

# 18.6501x Fundamentals of Statistics

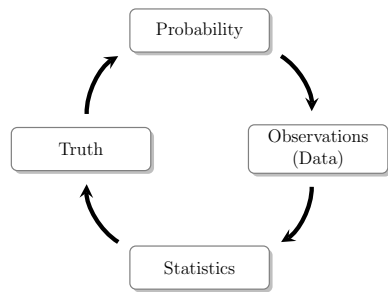
This is a cheat sheet for statistics based on the online course given by Prof. Philippe Rigollet. Compiled by Janus B. Advincula.

Last Updated December 15, 2019

## Introduction to Statistics

### What is Statistics?

**Statistical view** Data comes from a *random process*. The goal is to learn how this process works in order to make predictions or to understand what plays a role in it.



### Statistics vs. Probability

**Probability** Previous studies showed that the drug was 80% effective. Then we can anticipate that for a study on 100 patients, in average 80 will be cured and at least 65 will be cured with 99.99% chances.

**Statistics** Observe that  $\frac{78}{100}$  patients were cured. We (will be able to) conclude that we are 95% confident that for other studies, the drug will be effective on between 69.88% and 86.11% of patients.

### Probability Redux

Let  $X_1, \dots, X_n$  be i.i.d. random variables with  $\mathbb{E}[X] = \mu$  and  $\text{Var}(X) = \sigma^2$ .

#### Law of Large Numbers

$$\bar{X}_n = \frac{1}{n} \sum_{i=1}^n X_i \xrightarrow[n \rightarrow \infty]{\mathbb{P}, a.s.} \mu.$$

#### Central Limit Theorem

$$\sqrt{n} \frac{\bar{X}_n - \mu}{\sigma} \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, 1).$$

Equivalently,

$$\sqrt{n} (\bar{X}_n - \mu) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, \sigma^2).$$

**Hoeffding's Inequality** Let  $n$  be a positive integer and  $X, X_1, \dots, X_n$  be i.i.d. random variables such that  $\mathbb{E}[X] = \mu$  and  $X \in [a, b]$  almost surely. Then,

$$\mathbb{P}(|\bar{X}_n - \mu| \geq \epsilon) \leq 2e^{-\frac{2n\epsilon^2}{(b-a)^2}} \quad \forall \epsilon > 0$$

### The Gaussian Distribution

Because of the CLT, the Gaussian (a.k.a. normal) distribution is ubiquitous in statistics.

- $X \sim \mathcal{N}(\mu, \sigma^2)$
- $\mathbb{E}[X] = \mu$
- $\text{Var}(X) = \sigma^2 > 0$

**Gaussian density** (PDF)

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

**Useful Properties of Gaussian**

It is invariant under *affine transformation*.

- If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then for any  $a, b \in \mathbb{R}$ ,

$$aX + b \sim \mathcal{N}(a\mu + b, a^2\sigma^2).$$

- Standardization:** If  $X \sim \mathcal{N}(\mu, \sigma^2)$ , then

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

We can compute probabilities from the CDF of  $Z \sim \mathcal{N}(0, 1)$ :

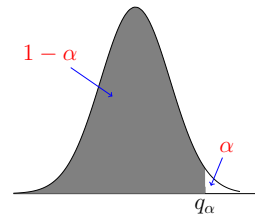
$$\mathbb{P}(u \leq X \leq v) = \mathbb{P}\left(\frac{u - \mu}{\sigma} \leq Z \leq \frac{v - \mu}{\sigma}\right)$$

- Symmetry: If  $X \sim \mathcal{N}(0, \sigma^2)$ , then  $-X \sim \mathcal{N}(0, \sigma^2)$ . If  $x > 0$ ,

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

**Quantiles** Let  $\alpha \in (0, 1)$ . The quantile of order  $1 - \alpha$  of a random variable  $X$  is the number  $q_\alpha$  such that

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



Let  $F$  denote the CDF of  $X$ .

- $F(q_\alpha) = 1 - \alpha$
- If  $F$  is invertible, then  $q_\alpha = F^{-1}(1 - \alpha)$
- $\mathbb{P}(X > q_\alpha) = \alpha$
- If  $X \sim \mathcal{N}(0, 1)$ ,  $\mathbb{P}(|X| > q_{\alpha/2}) = \alpha$

### Three Types of Convergence

**Almost Surely (a.s.) Convergence**

$$T_n \xrightarrow[n \rightarrow \infty]{a.s.} T \iff \mathbb{P}\left[\left\{\omega : T_n(\omega) \xrightarrow[n \rightarrow \infty]{} T(\omega)\right\}\right] = 1$$

**Convergence in Probability**

$$T_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} T \iff \mathbb{P}(|T_n - T| \geq \epsilon) \xrightarrow[n \rightarrow \infty]{} 0 \quad \forall \epsilon > 0$$

**Convergence in Distribution**

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} T \iff \mathbb{E}[f(T_n)] \xrightarrow[n \rightarrow \infty]{} \mathbb{E}[f(T)]$$

for all continuous and bounded function  $f$ .

### Properties

- If  $(T_n)_{n \geq 1}$  converges a.s., then it also converges in probability, and the two limits are equal.
- If  $(T_n)_{n \geq 1}$  converges in probability, then it also converges in distribution.
- Convergence in distribution implies convergence in probability if the limit has a density (e.g. Gaussian):

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} T \implies \mathbb{P}(a \leq T_n \leq b) \xrightarrow[n \rightarrow \infty]{} \mathbb{P}(a \leq T \leq b)$$

### Addition, Multiplication, Division

Assume

$$T_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}} T \quad \text{and} \quad U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}} U.$$

- $T_n + U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}} T + U$
- $T_n U_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}} TU$
- If, in addition,  $U \neq 0$  a.s., then

$$\frac{T_n}{U_n} \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}} \frac{T}{U}$$

### Slutsky's Theorem

Let  $(X_n), (Y_n)$  be two sequences of random variables such that

$$(i) T_n \xrightarrow[n \rightarrow \infty]{(d)} T \quad \text{and} \quad (ii) U_n \xrightarrow[n \rightarrow \infty]{\mathbb{P}} u$$

where  $T$  is a random variable and  $u$  is a given real number. Then,

- $T_n + U_n \xrightarrow[n \rightarrow \infty]{(d)} T + u$
- $T_n U_n \xrightarrow[n \rightarrow \infty]{(d)} Tu$
- If, in addition,  $u \neq 0$ , then  $\frac{T_n}{U_n} \xrightarrow[n \rightarrow \infty]{(d)} \frac{T}{u}$ .

### Continuous Mapping Theorem

If  $f$  is a continuous function, then

$$T_n \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}/(d)} T \implies f(T_n) \xrightarrow[n \rightarrow \infty]{a.s./\mathbb{P}/(d)} f(T).$$

## Foundation of Inference

### Statistical Model

Let the observed outcome of a statistical experiment be a *sample*  $X_1, \dots, X_n$  of  $n$  i.i.d. random variables in some measurable space  $E$  (usually  $E \subseteq \mathbb{R}$ ) and denote by  $\mathbb{P}$  their common distribution. A *statistical model* associated to that statistical experiment is a *pair*

$$(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$$

where

- $E$  is called *sample space*;
- $(\mathbb{P}_\theta)_{\theta \in \Theta}$  is a family of probability measures on  $E$ ;
- $\Theta$  is any set, called *parameter set*.

Parametric, Nonparametric and Semiparametric Models

- Usually, we will assume that the statistical model is **well-specified**, i.e., defined such that  $\exists \theta$  such that  $\mathbb{P} = \mathbb{P}_\theta$ . This particular  $\theta$  is called the **true parameter** and is unknown.
- We often assume that  $\Theta \subseteq \mathbb{R}^d$  for some  $d \geq 1$ . The model is called **parametric**.
- Sometimes we could have  $\Theta$  be infinite dimensional, in which case the model is called **nonparametric**.
- If  $\Theta = \Theta_1 \times \Theta_2$ , where  $\Theta_1$  is finite dimensional and  $\Theta_2$  is infinite dimensional, then we have a **semiparametric** model. In these models, we only care to estimate the finite dimensional parameter and the infinite dimensional one is called **nuisance parameter**.

