dots, p\}$. Let us consider

$H_0: \forall j \in S, \beta_j = 0, \quad H_1: \exists j \in S, \beta_j \neq 0$

The **Bonferroni's test** with significance level α is given by:

$$\psi_\alpha^{(S)} = \max_{j \in S} \psi_{\alpha/K}^{(j)}$$

where $K = |S|$. The rejection region therefore is the union of all rejection regions:

$$R_\alpha^{(S)} = \bigcup_{j \in S} R_{\alpha/K}^{(j)}$$

This test has nonasymptotic level at most α :

$$\mathbb{P}_{H_0}[R_\alpha^{(S)}] \leq \sum_{j \in S} \mathbb{P}_{H_0}[R_{\alpha/K}^{(j)}] = \alpha$$

This test also works for implicit testing (for example, $\beta_1 \geq \beta_2$).

14 Generalized Linear Models

We relax the assumption that μ is linear. Instead, we assume that $g \circ \mu$ is linear, for some function g :

$$g(\mu(x)) = x^\top \beta$$

The function g is assumed to be known, and is referred to as the link function. It maps the domain of the dependent variable to the entire real line.

It has to be strictly increasing, it has to be continuously differentiable and its range is all of \mathbb{R} .

14.1 The Exponential Family

A family of distribution $\{P_\theta: \theta \in \Theta\}$, where the parameter space $\Theta \subset \mathbb{R}^k$ is $-k$ dimensional, is called a k -parameter exponential family on \mathbb{R}^1 if the pmf or pdf $f_\theta: \mathbb{R}^d \rightarrow \mathbb{R}$ of P_θ can be written in the form:

$$\begin{aligned} f_\theta(y) &= h(y) \exp(\eta(\theta) \cdot T(y) - B(\theta)) \quad \text{where} \\ \eta(\theta) &= \begin{pmatrix} \eta_1(\theta) \\ \vdots \\ \eta_k(\theta) \end{pmatrix} : \mathbb{R}^k \rightarrow \mathbb{R}^k \\ T(y) &= \begin{pmatrix} T_1(y) \\ \vdots \\ T_k(y) \end{pmatrix} : \mathbb{R}^d \rightarrow \mathbb{R}^k \\ B(\theta) &: \mathbb{R}^k \rightarrow \mathbb{R} \\ h(y) &: \mathbb{R}^d \rightarrow \mathbb{R} \$$

Canonical exponential form:

$$f_0(y) = \exp\left(\underbrace{y\theta - \ln(1 + e^\theta)}_{b(\theta)} + \underbrace{0}_{c(y,\phi)}\right)$$
$$\theta = \ln\left(\frac{p}{1-p}\right)$$
$$\phi = 1$$

Binomial

Parameters p and n , discrete.
Describes the number of successes in n independent Bernoulli trials.

$$p_x(k) = \binom{n}{k} p^k (1-p)^{n-k}, k = 0, \dots, n$$

$$\mathbb{E}[X] = np$$

$$\text{Var}(X) = np(1-p)$$

Likelihood:

$$L_n(X_1, \dots, X_n, \theta) = \left(\prod_{i=1}^n \binom{n}{X_i}\right) \theta^{\sum_{i=1}^n X_i} (1-\theta)^{nK - \sum_{i=1}^n X_i}$$

Loglikelihood:

$$\ell_n(\theta) = C + \left(\sum_{i=1}^n X_i\right) \log \theta + (nK - \sum_{i=1}^n X_i) \log(1-\theta)$$

MLE:

Fisher Information:

$$I(p) = \frac{n}{p(1-p)}$$

Canonical exponential form:

$$f_p(y) = \exp\left(\underbrace{y(\ln p - \ln(1-p))}_{\theta} + \underbrace{n \ln(1-p)}_{-b(\theta)} + \underbrace{\ln\left(\prod_{i=1}^n \binom{n}{X_i}\right)}_{c(y,\phi)}\right)$$

Geometric

Number of T trials up to (and including) the first success.

$$p_T(t) = (1-p)^{t-1} \cdot p, t = 1, 2, \dots$$

$$\mathbb{E}[T] = \frac{1}{p}$$

$$\text{var}(T) = \frac{1-p}{p^2}$$

Pascal

The negative binomial or Pascal distribution is a generalization of the geometric distribution. It relates to the random experiment of repeated independent trials until observing m successes. I.e. the time of the k th arrival.