Identifiability

The parameter  $\theta$  is called *identifiable* if and only if the map  $\theta \in \Theta \mapsto \mathbb{P}_\theta$  is injective, i.e.,

$$\theta \neq \theta' \implies \mathbb{P}_\theta \neq \mathbb{P}_{\theta'}$$

or equivalently,

$$\mathbb{P}_\theta = \mathbb{P}_{\theta'} \implies \theta = \theta'.$$

Parameter Estimation

**Statistic** Any *measurable* function of the sample, e.g.,  $\bar{X}_n, \max_i X_i$ , etc.

**Estimator of  $\theta$**  Any statistic whose expression does not depend on  $\theta$

- An estimator  $\hat{\theta}_n$  of  $\theta$  is weakly (resp.strongly) **consistent** if

$$\hat{\theta}_n \overset{\mathbb{P} \text{ (resp. a.s.)}}{\underset{n \rightarrow \infty}{\rightarrow}} \theta \quad (\text{w.r.t. } \mathbb{P}).$$

- An estimator  $\hat{\theta}_n$  of  $\theta$  is **asymptotically normal** if

$$\sqrt{n} \left( \hat{\theta}_n - \theta \right) \overset{(d)}{\underset{n \rightarrow \infty}{\rightarrow}} \mathcal{N} \left( 0, \sigma^2 \right)$$

Bias of an Estimator

- Bias** of an estimator of  $\hat{\theta}_n$  of  $\theta$ :

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

- If bias  $\left( \hat{\theta}_n \right) = 0$ , we say that  $\hat{\theta}_n$  is **unbiased**.

Jensen’s Inequality

- If the function  $f(x)$  is convex,

$$\mathbb{E} \left[ f \left( X \right) \right] \geq f \left( \mathbb{E} \left[ X \right] \right).$$

- If the function  $g(x)$  is concave,

$$\mathbb{E} \left[ g \left( X \right) \right] \leq g \left( \mathbb{E} \left[ X \right] \right).$$

Quadratic Risk

- We want estimators to have low bias and low variance at the same time.
- The **risk** (or **quadratic risk**) of an estimator  $\hat{\theta}_n \in \mathbb{R}$  is

$$R \left( \hat{\theta}_n \right) = \mathbb{E} \left[ \left| \hat{\theta}_n - \theta \right|^2 \right] = \text{variance} + \text{bias}^2$$

- Low quadratic risk means that both bias and variance are small.

Confidence Intervals

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

- Confidence interval (C.I.) of level  $1 - \alpha$  for  $\theta$ : Any random (depending on  $X_1, \dots, X_n$ ) interval  $\mathcal{I}$  whose boundaries do not depend on  $\theta$  and such that

$$\mathbb{P}_\theta \left[ \mathcal{I} \ni \theta \right] \geq 1 - \alpha, \quad \forall \theta \in \Theta.$$

- C.I. of asymptotic level  $1 - \alpha$  for  $\theta$ : Any random interval  $\mathcal{I}$  whose boundaries do not depend on  $\theta$  and such that

$$\lim_{n \rightarrow \infty} \mathbb{P}_\theta \left[ \mathcal{I} \ni \theta \right] \geq 1 - \alpha, \quad \forall \theta \in \Theta.$$

**Example** We observe  $R_1, \dots, R_n \overset{\text{iid}}{\sim} \text{Ber}(p)$  for some unknown  $p \in (0, 1)$ .

- Statistical model:  $\left( \{0, 1\}, (\text{Ber}(p))_{p \in (0, 1)} \right)$
- From CLT:

$$\sqrt{n} \frac{\bar{R}_n - p}{\sqrt{p(1-p)}} \overset{(d)}{\underset{n \rightarrow \infty}{\rightarrow}} \mathcal{N}(0, 1)$$

- It yields

$$\mathcal{I} = \left[ \bar{R}_n - \frac{q_{\frac{\alpha}{2}} \sqrt{p(1-p)}}{\sqrt{n}}, \bar{R}_n + \frac{q_{\frac{\alpha}{2}} \sqrt{p(1-p)}}{\sqrt{n}} \right]$$

- But this is **not** a confidence interval because it depends on  $p$ !

**Three solutions:**

- Conservative bound
- Solving the (quadratic) equation for  $p$
- Plug-in

The Delta Method

Let  $(Z_n)_{n \geq 1}$  be a sequence of random variables that satisfies

$$\sqrt{n} (Z_n - \theta) \overset{(d)}{\underset{n \rightarrow \infty}{\rightarrow}} \mathcal{N} \left( 0, \sigma^2 \right)$$

for some  $\theta \in \mathbb{R}$  and  $\sigma^2 > 0$  (the sequence  $(Z_n)_{n \geq 1}$  is said to be **asymptotically normal around  $\theta$** ). Let  $g : \mathbb{R} \rightarrow \mathbb{R}$  be continuously differentiable at the point  $\theta$ . Then,

- $(g(Z_n))_{n \geq 1}$  is also asymptotically normal around  $g(\theta)$ .
- More precisely,

$$\sqrt{n} (g(Z_n) - g(\theta)) \overset{(d)}{\underset{n \rightarrow \infty}{\rightarrow}} \mathcal{N} \left( 0, (g'(\theta))^2 \sigma^2 \right).$$

Introduction to Hypothesis Testing

**Statistical Formulation** Consider a sample  $X_1, \dots, X_n$  of i.i.d. random variables and a statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ . Let  $\Theta_0$  and  $\Theta_1$  be disjoint subsets of  $\Theta$ .

Consider the two hypotheses:

- $H_0 : \theta \in \Theta_0$
- $H_1 : \theta \in \Theta_1$

$H_0$  is the **null hypothesis** and  $H_1$  is the **alternative hypothesis**.

**Asymmetry in the hypotheses**  $H_0$  and  $H_1$  do not play a symmetric role: the data is only used to try to disprove  $H_0$ . Lack of evidence does not mean that  $H_0$  is true.

A test is a statistic  $\psi \in \{0, 1\}$  such that:

- If  $\psi = 0$ ,  $H_0$  is not rejected.
- If  $\psi = 1$ ,  $H_0$  is rejected.

Errors

- Rejection region** of a test  $\psi$ :

$$R_\psi = \{x \in E^n : \psi(x) = 1\}.$$

- Type 1 error** of a test  $\psi$ :

$$\begin{aligned} \alpha_\psi : \Theta_0 &\rightarrow \mathbb{R} \quad (\text{or } [0, 1]) \\ \theta &\mapsto \mathbb{P}_\theta [\psi = 1] \end{aligned}$$

- Type 2 error** of a test  $\psi$ :

$$\begin{aligned} \beta_\psi : \Theta_1 &\rightarrow \mathbb{R} \\ \theta &\mapsto \mathbb{P}_\theta [\psi = 0] \end{aligned}$$

- Power** of a test  $\psi$ :

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

Level, test statistic and rejection region

- A test  $\psi$  has level  $\alpha$  if

$$\alpha_\psi(\theta) \leq \alpha, \quad \forall \theta \in \Theta_0.$$

- A test  $\psi$  has asymptotic level  $\alpha$  if

$$\lim_{n \rightarrow \infty} \alpha_\psi(\theta) \leq \alpha, \quad \forall \theta \in \Theta_0.$$

- In general, a test has the form

$$\psi = \mathbb{1} \{T_n > c\}$$

for some statistic  $T_n$  and threshold  $c \in \mathbb{R}$ .  $T_n$  is called the **test statistic**. The rejection region is  $R_\psi = \{T_n > c\}$ .

**p-value** The (asymptotic)  $p$ -value of a test  $\psi_\alpha$  is the smallest (asymptotic) level  $\alpha$  at which  $\psi_\alpha$  rejects  $H_0$ .

Methods of Estimation

Total Variation Distance

Let  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$  be a statistical model associated with a sample of i.i.d. r.v.  $X_1, \dots, X_n$ . Assume that there exists  $\theta^* \in \Theta$  such that  $X_1 \sim \mathbb{P}_{\theta^*}$ .

**Statistician’s goal:** Given  $X_1, \dots, X_n$ , find an estimator  $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$  such that  $\mathbb{P}_{\hat{\theta}}$  is close to  $\mathbb{P}_{\theta^*}$  for the true parameter  $\theta^*$ .

The **total variation distance** between two probability measures  $\mathbb{P}_\theta$  and  $\mathbb{P}_{\theta'}$  is defined by

$$\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \max_{A \subseteq E} |\mathbb{P}_\theta(A) - \mathbb{P}_{\theta'}(A)|$$

**Total Variation Distance between Discrete Measures** Assume that  $E$  is discrete (i.e., finite or countable). The total variation distance between  $\mathbb{P}_\theta$  and  $\mathbb{P}_{\theta'}$  is

$$\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \frac{1}{2} \sum_{x \in E} |p_\theta(x) - p_{\theta'}(x)|$$

**Total Variation Distance between Continuous Measures** Assume that  $E$  is continuous. The total variation distance between  $\mathbb{P}_\theta$  and  $\mathbb{P}_{\theta'}$  is

$$\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \frac{1}{2} \int |f_\theta(x) - f_{\theta'}(x)| \, dx$$

Properties of Total Variation

- $\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \text{TV}(\mathbb{P}_{\theta'}, \mathbb{P}_\theta)$  **symmetric**
- $\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \geq 0, \text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq 1$  **positive**
- If  $\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = 0$ , then  $\mathbb{P}_\theta = \mathbb{P}_{\theta'}$  **definite**
- $\text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \leq \text{TV}(\mathbb{P}_\theta, \mathbb{P}_{\theta''}) + \text{TV}(\mathbb{P}_{\theta''}, \mathbb{P}_{\theta'})$  **triangle inequality**

These imply that the total variation is a **distance** between probability distributions.

## Kullback-Leibler (KL) Divergence

The Kullback-Leibler (KL) divergence between two probability measures  $\mathbb{P}_\theta$  and  $\mathbb{P}_{\theta'}$  is defined by

$$\text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = \begin{cases} \sum_{x \in E} p_\theta(x) \log \left( \frac{p_\theta(x)}{p_{\theta'}(x)} \right) & \text{if } E \text{ is discrete} \\ \int_E f_\theta(x) \log \left( \frac{f_\theta(x)}{f_{\theta'}(x)} \right) dx & \text{if } E \text{ is continuous} \end{cases}$$

KL-divergence is also known as **relative entropy**.

**Properties of KL-divergence**

- $\text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \neq \text{KL}(\mathbb{P}_{\theta'}, \mathbb{P}_\theta)$  in general
- $\text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \geq 0$
- If  $\text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) = 0$ , then  $\mathbb{P}_\theta = \mathbb{P}_{\theta'}$  (definite)
- $\text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta'}) \not\leq \text{KL}(\mathbb{P}_\theta, \mathbb{P}_{\theta''}) + \text{KL}(\mathbb{P}_{\theta''}, \mathbb{P}_{\theta'})$  in general

## Maximum Likelihood Estimation

**Likelihood, Discrete Case** Let  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$  be a statistical model associated with a sample of i.i.d. r.v.  $X_1, \dots, X_n$ . Assume that  $E$  is discrete (i.e., finite or countable).

**Definition** The likelihood of the model is the map  $L_n$  (or just  $L$ ) defined as

$$\begin{aligned} L_n : E^n \times \Theta &\rightarrow \mathbb{R} \\ (x_1, \dots, x_n; \theta) &\mapsto \mathbb{P}_\theta[X_1 = x_1, \dots, X_n = x_n] \\ &= \prod_{i=1}^n \mathbb{P}_\theta[X_i = x_i] \end{aligned}$$

**Likelihood, Continuous Case** Let  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$  be a statistical model associated with a sample of i.i.d. r.v.  $X_1, \dots, X_n$ . Assume that all the  $\mathbb{P}_\theta$  have density  $f_\theta$ .

**Definition** The likelihood of the model is the map  $L$  defined as

$$\begin{aligned} L : E^n \times \Theta &\rightarrow \mathbb{R} \\ (x_1, \dots, x_n; \theta) &\mapsto \prod_{i=1}^n f_\theta(x_i) \end{aligned}$$

**Maximum Likelihood Estimator** Let  $X_1, \dots, X_n$  be an i.i.d. sample associated with a statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$  and let  $L$  be the corresponding likelihood.

**Definition** The maximum likelihood estimator of  $\theta$  is defined as

$$\hat{\theta}_n^{\text{MLE}} = \underset{\theta \in \Theta}{\operatorname{argmax}} L(X_1, \dots, X_n, \theta),$$

provided it exists.

**Log-likelihood Estimator** In practice, we use the fact that

$$\hat{\theta}_n^{\text{MLE}} = \underset{\theta \in \Theta}{\operatorname{argmax}} \log L(X_1, \dots, X_n, \theta),$$

## Concave and Convex Functions

A twice-differentiable function  $h : \Theta \subset \mathbb{R} \rightarrow \mathbb{R}$  is said to be **concave** if its second derivative satisfies

$$h''(\theta) \leq 0, \quad \forall \theta \in \Theta.$$

It is said to be **strictly concave** if the inequality is strict:  $h''(\theta) < 0$ . Moreover,  $h$  is said to be (strictly) **convex** if  $-h$  is (strictly) concave, i.e.  $h''(\theta) \geq 0$  ( $h''(\theta) > 0$ ).

**Multivariate Concave Functions** More generally, for a multivariate function:

$h : \Theta \subset \mathbb{R}^d \rightarrow \mathbb{R}$ ,  $d \geq 2$ , define the

- gradient vector:**

$$\nabla h(\theta) = \begin{pmatrix} \frac{\partial h(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial h(\theta)}{\partial \theta_d} \end{pmatrix} \in \mathbb{R}^d$$

- Hessian matrix:**

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

$h$  is concave  $\iff x^\top \mathbb{H}h(\theta)x \leq 0, \forall x \in \mathbb{R}^d, \theta \in \Theta$

$h$  is strictly concave  $\iff x^\top \mathbb{H}h(\theta)x < 0, \forall x \in \mathbb{R}^d, \theta \in \Theta$

**Consistency of Maximum Likelihood Estimator** Under mild regularity conditions, we have

$$\hat{\theta}_n^{\text{MLE}} \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \theta^*$$

**Covariance** In general, when  $\theta \in \mathbb{R}^d$ ,  $d \geq 2$ , its coordinates are not necessarily independent. The covariance between two random variables  $X$  and  $Y$  is

$$\begin{aligned} \text{Cov}(X, Y) &:= \mathbb{E}[(X - \mathbb{E}[X])(Y - \mathbb{E}[Y])] \\ &= \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y] \end{aligned}$$

**Properties**

- $\text{Cov}(X, X) = \text{Var}(X)$
- $\text{Cov}(X, Y) = \text{Cov}(Y, X)$
- If  $X$  and  $Y$  are independent, then  $\text{Cov}(X, Y) = 0$

**Covariance Matrix** The covariance matrix of a random vector

$$X = \left( X^{(1)}, \dots, X^{(d)} \right)^\top \in \mathbb{R}^d$$

is given by

$$\Sigma = \text{Cov}(X) = \mathbb{E}[(X - \mathbb{E}[X])(X - \mathbb{E}[X])^\top].$$

This is a matrix of size  $d \times d$ .

If  $X \in \mathbb{R}^d$  and  $A, B$  are matrices:

$$\text{Cov}(AX + B) = \text{Cov}(AX) = A \text{Cov}(X) A^\top = A \Sigma_X A^\top$$

**The Multivariate Gaussian Distribution** If  $(X, T)^\top$  is a Gaussian vector then its PDF depends on 5 parameters:

$$\mathbb{E}[X], \text{Var}(X), \mathbb{E}[Y], \text{Var}(Y), \text{ and } \text{Cov}(X, Y).$$

A Gaussian vector  $X \in \mathbb{R}^d$  is completely determined by its expected value and covariance matrix  $\Sigma$ :

$$X \sim \mathcal{N}_d(\mu, \Sigma).$$

It has PDF over  $\mathbb{R}^d$  given by:

$$f(x) = \frac{1}{((2\pi)^d \det(\Sigma))^{\frac{1}{2}}} \exp \left( -\frac{1}{2} (x - \mu)^\top \Sigma^{-1} (x - \mu) \right)$$

**The Multivariate CLT** Let  $X_1, \dots, X_n \in \mathbb{R}^d$  be independent copies of a random vector  $X$  such that  $\mathbb{E}[X] = \mu$ ,  $\text{Cov}(X) = \Sigma$ , then

$$\sqrt{n} \left( \bar{X}_n - \mu \right) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}_d(0, \Sigma)$$

**Multivariate Delta Method** Let  $(T_n)_{n \geq 1}$  sequence of random vectors in  $\mathbb{R}^d$  such that

$$\sqrt{n} (T_n - \theta) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}_d(0, \Sigma),$$

for some  $\theta \in \mathbb{R}^d$  and some covariance  $\Sigma \in \mathbb{R}^{d \times d}$ . Let  $g : \mathbb{R}^d \rightarrow \mathbb{R}^k$  ( $k \geq 1$ ) be continuously differentiable at  $\theta$ . Then,

$$\sqrt{n} (g(T_n) - g(\theta)) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, \nabla g(\theta)^\top \Sigma \nabla g(\theta)),$$

where  $\nabla g(\theta) = \frac{\partial g(\theta)}{\partial \theta} = \left( \frac{\partial g_j}{\partial \theta_i} \right)_{\substack{1 \leq i \leq d \\ 1 \leq j \leq k}} \in \mathbb{R}^{d \times k}$

## Fisher Information

Define the log-likelihood for one observation as

$$\ell(\theta) = \log L_1(X, \theta), \quad \theta \in \Theta \subset \mathbb{R}^d.$$

Assume that  $\ell$  is a.s. twice differentiable. Under some regularity conditions, the Fisher information of the statistical model is defined as

$$I(\theta) = \mathbb{E}[\nabla \ell(\theta) \nabla \ell(\theta)^\top] - \mathbb{E}[\nabla \ell(\theta)] \mathbb{E}[\nabla \ell(\theta)]^\top = -\mathbb{E}[\mathbb{H} \ell(\theta)].$$

If  $\Theta \subset \mathbb{R}$ , we get

$$I(\theta) = \text{Var}[\ell'(\theta)] = -\mathbb{E}[\ell''(\theta)].$$

## Asymptotic Normality of the MLE

**Theorem** Let  $\theta^* \in \Theta$  (the true parameter). Assume the following:

- The parameter is identifiable.
- For all  $\theta \in \Theta$ , the support of  $\mathbb{P}_\theta$  does not depend on  $\theta$ .
- $\theta^*$  is not on the boundary of  $\Theta$ .
- $I(\theta)$  is invertible in a neighborhood of  $\theta^*$ .
- A few more technical conditions.