$$Y_k = T_1 + \dots T_k$$

$$T_i \sim \text{iid Geometric}(p)$$

$$\mathbb{E}[Y_k] = \frac{k}{p}$$

$$\text{Var}(Y_k) = \frac{k(1-p)}{p^2}$$

$$p_{X_k}(t) = \binom{t-1}{k-1} p^k (1-p)^{t-k}$$

$$t = k, k+1, \dots$$

Multinomial

Parameters $n > 0$ and p_1, \dots, p_r .

$$p_X(x) = \frac{n!}{x_1! \dots x_r!} p_1^{x_1} \dots p_r^{x_r}$$

$$\mathbb{E}[X_i] = n \cdot p_i$$

$$\text{Var}(X_i) = np_i(1-p_i)$$

Likelihood:

$$p_X(x) = \prod_{j=1}^r p_j^{T_j}, \text{ where } T^j = \mathbf{1}(X_i = j) \text{ is the count how often an outcome is seen in trials.}$$

Loglikelihood:

$$\ell_n = \sum_{j=2}^r T_j \ln(p_j)$$

Poisson

Parameter λ , discrete, approximates the binomial PMF when n is large, p is small, and $\lambda = np$.

$$p_X(k) = \exp(-\lambda) \frac{\lambda^k}{k!} \text{ for } k = 0, 1, \dots,$$

$$\mathbb{E}[X] = \lambda$$

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

Likelihood:

$$L_n(x_1, \dots, x_n, \lambda) = \prod_{i=1}^n \frac{\lambda^{x_i}}{x_i!} e^{-\lambda}$$

Loglikelihood:

$$\ell_n(\lambda) = -n\lambda + \log(\lambda) \left(\sum_{i=1}^n x_i\right) - \log\left(\prod_{i=1}^n x_i!\right)$$

MLE:

$$\hat{\lambda}_{MLE} = \frac{1}{n} \sum_{i=1}^n (X_i)$$

Fisher Information:

$$I(\lambda) = \frac{1}{\lambda}$$

Canonical exponential form:

$$f_0(y) = \exp\left(y\theta - \underbrace{e^\theta}_{b(\theta)} - \underbrace{\ln y!}_{c(y,\phi)}\right)$$

$$\theta = \ln \lambda$$

$$\phi = 1$$

Poisson process:
k arrivals in t slots

$$p_X(k, t) = \mathbb{P}(N_t = k) = e^{-\lambda t} \frac{(\lambda t)^k}{k!}$$

$$\mathbb{E}[N_t] = \lambda t$$

$$\text{Var}(N_t) = \lambda t$$

Memoryless property: The distance between two consecutive points of a point process on the real line will be an exponential random variable with parameter λ (or equivalently, mean $\frac{1}{\lambda}$). This implies that the points have the memoryless property: the existence of one point existing in a finite interval does not affect the probability (distribution) of other points existing

Interarrival Times for Poisson Processes

If $N(t)$ is a Poisson process with rate λ , then the interarrival times X_1, X_2, \dots are independent and

$$X_i \sim \text{Exponential}(\lambda), \text{ for } i = 1, 2, 3, \dots$$

Now that we know the distribution of the interarrival times, we can find the distribution of arrival times

$$L(X_1 \dots X_n; \lambda, \theta) = \lambda^n \exp(-\lambda \sum_{i=1}^n (X_i - a)) \mathbf{1}_{\min_{i=1, \dots, n} (X_i) \geq a}$$

Loglikelihood:

$$\ell(\lambda, a) := n \ln \lambda - \lambda \sum_{i=1}^n X_i + n \ln a$$

MLE:

$$\hat{\lambda}_{MLE} = \frac{1}{\bar{X}_n - a}$$

$$\hat{a}_{MLE} = \min_{i=1, \dots, n} (X_i)$$

Univariate Gaussians

Parameters μ and $\sigma^2 > 0$, continuous

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

$$\mathbb{E}[X] = \mu$$

$$\text{Var}(X) = \sigma^2$$

CDF of standard gaussian:

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

Likelihood:

$$L(x_1 \dots x_n; \mu, \sigma^2) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right)$$

Loglikelihood:

$$\ell_n(\mu, \sigma^2) = -n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2$$

$P(H) = p$ is tossed. If the coin lands heads up, the arrival is sent to the first process ($N_1(t)$), otherwise – to second. The coin tosses are independent of each other and are independent of $N(t)$. Then, $N_1(t)$ is a Poisson process with rate λp ; $N_2(t)$ is a Poisson process with rate $\lambda(1-p)$; $N_1(t)$ and $N_2(t)$ are independent.