Then,  $\hat{\theta}_n^{\text{MLE}}$  satisfies

- $\hat{\theta}_n^{\text{MLE}} \xrightarrow[n \rightarrow \infty]{\mathbb{P}} \theta^*$  w.r.t.  $\mathbb{P}_{\theta^*}$ ;
- $\sqrt{n} \left( \hat{\theta}_n^{\text{MLE}} - \theta^* \right) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}_d(0, I^{-1}(\theta^*))$  w.r.t.  $\mathbb{P}_{\theta^*}$ .

## The Method of Moments

### Moments

Let  $X_1, \dots, X_n$  be an i.i.d. sample associated with a statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ . Assume that  $E \subseteq \mathbb{R}$  and  $\Theta \subseteq \mathbb{R}^d$ , for some  $d \geq 1$ .

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

**Empirical Moments** Let  $\hat{m}_k = \overline{X_n^k} = \frac{1}{n} \sum_{i=1}^n X_i^k$ ,  $1 \leq k \leq d$ .

From LLN,

$$\hat{m}_k \xrightarrow[n \rightarrow \infty]{\mathbb{P}/a.s.} m_k(\theta)$$

More compactly, we say that the whole vector converges:

$$(\hat{m}_1, \dots, \hat{m}_d) \xrightarrow[n \rightarrow \infty]{\mathbb{P}/a.s.} (m_1(\theta), \dots, m_d(\theta))$$

Moments Estimator

Let

M : \Theta \to \mathbb{R}^d  
\theta \mapsto M(\theta) = (m\_1(\theta), \dots, m\_d(\theta))

Assume M is one-to-one:

\theta = M^{-1}(m\_1(\theta), \dots, m\_d(\theta))

Moments estimator of \theta:

\hat{\theta}\_n^{MM} = M^{-1}(\hat{m}\_1, \dots, \hat{m}\_d)

provided it exists.

Generalized Method of Moments

Applying the multivariate CLT and Delta method yields:

Theorem

\sqrt{n}(\hat{\theta}\_n^{MM} - \theta) \xrightarrow[n \to \infty]{(d)} \mathcal{N}(0, \Gamma(\theta)),

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]

MLE vs. Moment Estimator

- Comparison of the quadratic risks: In general, the MLE is more accurate.
- MLE still gives good results if the model is misspecified.
- Computational issues: Sometimes, the MLE is intractable but MM is easier (polynomial equations).

M-Estimation

- Let X\_1, \dots, X\_n be i.i.d. with some unknown distribution \mathbb{P} in some sample space E (E \subseteq \mathbb{R}^d for some d \ge 1).
- No statistical model needs to be assumed (similar to ML).
- The goal is to estimate some parameter \mu^\* associated with \mathbb{P}, e.g. its mean, variance, median, other quantiles, the true parameter in some statistical model, etc.
- We want to find a function \rho : E \times \mathcal{M} \to \mathbb{R}, where \mathcal{M} is the set of all possible values for the unknown \mu^\*, such that

Q(\mu) := \mathbb{E}[\rho(X\_1, \mu)]

achieves its minimum at \mu = \mu^\*.

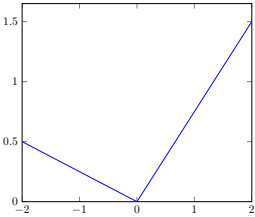
Examples (1)

- If E = \mathcal{M} = \mathbb{R} and \rho(x, \mu) = (x - \mu)^2, for all x, \mu \in \mathbb{R}: \mu^\* = \mathbb{E}[X].
- If E = \mathcal{M} = \mathbb{R}^d and \rho(x, \mu) = \|x - \mu\|\_2^2, for all x, \mu \in \mathbb{R}^d: \mu^\* = \mathbb{E}[X] \in \mathbb{R}^d.
- If E = \mathcal{M} = \mathbb{R} and \rho(x, \mu) = |x - \mu|, for all x, \mu \in \mathbb{R}: \mu^\* is a **median** of \mathbb{P}.

Example (2) If E = \mathcal{M} = \mathbb{R}, \alpha \in (0, 1) is fixed and \rho(x, \mu) = C\_\alpha(x - \mu), for all x, \mu \in \mathbb{R}: \mu^\* is a \alpha-quantile of \mathbb{P}.

Check Function

C\_\alpha = \begin{cases} -(1 - \alpha)x & \text{if } x < 0 \\ \alpha x & \text{if } x \geq 0. \end{cases}



MLE is an M-estimator Assume that (E, (\mathbb{P}\_\theta)\_{\theta \in \Theta}) is a statistical model associated with the data.

Theorem Let \mathcal{M} = \Theta and \rho(x, \theta) = -\log L\_1(x, \theta), provided the likelihood is positive everywhere. Then,

\mu^\* = \theta^\*,

where \mathbb{P} = \mathbb{P}\_{\theta^\*} (i.e., \theta^\* is the true value of the parameter).

Statistical Analysis

- Define \hat{\mu}\_n as a minimizer of
- Q\_n(\mu) := \frac{1}{n} \sum\_{i=1}^n \rho(X\_i, \mu).
- Let J(\mu) = \frac{\partial^2 Q(\mu)}{\partial \mu \partial \mu^T}.
- Under some regularity conditions, J(\mu) = \mathbb{E} \left[ \frac{\partial^2 \rho(X\_1, \mu)}{\partial \mu \partial \mu^T} \right]
- Let K(\mu) = \text{Cov} \left( \frac{\partial \rho(X\_1, \mu)}{\partial \mu} \right)
- Remark: In the log-likelihood case,
- J(\theta) = K(\theta) = I(\theta) (Fisher information)

Asymptotic Normality Let \mu^\* \in \mathcal{M} (the true parameter). Assume the following:

1. \mu^\* is the only minimizer of the function Q,
2. J(\mu) is invertible for all \mu \in \mathcal{M},
3. A few more technical conditions.

Then, \hat{\mu}\_n satisfies

- \hat{\mu}\_n \xrightarrow[n \to \infty]{\mathbb{P}} \mu^\*
- \sqrt{n}(\hat{\mu}\_n - \mu^\*) \xrightarrow[n \to \infty]{(d)} \mathcal{N}(0, J(\mu^\*)^{-1} K(\mu^\*) J(\mu^\*)^{-1})

Hypothesis Testing

Parametric Hypothesis Testing

Hypotheses

H\_0 : \Delta\_c = \Delta\_d \quad \text{vs.} \quad H\_1 : \Delta\_d > \Delta\_c

Since the data is Gaussian by assumption, we don't need the CLT.

\bar{X}\_n \sim \mathcal{N}\left(\Delta\_d, \frac{\sigma\_d^2}{n}\right) \quad \text{and} \quad \bar{Y}\_m \sim \mathcal{N}\left(\Delta\_c, \frac{\sigma\_c^2}{m}\right)

Then,

\frac{\bar{X}\_n - \bar{Y}\_m - (\Delta\_d - \Delta\_c)}{\sqrt{\frac{\sigma\_d^2}{n} + \frac{\sigma\_c^2}{m}}} \sim \mathcal{N}(0, 1)

Asymptotic test Assume that m = cn and n \to \infty

Using Slutsky's theorem, we also have

\frac{\bar{X}\_n - \bar{Y}\_m - (\Delta\_d - \Delta\_c)}{\sqrt{\frac{\hat{\sigma}\_d^2}{n} + \frac{\hat{\sigma}\_c^2}{m}}} \xrightarrow[n \to \infty]{(d)} \mathcal{N}(0, 1)

where \hat{\sigma}\_d^2 = \frac{1}{n-1} \sum\_{i=1}^n (X\_i - \bar{X}\_n)^2 and \hat{\sigma}\_c^2 = \frac{1}{m-1} \sum\_{i=1}^m (Y\_i - \bar{Y}\_m)^2

We get the following test at asymptotic level \alpha:

R\_\psi = \left\{ \frac{\bar{X}\_n - \bar{Y}\_m}{\sqrt{\frac{\hat{\sigma}\_d^2}{n} + \frac{\hat{\sigma}\_c^2}{m}}} > q\_\alpha \right\}

The \chi^2 Distribution

Definition For a positive integer d, the \chi^2 distribution with d degrees of freedom is the law of the random variable Z\_1^2 + \dots + Z\_d^2, where Z\_1, \dots, Z\_d \stackrel{iid}{\sim} \mathcal{N}(0, 1).

Properties If V \sim \chi\_k^2, then

- \mathbb{E}[V] = \mathbb{E}[Z\_1^2] + \dots + \mathbb{E}[Z\_d^2] = d
- \text{Var}(V) = \text{Var}(Z\_1^2) + \dots + \text{Var}(Z\_d^2) = 2d

Sample Variance 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

Cochran's Theorem If X\_1, \dots, X\_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2), then

- \bar{X}\_n \perp\!\!\!\perp S\_n, for all n.
- \frac{n S\_n}{\sigma^2} \sim \chi\_{n-1}^2

We often prefer the unbiased estimator of \sigma^2:

\tilde{S}\_n = \frac{1}{n-1} \sum\_{i=1}^n (X\_i - \bar{X}\_n)^2 = \frac{n}{n-1} S\_n

Student's T Distribution

Definition For a positive integer d, the Student's T distribution with d degrees of freedom (denoted by t\_d) is the law of the random variable \frac{Z}{\sqrt{V/d}}, where

Z \sim \mathcal{N}(0, 1), V \sim \chi\_d^2 and Z \perp\!\!\!\perp V.

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

Let X\_1, \dots, X\_n \stackrel{iid}{\sim} \mathcal{N}(\mu, \sigma^2) where both \mu and \sigma^2 are unknown. We want to test:

H\_0 : \mu = 0 \quad \text{vs.} \quad H\_1 : \mu \neq 0

Test statistic:

T\_n = \sqrt{n} \frac{\bar{X}\_n}{\sqrt{\tilde{S}\_n}} = \frac{\sqrt{n} \bar{X}\_n - \mu}{\sqrt{\frac{\sigma}{\tilde{S}\_n}}}

Since \sqrt{n} \frac{\bar{X}\_n}{\sigma} \sim \mathcal{N}(0, 1) (under H\_0) and \frac{\tilde{S}\_n}{\sigma^2} \sim \frac{\chi\_{n-1}^2}{n-1} are independent by Cochran's theorem, we have

T\_n \sim t\_{n-1}.

Student's test with (non-asymptotic) level \alpha \in (0, 1):

\psi\_\alpha = \mathbb{1} \left\{ |T\_n| > q\_{\frac{\alpha}{2}} \right\},

where q\_{\frac{\alpha}{2}} is the (1 - \frac{\alpha}{2})-quantile of t\_{n-1}.

Student’s T test (one-sample, one-sided)

$H_0 : \mu \leq \mu_0 \quad \text{vs.} \quad H_1 : \mu > \mu_0$

Test statistic:

$$T_n = \sqrt{n} \frac{\overline{X}_n - \mu_0}{\sqrt{\widehat{S}_n}} \sim t_{n-1} \quad (\text{under } H_0)$$

Student’s test with (non-asymptotic) level  $\alpha \in (0, 1)$ :

$\psi_\alpha = \mathbb{1} \{T_n > q_\alpha\}$

where  $q_\alpha$  is the  $(1 - \alpha)$ -quantile of  $t_{n-1}$ .

Two-sample T-test

$$\frac{\overline{X}_n - \overline{Y}_m - (\Delta_d - \Delta_c)}{\sqrt{\frac{\widehat{\sigma}_d^2}{n} + \frac{\widehat{\sigma}_c^2}{m}}} \sim t_N$$

Welch-Satterthwaite formula

$$N = \frac{\left(\frac{\widehat{\sigma}_d^2}{n} + \frac{\widehat{\sigma}_c^2}{m}\right)^2}{\frac{\widehat{\sigma}_d^4}{n^2(n-1)} + \frac{\widehat{\sigma}_c^4}{m^2(m-1)}} \geq \min(n, m)$$

Wald’s Test

**A test based on the MLE** Consider an i.i.d. sample  $X_1, \dots, X_n$  with statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ , where  $\Theta \subseteq \mathbb{R}^d$  ( $d \geq 1$ ) and let  $\theta_0 \in \Theta$  be fixed and given.  $\theta^*$  is the true parameter.

Consider the following hypotheses:

$H_0 : \theta^* = \theta_0 \quad \text{vs.} \quad H_1 : \theta^* \neq \theta_0$

Let  $\widehat{\theta}_n^{\text{MLE}}$  be the MLE. Assume the MLE technical conditions are satisfied.

If  $H_0$  is true, then

$$\sqrt{n} \, I \left( \widehat{\theta}^{\text{MLE}} \right)^{\frac{1}{2}} \left( \widehat{\theta}_n^{\text{MLE}} - \theta_0 \right) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}_d(0, \mathbb{I}_d)$$

Wald’s test

$$T_n := n \left( \widehat{\theta}_n^{\text{MLE}} - \theta_0 \right)^\top I \left( \widehat{\theta}_n^{\text{MLE}} \right) \left( \widehat{\theta}_n^{\text{MLE}} - \theta_0 \right) \xrightarrow[n \rightarrow \infty]{(d)} \chi_d^2$$

Wald’s test with asymptotic level  $\alpha \in (0, 1)$ :

$\psi = \mathbb{1} \{T_n > q_\alpha\},$

where  $q_\alpha$  is the  $(1 - \alpha)$ -quantile of  $\chi_d^2$ .

**Wald’s Test in 1 dimension** In one dimension, Wald’s test coincides with the two-sided test based on the asymptotic normality of the MLE.

Likelihood Ratio Test

**Basic Form of the Likelihood Ratio Test** Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathbb{P}_{\theta^*}$ , and consider the associated statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \mathbb{R}^d})$ . Suppose that  $\mathbb{P}_\theta$  is a discrete probability distribution with pmf given by  $p_\theta$ .

In its most basic form, the likelihood ratio test can be used to decide between two hypotheses of the following form:

$H_0 : \theta^* = \theta_0 \quad \text{vs.} \quad H_1 : \theta^* = \theta_1$

Likelihood function

$$L_n : \mathbb{R}^n \times \mathbb{R}^d \rightarrow \mathbb{R}$$
$$(x_1, \dots, x_n; \theta) \mapsto \prod_{i=1}^n p_\theta(x_i)$$

The likelihood ratio test in this set-up is of the form

$$\psi_C = \mathbb{1} \left( \frac{L_n(x_1, \dots, x_n; \theta_1)}{L_n(x_1, \dots, x_n; \theta_0)} > C \right)$$

where  $C$  is a threshold to be specified.

**A test based on the log-likelihood** Consider an i.i.d. sample  $X_1, \dots, X_n$  with statistical model  $(E, (\mathbb{P}_\theta)_{\theta \in \Theta})$ , where  $\Theta \subseteq \mathbb{R}^d$  ( $d \geq 1$ ). Suppose the null hypothesis has the form

$$H_0 : (\theta_{r+1}, \dots, \theta_d) = \left( \theta_{r+1}^{(0)}, \dots, \theta_d^{(0)} \right),$$

for some fixed and given numbers  $\theta_{r+1}^{(0)}, \dots, \theta_d^{(0)}$ .

Let

$$\widehat{\theta}_n = \operatorname{argmax}_{\theta \in \Theta} \ell_n(\theta) \quad (\text{MLE})$$

and

$$\widehat{\theta}_n^c = \operatorname{argmax}_{\theta \in \Theta_0} \ell_n(\theta) \quad (\text{constrained MLE})$$

where  $\Theta_0 = \left\{ \theta \in \Theta : (\theta_{r+1}, \dots, \theta_d) = \left( \theta_{r+1}^{(0)}, \dots, \theta_d^{(0)} \right) \right\}$

Test statistic:

$$T_n = 2 \left( \ell_n \left( \widehat{\theta}_n \right) - \ell_n \left( \widehat{\theta}_n^c \right) \right).$$

**Wilk’s Theorem** Assume  $H_0$  is true and the MLE technical conditions are satisfied. Then,

$$T_n \xrightarrow[n \rightarrow \infty]{(d)} \chi_{d-r}^2$$

Likelihood ratio test with asymptotic level  $\alpha \in (0, 1)$ :

$\psi = \mathbb{1} \{T_n > q_\alpha\},$

where  $q_\alpha$  is the  $(1 - \alpha)$ -quantile of  $\chi_{d-r}^2$ .

Goodness of Fit Tests

Let  $X$  be a r.v. We want to know if the hypothesized distribution is a good fit for the data.

Key characteristic of Goodness of Fit tests: no parametric modeling.