Exponential

Parameter λ , continuous

$$f_X(x) = \begin{cases} \lambda \exp(-\lambda x), & \text{if } x \geq 0 \\ 0, & \text{o.w.} \end{cases}$$

$$P(X > a) = \exp(-\lambda a)$$

$$F_X(x) = \begin{cases} 1 - \exp(-\lambda x), & \text{if } x \geq 0 \\ 0, & \text{o.w.} \end{cases}$$

$$\mathbb{E}[X] = \frac{1}{\lambda}$$

$$\mathbb{E}[X^2] = \frac{2}{\lambda^2}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

Likelihood:

$$L(X_1 \dots X_n; \lambda) = \lambda^n \exp(-\lambda \sum_{i=1}^n X_i)$$

Loglikelihood:

$$\ell_n(\lambda) = n \ln \lambda - \lambda \sum_{i=1}^n (X_i)$$

MLE:

$$\hat{\lambda}_{MLE} = \frac{n}{\sum_{i=1}^n (X_i)}$$

Fisher Information:

$$I(\lambda) = \frac{1}{\lambda^2}$$

Canonical exponential form:

$$f_0(y) = \exp\left(y\theta - \underbrace{\ln(-\theta)}_{b(\theta)} + \underbrace{0}_{c(y,\phi)}\right)$$

$$\theta = -\lambda = -\frac{1}{\mu}$$

$$\phi = 1$$

Shifted Exponential

Parameters $\lambda, a \in \mathbb{R}$, continuous

$$f_X(x) = \begin{cases} \lambda \exp(-\lambda(x-a)), & x \geq a \\ 0, & x < a \end{cases}$$

$$F_X(x) = \begin{cases} 1 - \exp(-\lambda(x-a)), & \text{if } x \geq a \\ 0, & x < a \end{cases}$$

$$\mathbb{E}[X] = a + \frac{1}{\lambda}$$

$$\text{Var}(X) = \frac{1}{\lambda^2}$$

Likelihood:

$$L(X_1 \dots X_n; \lambda, \theta) = \lambda^n \exp(-\lambda \sum_{i=1}^n (X_i - a)) \mathbf{1}_{\min_{i=1, \dots, n} (X_i) \geq a}$$

Loglikelihood:

$$\ell(\lambda, a) := n \ln \lambda - \lambda \sum_{i=1}^n X_i + n \ln a$$

MLE:

$$\hat{\lambda}_{MLE} = \frac{1}{\bar{X}_n - a}$$

$$\hat{a}_{MLE} = \min_{i=1, \dots, n} (X_i)$$

Univariate Gaussians

Parameters μ and $\sigma^2 > 0$, continuous

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

$$\mathbb{E}[X] = \mu$$

$$\text{Var}(X) = \sigma^2$$

CDF of standard gaussian:

$$\Phi(z) = \int_{-\infty}^z \frac{1}{\sqrt{2\pi}} e^{-x^2/2} dx$$

Likelihood:

$$L(x_1 \dots x_n; \mu, \sigma^2) = \left(\frac{1}{\sigma\sqrt{2\pi}}\right)^n \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2\right)$$

Loglikelihood:

$$\ell_n(\mu, \sigma^2) = -n \log(\sigma\sqrt{2\pi}) - \frac{1}{2\sigma^2} \sum_{i=1}^n (X_i - \mu)^2$$

MLE:

$$\hat{\mu}_{MLE} = \bar{X}_n$$

$$\hat{\sigma}_{MLE}^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X}_n)^2$$

Fisher Information:

$$I(\mu, \sigma^2) = \begin{pmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{2\sigma^4} \end{pmatrix}$$

Canonical exponential form:

Gaussians are invariant under affine transformation:

$$aX + b \sim N(X + b, a^2 \sigma^2)$$

Sum of independent gaussians:

$$\text{Let } X \sim N(\mu_X, \sigma_X^2) \text{ and } Y \sim N(\mu_Y, \sigma_Y^2)$$

$$\text{If } Y = X + Z, \text{ then } Y \sim N(\mu_X + \mu_Y, \sigma_X + \sigma_Y)$$

$$\text{If } U = X - Y, \text{ then } U \sim N(\mu_X - \mu_Y, \sigma_X + \sigma_Y)$$

Symmetry:

$$\text{If } X \sim N(0, \sigma^2), \text{ then } -X \sim N(0, \sigma^2)$$

$$\mathbb{P}(|X| > x) = 2\mathbb{P}(X > x)$$

Standardization:

$$Z = \frac{X - \mu}{\sigma} \sim N(0, 1)$$

$$\mathbb{P}(X \leq t) = \mathbb{P}\left(Z \leq \frac{t - \mu}{\sigma}\right)$$

Higher moments:

$$\mathbb{E}[X^2] = \mu^2 + \sigma^2$$

$$\mathbb{E}[X^3] = \mu^3 + 3\mu\sigma^2$$

$$\mathbb{E}[X^4] = \mu^4 + 6\mu^2\sigma^2 + 3\sigma^4$$

Quantiles:

Uniform

Parameters a and b , continuous.

$$f_X(x) = \begin{cases} \frac{1}{b-a}, & \text{if } a < x < b \\ 0, & \text{o.w.} \end{cases}$$

$$F_X(x) = \begin{cases} 0, & \text{for } x \leq a \\ \frac{x-a}{b-a}, & x \in (a, b) \\ 1, & x \geq b \end{cases}$$

$$\mathbb{E}[X] = \frac{a+b}{2}$$

$$\text{Var}(X) = \frac{(b-a)^2}{12}$$

Likelihood:

$$L(X_1 \dots X_n; b) = \frac{1(\max_i(X_i) \leq b)}{b^n}$$

Loglikelihood:

Chi squared

The χ^2_d distribution with d degrees of freedom is given by the distribution of $Z_1^2 + Z_2^2 + \dots + Z_d^2$, where $Z_1, \dots, Z_d \stackrel{\text{iid}}{\sim} \mathcal{N}(0, 1)$

$$\text{If } V \sim \chi^2_k:$$

$$\mathbb{E}[\mathbb{E}[Z_1^2] + \mathbb{E}[Z_2^2] + \dots + \mathbb{E}[Z_d^2]] = d$$

$$\text{Var}(V) = \text{Var}(Z_1^2) + \text{Var}(Z_2^2) + \dots + \text{Var}(Z_d^2) = 2d$$

Student's T Distribution

$$T_n := \frac{Z}{\sqrt{V/n}} \text{ where } Z \sim \mathcal{N}(0, 1), \text{ and } Z \text{ and } V \text{ are independent}$$

15.1 Useful to know

15.1.1 Min of iid exponential r.v.

Let X_1, \dots, X_n be i.i.d. $\text{Exp}(\lambda)$ random variables.

Distribution of $\min_i(X_i)$

$$\mathbb{P}(\min_i(X_i) \leq t) = 1 - \mathbb{P}(\min_i(X_i) > t)$$

$$= 1 - (\mathbb{P}(X_1 > t))(\mathbb{P}(X_2 > t)) \dots (\mathbb{P}(X_n > t))$$

$$= 1 - (1 - F_X(t))^n = 1 - e^{-n\lambda t}$$

Differentiate w.r.t X to get the pdf of $\min_i(X_i)$:

$$f_{\min}(x) = (n\lambda)e^{-(n\lambda)x}$$

15.1.2 Counting Committees

Out of $2n$ people, we want to choose a committee of n people, one of whom will be its chair. In how many different ways can this be done?"

$$\binom{2n}{n} = 2n \binom{2n-1}{n-1}$$

"In a group of $2n$ people, consisting of n boys and n girls, we want to select a committee of n people. In how many ways can this be done?"

$$\binom{2n}{n} = \sum_{i=0}^n \binom{n}{i} \binom{n}{n-i}$$

"How many subsets does a set with $2n$ elements have?"

$$2^{2n} = \sum_{i=0}^{2n} \binom{2n}{i}$$

"Out of n people, we want to form a committee consisting of a chair and other members. We allow the committee size to be any integer in the range $1, 2, \dots, n$. How many choices do we have in selecting a committee-chair combination?"