**Discrete distribution** Let  $E = \{a_1, \dots, a_K\}$  be a finite space and  $(\mathbb{P}_\mathbf{p})_{\mathbf{p} \in \Delta_K}$  be the family of all probability distributions on  $E$ .

$$\bullet \quad \Delta_K = \left\{ \mathbf{p} = (p_1, \dots, p_K) \in (0, 1)^K : \sum_{j=1}^K p_j = 1 \right\}$$

$$\bullet \quad \text{For } \mathbf{p} \in \Delta_K \text{ and } X \sim \mathbb{P}_\mathbf{p},$$

$$\mathbb{P}_\mathbf{p} [X = a_j] = p_j, \quad j = 1, \dots, K.$$

Let  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \mathbb{P}_\mathbf{p}$ , for some unknown  $\mathbf{p} \in \Delta_K$ , and let  $\mathbf{p}^0 \in \Delta_K$  be fixed.

We want to test:

$$H_0 : \mathbf{p} = \mathbf{p}^0 \quad \text{vs.} \quad H_1 : \mathbf{p} \neq \mathbf{p}^0$$

with asymptotic level  $\alpha \in (0, 1)$ .

**The Probability Simplex in  $K$  Dimensions** The probability simplex in  $\mathbb{R}^K$ , denoted by  $\Delta_K$ , is the set of all vectors  $\mathbf{p} = [p_1, \dots, p_K]^\top$  such that

$$\mathbf{p} \cdot \mathbf{1} = \mathbf{p}^\top \mathbf{1} = 1, \quad p_i \geq 0 \quad \text{for all } K$$

where  $\mathbf{1}$  denotes the vector  $\mathbf{1} = (1, \dots, 1)^\top$

Categorical Likelihood

- Likelihood of the model:

$$L_n(X_1, \dots, X_n; \mathbf{p}) = p_1^{N_1} p_2^{N_2} \dots p_K^{N_K}$$

where  $N_j = \# \{i = 1, \dots, n : X_i = a_j\}$ .

- Let  $\widehat{\mathbf{p}}$  be the MLE:

$$\widehat{\mathbf{p}}_j = \frac{N_j}{n}, \quad j = 1, \dots, K.$$

$\widehat{\mathbf{p}}$  maximizes  $\log L_n(X_1, \dots, X_n, \mathbf{p})$  under the constraint.

**$\chi^2$  test** If  $H_0$  is true, then  $\sqrt{n} (\widehat{\mathbf{p}} - \mathbf{p}^0)$  is asymptotically normal, and the following holds:

**Theorem** Under  $H_0$ :

$$T_n = n \sum_{j=1}^n \frac{\left( \widehat{\mathbf{p}}_j - \mathbf{p}_j^0 \right)^2}{\mathbf{p}_j^0} \xrightarrow[n \rightarrow \infty]{(d)} \chi_{K-1}^2$$

**CDF and empirical CDF** Let  $X_1, \dots, X_n$  be i.i.d. real random variables. The CDF of  $X_1$  is defined as

$$F(t) = \mathbb{P} [X_1 \leq 1], \quad \forall t \in \mathbb{R}.$$

It completely characterizes the distribution of  $X_1$ .

The **empirical CDF** of the sample  $X_1, \dots, X_n$  is defined as

$$F_n(t) = \frac{1}{n} \sum_{i=1}^n \mathbb{1} \{X_i \leq 1\}$$
$$= \frac{\# \{i = 1, \dots, n : X_i \leq t\}}{n}, \quad \forall t \in \mathbb{R}.$$

**Consistency** By the LLN, for all  $t \in \mathbb{R}$ ,

$$F_n(t) \xrightarrow[n \rightarrow \infty]{a.s.} F(t).$$

**Glivenko-Cantelli Theorem (Fundamental theorem of statistics)**

$$\sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow \infty]{a.s.} 0$$

**Asymptotic normality** By the CLT, for all  $t \in \mathbb{R}$ ,

$$\sqrt{n} (F_n(t) - F(t)) \xrightarrow[n \rightarrow \infty]{(d)} \mathcal{N}(0, F(t) (1 - F(t)))$$

**Donsker’s Theorem** If  $F$  is continuous, then

$$\sqrt{n} \sup_{t \in \mathbb{R}} |F_n(t) - F(t)| \xrightarrow[n \rightarrow \infty]{a.s.} \sup_{0 \leq t \leq 1} |\mathbf{B}(t)|,$$

where  $\mathbf{B}(t)$  is a Brownian bridge on  $[0, 1]$ .

Kolmogorov-Smirnov Test

Let  $T_n = \sup_{t \in \mathbb{R}} \sqrt{n} |F_n(t) - F(t)|$ . By Donsker’s theorem, if  $H_0$  is true, then

$T_n \xrightarrow[n \rightarrow \infty]{(d)} Z$ , where  $Z$  has a known distribution (supremum of the absolute value of a Brownian bridge).

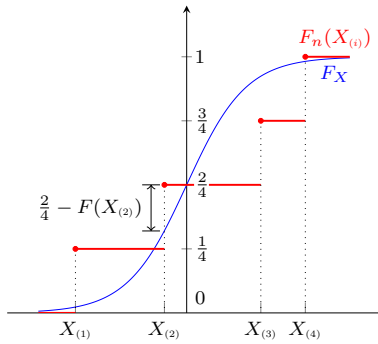
**KS test with asymptotic level  $\alpha$ :**

$$\delta_\alpha^{\text{KS}} = \mathbb{1} \{T_n > q_\alpha\}$$

where  $q_\alpha$  is the  $(1 - \alpha)$ -quantile of  $Z$ .

Let  $X_{(1)} \leq X_{(2)} \leq \dots \leq X_{(n)}$  be the reordered sample. The expression for  $T_n$  reduces to

$$T_n = \sqrt{n} \max_{i=1, \dots, n} \left\{ \max \left( \left| \frac{i-1}{n} - F^0(X_{(i)}) \right|, \left| \frac{i}{n} - F^0(X_{(i)}) \right| \right) \right\}.$$



**Pivotal Distribution**  $T_n$  is called a **pivotal statistic**: If  $H_0$  is true, the distribution of  $T_n$  does not depend on the distribution of the  $X_i$ 's.

## Other Goodness of Fit Tests

### Kolmogorov-Smirnov

$$d(F_n, F) = \sup_{t \in \mathbb{R}} |F_n(t) - F(t)|$$

### Cramér-Von Mises

$$\begin{aligned} d^2(F_n, F) &= \int_{\mathbb{R}} [F_n(t) - F(t)]^2 dF(t) \\ &= \mathbb{E}_{X \sim F} [|F_n(X) - F(X)|^2] \end{aligned}$$

### Anderson-Darling

$$d^2(F_n, F) \int_{\mathbb{R}} \frac{[F_n(t) - F(t)]^2}{F(t)(1 - F(t))} dF(t)$$

### Kolmogorov-Lilliefors Test

We want to test if  $X$  has a Gaussian distribution with unknown parameters. In this case, Donsker's theorem is *no longer valid*. Instead, we compute the quantiles for the test statistic

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

where  $\hat{\mu} = \bar{X}_n$ ,  $\hat{\sigma}^2 = S_n^2$  and  $\Phi_{\hat{\mu}, \hat{\sigma}^2}(t)$  is the CDF of  $\mathcal{N}(\hat{\mu}, \hat{\sigma}^2)$ .

They do not depend on unknown parameters.

### Quantile-Quantile (QQ) plots

- Provide a visual way to perform goodness of fit tests.
- Not a formal test but quick and easy check to see if a distribution is plausible.
- Main idea: We want to check visually if the plot of  $F_n$  is close to that of  $F$  or, equivalently, if the plot of  $F_n^{-1}$  is close to  $F^{-1}$ .
- Check if the points

$$\left(F^{-1}\left(\frac{1}{n}\right), F_n^{-1}\left(\frac{1}{n}\right)\right), \dots, \left(F^{-1}\left(\frac{n-1}{n}\right), F_n^{-1}\left(\frac{n-1}{n}\right)\right)$$

are near the line  $y = x$ .

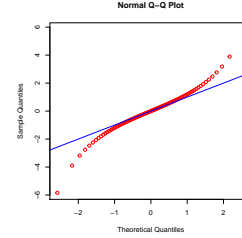
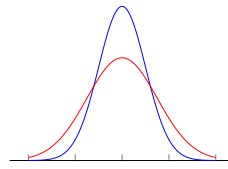
- $F_n$  is not technically invertible but we define

$$F_n^{-1}\left(\frac{i}{n}\right) = X_i,$$

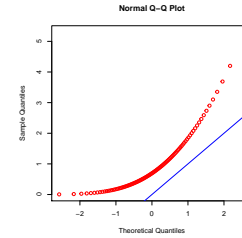
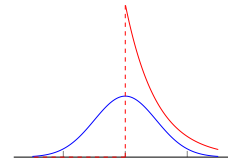
the  $i^{\text{th}}$  largest observation.