$$n2^{n-1} = \sum_{i=0}^n \binom{n}{i} i$$

15.2 Finding Joint PDFs

$$f_{X,Y}(x,y) = f_X(x)f_Y(y|x)$$

16 Expectation

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} x \cdot f_X(x) dx$$

$$\mathbb{E}[g(X)] = \int_{-\infty}^{+\infty} g(x) \cdot f_X(x) dx$$

$$\mathbb{E}[X|Y=y] = \int_{-\infty}^{+\infty} x \cdot f_{X|Y}(x|y) dx$$

Integration limits only have to be over the support of the pdf. Discrete r.v. same as continuous but with sums and pmfs.

$$\mathbb{E}[X|Y=y] = \int_{-\infty}^{+\infty} x \cdot f_{X|Y}(x|y) dx$$

Total expectation theorem:

$$\mathbb{E}[X] = \int_{-\infty}^{+\infty} f_Y(y) \cdot \mathbb{E}[X|Y=y] dy$$

Law of iterated expectation:

$$\mathbb{E}[Y] = \mathbb{E}[\mathbb{E}[Y|X]]$$

Expectation of constant a :

$$\mathbb{E}[a] = a$$

Product of independent r.v.s X and Y :

$$\mathbb{E}[X \cdot Y] = \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

Product of dependent r.v.s X and Y :

$$\mathbb{E}[X \cdot Y] \neq \mathbb{E}[X] \cdot \mathbb{E}[Y]$$

$$\mathbb{E}[X \cdot Y] = \mathbb{E}[\mathbb{E}[Y \cdot X|Y]] = \mathbb{E}[Y \cdot \mathbb{E}[X|Y]]$$

Linearity of Expectation where a and c are given scalars:

$$\mathbb{E}[aX + cY] = a\mathbb{E}[X] + c\mathbb{E}[Y]$$

If Variance of X is known:

$$\mathbb{E}[X^2] = \text{Var}(X) + \mathbb{E}[X]^2$$

17 Variance

Variance is the squared distance from the mean.

$$\text{Var}(X) = \mathbb{E}[(X - \mathbb{E}(X))^2]$$

$$\text{Var}(X) = \mathbb{E}[X^2] - (\mathbb{E}[X])^2$$

Variance of a product with constant a :

Multivariate Gaussians

The distribution of X the d -dimensional Gaussian or normal distribution, is completely specified by the vector mean $\mu = \mathbb{E}[X] = (\mathbb{E}[X^{(1)}], \dots, \mathbb{E}[X^{(d)}])^T$ and the $d \times d$ covariance matrix Σ . If Σ is invertible, then the pdf of X is:

$$f_X(\mathbf{x}) = \frac{1}{\sqrt{(2\pi)^d \det(\Sigma)}} e^{-\frac{1}{2}(\mathbf{x}-\mu)^T \Sigma^{-1}(\mathbf{x}-\mu)},$$
$$\mathbf{x} \in \mathbb{R}^d$$

Where $\det(\Sigma)$ is the determinant of Σ , which is positive when Σ is invertible. If $\mu = 0$ and Σ is the identity matrix, then X is called a standard normal random vector.

If the covariant matrix Σ is diagonal, the pdf factors into pdfs of univariate Gaussians, and hence the components are independent.

The linear transform of a gaussian $X \sim N_d(\mu, \Sigma)$ with conformable matrices A and B is a gaussian:

$$AX + B \sim N_d(A\mu + B, A\Sigma A^T)$$

Multivariate CLT

Let $X_1, \dots, X_d \in \mathbb{R}^d$ be independent copies of a random vector X such that $\mathbb{E}[x] = \mu$ ($d \times 1$ vector of expectations) and $\text{Cov}(X) = \Sigma$

$$\sqrt{n}(\overline{X_n} - \mu) \xrightarrow{n \rightarrow \infty} N(0, \Sigma)$$

$$\sqrt{n}\Sigma^{-1/2}\overline{X_n} - \mu \xrightarrow{n \rightarrow \infty} N(0, I_d)$$

Where $\Sigma^{-1/2}$ is the $d \times d$ matrix such that $\Sigma^{-1/2}\Sigma^{-1/2} = \Sigma^{-1}$ and I_d is the identity matrix.