### Four patterns

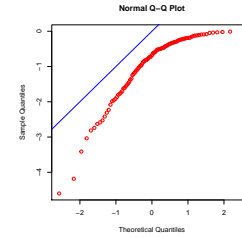
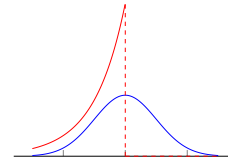
1. heavy tails



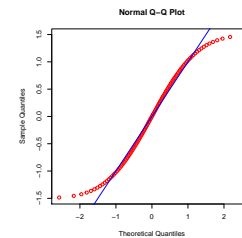
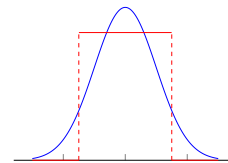
2. right skewed



3. left skewed



4. light tails



- Let  $X_1, \dots, X_n$  be a sample of  $n$  random variables.
- Denote by  $L_n(\cdot|\theta)$  the joint PDF of  $X_1, \dots, X_n$  conditionally on  $\theta$ , where  $\theta \sim \pi$ .
- **Remark:**  $L_n(X_1, \dots, X_n|\theta)$  is the likelihood used in the frequentist approach.
- The conditional distribution of  $\theta$  given  $X_1, \dots, X_n$  is called the **posterior distribution**. Denote by  $\pi(\cdot|X_1, \dots, X_n)$  its PDF.

### Bayes' formula

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

### Bernoulli experiment with a Beta prior

- $p \sim \text{Beta}(a, a)$ :

$$\pi(p) \propto p^{a-1}(1-p)^{a-1}, \quad p \in (0, 1)$$

- Given  $p$ ,  $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} \text{Ber}(p)$ , so

$$L_n(X_1, \dots, X_n|p) = p^{\sum_{i=1}^n X_i} (1-p)^{n - \sum_{i=1}^n X_i}.$$

- Hence,

$$\pi(p|X_1, \dots, X_n) \propto p^{a-1 + \sum_{i=1}^n X_i} (1-p)^{a-1 + n - \sum_{i=1}^n X_i}$$

- The posterior distribution is

$$\text{Beta}\left(a + \sum_{i=1}^n X_i, a + n - \sum_{i=1}^n X_i\right) \quad \text{conjugate prior}$$

### Non-informative Priors

- We can still use a Bayesian approach if we have no prior information about the parameter.
- Good candidate:  $\pi(\theta) \propto 1$ , i.e., constant PDF on  $\Theta$ .
- If  $\Theta$  is bounded, this is the uniform prior on  $\Theta$ .
- If  $\Theta$  is unbounded, this define a proper PDF on  $\Theta$ .
- An **improper prior** on  $\Theta$  is a measurable, non-negative function  $\pi(\cdot)$  defined on  $\Theta$  that is not integrable:

$$\int \pi(\theta) d\theta = \infty.$$

- In general, one can still define a posterior distribution using an improper prior, using Bayes' formula.

## Jeffreys Prior and Bayesian Confidence Interval

Jeffreys prior is an attempt to incorporate frequentist ideas of likelihood in the Bayesian framework, as well as an example of a **non-informative prior**:

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

where  $I(\theta)$  is the Fisher information matrix of the statistical model associated with  $X_1, \dots, X_n$  in the frequentist approach (provided it exists).

### Examples

- Bernoulli experiment:  $\pi_J(\theta) \propto \frac{1}{\sqrt{p(1-p)}}$ ,  $p \in (0, 1)$ : the prior is  $\text{Beta}(\frac{1}{2}, \frac{1}{2})$
- Gaussian experiment:  $\pi_J(\theta) \propto 1$ ,  $\theta \in \mathbb{R}$ , is an improper prior

## Bayesian Statistics

### Introduction to Bayesian Statistics

#### Prior and Posterior

- Consider a probability distribution on a parameter space  $\Theta$  with some PDF  $\pi(\cdot)$ : the **prior distribution**.



Jeffreys prior satisfies a **reparametrization invariance principle**: If  $\eta$  is a reparametrization of  $\theta$  (i.e.,  $\eta = \phi(\theta)$  for some one-to-one map  $\phi$ ), then the PDF  $\tilde{\pi}(\cdot)$  of  $\eta$  satisfies:

$$\tilde{\pi}(\eta) \propto \sqrt{\det \tilde{I}(\eta)},$$

where  $\tilde{I}(\eta)$  is the Fisher information of the statistical model parametrized by  $\eta$  instead of  $\theta$ .

**Bayesian confidence regions** For  $\alpha \in (0, 1)$ , a Bayesian confidence region with level  $\alpha$  is a random subset  $\mathcal{R}$  of the parameter space  $\Theta$ , which depends on the sample  $X_1, \dots, X_n$ , such that

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

Note that  $\mathcal{R}$  depends on the prior  $\pi(\cdot)$ .

*Bayesian confidence region* and *confidence interval* are two **distinct** notions.

**Bayesian estimation**

- **Posterior mean**:  $\hat{\theta}^{(\pi)} = \int_{\Theta} \theta \pi(\theta | X_1, \dots, X_n) d\theta$
  - **MAP (maximum a posteriori)**:  $\hat{\theta}^{\text{MAP}} = \underset{\theta \in \Theta}{\operatorname{argmax}} \pi(\theta | X_1, \dots, X_n)$
- It is the point that maximizes the posterior distribution, provided it is unique.

## Linear Regression

**Modeling Assumptions**  $(X_i, Y_i), i = 1, \dots, n$ , are i.i.d. from some *unknown joint distribution*  $\mathbb{P}$ .  $\mathbb{P}$  can be described entirely by (assuming all exist):

- either a joint PDF  $h(x, y)$
- the marginal density of  $X$ ,  $h(x) = \int h(x, y) dy$  **and** the conditional density

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

$h(y|x)$  answers all our questions. It contains all the information about  $Y$  given  $X$ .

**Partial Modeling** We can also describe the distribution only partially, e.g. using

- the expectation of  $Y$ :  $\mathbb{E}[Y]$
- the conditional expectation of  $Y$  given  $X = x$ :  $\mathbb{E}[X = x]$ . The function

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

is called **regression function**.

- other possibilities:
  - the conditional median:  $m(x)$  such that

$$\int_{-\infty}^{m(x)} h(y|x) dy = \frac{1}{2}$$

- conditional quantiles
- conditional variance (not information about location)

**Linear Regression** We focus on modeling the regression function

$$f(x) = \mathbb{E}[Y | X = x].$$

Restrict to *simple* functions. The simplest is

$$f(x) = a + bx \quad \text{linear (or affine) function}$$

**Probabilistic Analysis** Let  $X$  and  $Y$  be two r.v. (not necessarily independent) with two moments and such that  $\text{Var}(X) > 0$ . The theoretical linear regression of  $Y$  on  $X$  is the line  $x \mapsto a^* + b^* x$ , where

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

which gives

$$a^* = \frac{\text{Cov}(X, Y)}{\text{Var}(X)}$$

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

**Noise** The points are not exactly on the line  $x \mapsto a^* + b^* x$  if  $\text{Var}(Y | X = x) > 0$ . The random variable  $\varepsilon = Y - (a^* + b^* X)$  is called **noise** and satisfies

$$Y = a^* + b^* X + \varepsilon,$$

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

**Statistical Problem** In practice,  $a^*, b^*$  need to be estimated from data.

**Least Squares** The **least squares estimator** (LSE) of  $(a, b)$  is the minimizer of the sum of squared errors:

$$\sum_{i=1}^n (Y_i - a - bX_i)^2.$$

Then,

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

## Multivariate Regression

We have a vector of explanatory variables or **covariates**:

$$\mathbf{X}_i = \begin{pmatrix} X_i^{(1)} \\ \vdots \\ X_i^{(p)} \end{pmatrix} \in \mathbb{R}^p.$$

The **response** or **dependent variable** is  $Y_i$  with

$$Y_i = \mathbf{X}_i^\top \boldsymbol{\beta}^* + \varepsilon_i, \quad i = 1, \dots, n$$

and  $\boldsymbol{\beta}_1^*$  is called the **intercept**.