Multivariate Delta Method

21 Algebra

Absolute Value Inequalities:

$$|f(x)| < a \Rightarrow -a < f(x) < a$$
$$|f(x)| > a \Rightarrow f(x) > a \text{ or } f(x) < -a$$

22 Matrixalgebra

$$\|\mathbf{Ax}\|^2 = (\mathbf{Ax})^T(\mathbf{Ax}) = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax} = \mathbf{x}^T \mathbf{A}^T \mathbf{Ax}$$

23 Calculus

Differentiation under the integral sign

$$\frac{d}{dx} \left(\int_{a(x)}^{b(x)} f(x, t) dt \right) = f(x, b(x))b'(x) - f(x, a(x))a'(x) + \int_{a(x)}^{b(x)} f_x(x, t) dt.$$

Concavity in 1 dimension

If $g: I \rightarrow \mathbb{R}$ is twice differentiable in the interval I :

concave:

if and only if $g''(x) \leq 0$ for all $x \in I$

strictly concave:

if $g''(x) < 0$ for all $x \in I$

convex:

if and only if $g''(x) \geq 0$ for all $x \in I$

strictly convex if:

$g''(x) > 0$ for all $x \in I$

Multivariate Calculus

The Gradient ∇ of a twice differentiable function f is defined as:

$$\nabla f: \mathbb{R}^d \rightarrow \mathbb{R}^d$$
$$\theta = \begin{pmatrix} \theta_1 \\ \theta_2 \\ \vdots \\ \theta_d \end{pmatrix} \mapsto \begin{pmatrix} \frac{\partial f}{\partial \theta_1} \\ \frac{\partial f}{\partial \theta_2} \\ \vdots \\ \frac{\partial f}{\partial \theta_d} \end{pmatrix}_{\theta}$$

Hessian

The Hessian of f is a symmetric matrix of second partial derivatives of f

$$\mathbf{H}h(\theta) = \nabla^2 h(\theta) =$$

$$\begin{pmatrix} \frac{\partial^2 h}{\partial \theta_1 \partial \theta_1}(\theta) & \dots & \frac{\partial^2 h}{\partial \theta_1 \partial \theta_d}(\theta) \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 h}{\partial \theta_d \partial \theta_1}(\theta) & \dots & \frac{\partial^2 h}{\partial \theta_d \partial \theta_d}(\theta) \end{pmatrix} \in \mathbb{R}^{d \times d}$$

A symmetric (real-valued) $d \times d$ matrix \mathbf{A} is:

Positive semi-definite:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \geq 0 \quad \text{for all } \mathbf{x} \in \mathbb{R}^d.$$

Positive definite:

$$\mathbf{x}^T \mathbf{A} \mathbf{x} > 0 \text{ for all non-zero vectors } \mathbf{x} \in \mathbb{R}^d$$

Negative semi-definite (resp. negative definite):

$$\mathbf{x}^T \mathbf{A} \mathbf{x} \text{ is negative for all } \mathbf{x} \in \mathbb{R}^d - \{0\}.$$

Positive (or negative) definiteness implies positive (or negative) semi-definiteness.

If the Hessian is positive definite then f attains a local minimum at a (convex).

If the Hessian is negative definite at a , then f attains a local maximum at a (concave).

If the Hessian has both positive and negative eigenvalues then a is a saddle point for f .

24 Covariance Matrix

Let X be a random vector of dimension $d \times 1$ with expectation μ_X .

Matrix outer products!

$$\Sigma = \mathbb{E}[(X - \mu_X)(X - \mu_X)^T]$$
$$= \mathbb{E}[XX^T] - \mathbb{E}[X]\mathbb{E}[X]^T$$
$$= \mathbb{E}[XX^T] - \mu_X \mu_X^T$$

25 Inequalities

25.0.1 Markov's Inequality

Markov's Inequality. For a nonnegative random variable X and positive t , we have:

$$Pr[X \geq \lambda] \leq \frac{\mathbb{E}[X]}{\lambda}$$

25.0.2 Chebyshev's Inequality

We can use Markov's Inequality to derive Chebyshev's Inequality.

$$Pr[|X - \mathbb{E}(X)| \geq \lambda] \leq \frac{\text{Var}[X]}{\lambda^2}$$