**Least Squares Estimator** The least squares estimator of  $\boldsymbol{\beta}^*$  is the minimizer of the sum of squared errors

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \sum_{i=1}^n (Y_i - \mathbf{X}_i^\top \boldsymbol{\beta})^2$$

**LSE in Matrix Form**

- Let  $\mathbf{Y} = (Y_1, \dots, Y_n)^\top \in \mathbb{R}^n$ .
- Let  $\mathbb{X}$  be the  $n \times p$  matrix whose rows are  $\mathbf{X}_1^\top, \dots, \mathbf{X}_n^\top$ .  $\mathbb{X}$  is called the **design matrix**.
- Let  $\boldsymbol{\varepsilon} = (\varepsilon_1, \dots, \varepsilon_n)^\top \in \mathbb{R}^n$ , the unobserved noise. Then,

$$\mathbf{Y} = \mathbb{X} \boldsymbol{\beta}^* + \boldsymbol{\varepsilon}, \quad \boldsymbol{\beta}^* \text{ unknown.}$$

- The LSE  $\hat{\boldsymbol{\beta}}$  satisfies

$$\hat{\boldsymbol{\beta}} = \underset{\boldsymbol{\beta} \in \mathbb{R}^p}{\operatorname{argmin}} \|\mathbf{Y} - \mathbb{X} \boldsymbol{\beta}\|_2^2.$$

**Closed Form Solution** Assume that  $\text{rank}(\mathbb{X}) = p$ . Then,

$$\hat{\boldsymbol{\beta}} = (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top \mathbf{Y}.$$

**Geometric Interpretation of the LSE**  $\mathbb{X} \hat{\boldsymbol{\beta}}$  is the orthogonal projection of  $\mathbf{Y}$  onto the subspace spanned by the columns of  $\mathbb{X}$ :

$$\mathbb{X} \hat{\boldsymbol{\beta}} = P \mathbf{Y},$$

where  $P = \mathbb{X} (\mathbb{X}^\top \mathbb{X})^{-1} \mathbb{X}^\top$ .

**Statistical Inference** To make inference, we need more assumptions.

- The design matrix  $\mathbb{X}$  is deterministic and  $\text{rank}(\mathbb{X}) = p$ .
- The model is **homoscedastic**:  $\varepsilon_1, \dots, \varepsilon_n$  are i.i.d.
- The noise vector  $\boldsymbol{\varepsilon}$  is Gaussian:

$$\boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 \mathbb{I}_n)$$

for some known or unknown  $\sigma^2 > 0$ .

**Properties of LSE**

- LSE = MSE
- Distribution of  $\hat{\boldsymbol{\beta}}$ :

$$\hat{\boldsymbol{\beta}} \sim \mathcal{N}_p(\boldsymbol{\beta}^*, \sigma^2 (\mathbb{X}^\top \mathbb{X})^{-1})$$

- Quadratic Risk of  $\hat{\boldsymbol{\beta}}$ :

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

- Prediction Error:

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

- Unbiased estimator of  $\sigma^2$ :

$$\hat{\sigma}^2 = \frac{\|\mathbf{Y} - \mathbb{X} \hat{\boldsymbol{\beta}}\|_2^2}{n - p} = \frac{1}{n - p} \sum_{i=1}^n \hat{\varepsilon}_i^2$$

**Significance Tests**

- Test whether the  $j^{\text{th}}$  explanatory variable is significant in the linear regression.
- $H_0 : \beta_j = 0$  v.s.  $H_1 : \beta \neq 0$
- If  $\gamma_j$  ( $\gamma_j > 0$ ) is the  $j^{\text{th}}$  diagonal coefficient of  $(\mathbb{X}^\top \mathbb{X})^{-1}$ :

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

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

- Test with non-asymptotic level  $\alpha \in (0, 1)$ :

$$R_{j, \alpha} = \left\{ \left| T_n^{(j)} \right| > q_{\frac{\alpha}{2}}(t_{n-p}) \right\}$$

where  $q_{\frac{\alpha}{2}}(t_{n-p})$  is the  $(1 - \frac{\alpha}{2})$ -quantile of  $t_{n-p}$ .

**Bonferroni's test** Test whether a **group** of explanatory variables is significant in the linear regression.

- $H_0 : \beta_j = 0 \forall j \in S$  v.s.  $H_1 : \exists j \in S, \beta_j \neq 0$  where  $S \subseteq \{1, \dots, p\}$ .
- Bonferroni's test:

$$R_{S, \alpha} = \bigcup_{j \in S} R_{j, \frac{\alpha}{k}}, \quad \text{where } k = |S|$$

# Generalized Linear Model

**Generalization** A generalized linear model (GLM) generalizes normal linear regression models in the following directions:

1. **Random component:**  $Y|X = x \sim$  some distribution
2. **Regression function:**

$$g(\mu(x)) = x^\top \beta$$
where  $g$  is called **link function** and  $\mu(x) = \mathbb{E}[Y|X = x]$  is the **regression function**.

## Exponential Family

A family of distribution  $\{\mathbb{P}_\theta : \theta \in \Theta\}$ ,  $\Theta \subset \mathbb{R}^k$  is said to be a  **$k$ -parameter exponential family** on  $\mathbb{R}^q$ , if there exist real-valued functions

- $\eta_1, \dots, \eta_k$  and  $B(\theta)$
- $T_1, \dots, T_k$ , and  $h(y) \in \mathbb{R}^q$

such that the density function of  $\mathbb{P}_\theta$  can be written as

$$f_\theta(y) = \exp \left[ \sum_{i=1}^k \eta_i(\theta) T_i(y) - B(\theta) \right] h(y)$$

**Examples of discrete distributions** The following distributions for **discrete** exponential families of distributions with PMF:

- Bernoulli ( $p$ ):  $p^y (1 - p)^{1-y}$ ,  $y \in \{0, 1\}$
- Poisson ( $\lambda$ ):  $\frac{\lambda^y}{y!} e^{-\lambda}$ ,  $y = 0, 1, \dots$

**Examples of continuous distributions** The following distributions form **continuous** exponential families of distributions with PDF:

- Gamma ( $a, b$ ):  $\frac{1}{\Gamma(a)b^a} y^{a-1} e^{-\frac{y}{b}}$
- Inverse Gamma ( $\alpha, \beta$ ):  $\frac{\beta^\alpha}{\Gamma(\alpha)} y^{-\alpha-1} e^{-\frac{\beta}{y}}$
- Inverse Gaussian ( $\mu, \sigma^2$ ):  $\sqrt{\frac{\sigma^2}{2\pi y^3}} \exp \left( -\frac{\sigma^2(y - \mu)^2}{2\mu^2 y} \right)$

**One-parameter Canonical Exponential Family**

$$f_\theta(y) = \exp \left( \frac{y\theta - b(\theta)}{\phi} + c(y, \phi) \right)$$

for some known functions  $b(\theta)$  and  $c(y, \phi)$ .

- If  $\phi$  is known, this is a one-parameter exponential family with  $\theta$  being the **canonical parameter**.
- If  $\phi$  is unknown, this may / may not be a two-parameter exponential family.
- $\phi$  is called **dispersion parameter**.

**Expected value** Note that

$$\ell(\theta) = \frac{Y\theta - b(\theta)}{\phi} + c(Y; \phi),$$

which leads to

$$\mathbb{E}[Y] = b'(\theta).$$

**Variance**

$$\text{Var}(Y) = b''(\theta) \cdot \phi$$

In GLM, we have  $Y|X = x \sim$  distribution in exponential family. Then,

$$\mathbb{E}[Y|X = x] = f(X^\top \beta)$$

**Link function**  $\beta$  is the parameter of interest. A **link function**  $g$  relates the linear predictor  $X^\top \beta$  to the mean parameter  $\mu$ ,

$$X^\top \beta = g(\mu) = g(\mu(X)).$$

$g$  is required to be monotone increasing and differentiable

$$\mu = g^{-1}(X^\top \beta)$$

**Canonical Link** The function  $g$  that links the mean  $\mu$  to the canonical parameter  $\theta$  is called **canonical link**:

$$g(\mu) = \theta.$$

Since  $\mu = b'(\theta)$ , the canonical link is given by

$$g(\mu) = (b')^{-1}(\mu).$$

If  $\phi > 0$ , the canonical link function is strictly increasing.

**Example** Bernoulli distribution

$$\begin{aligned} p^y (1 - p)^{1-y} &= \exp \left( y \log \left( \frac{p}{1 - p} \right) + \log(1 - p) \right) \\ &= \exp \left( y\theta - \log(1 + e^\theta) \right) \end{aligned}$$

Hence,  $\theta = \log \left( \frac{p}{1 - p} \right)$  and  $b(\theta) = \log(1 + e^\theta)$ .

$$b'(\theta) = \frac{e^\theta}{1 + e^\theta} = \mu \iff \theta = \log \left( \frac{\mu}{1 - \mu} \right)$$

The canonical link for the Bernoulli distribution is the **logit link**.

## Model and Notation

Let  $(X_i, Y_i) \in \mathbb{R}^p \times \mathbb{R}$ ,  $i = 1, \dots, n$  be independent random pairs such that the conditional distribution of  $Y_i$  given  $X_i = x_i$  has density in the canonical exponential family:

$$f_{\theta_i}(y_i) = \exp \left[ \frac{y_i \theta_i - b(\theta_i)}{\phi} + c(y_i, \phi) \right]$$

**Back to  $\beta$ :** Given a link function  $g$ , note the following relationship between  $\beta$  and  $\theta$ :

$$\theta_i = (b')^{-1}(\mu_i) = (b')^{-1} \left( g^{-1}(X_i^\top \beta) \right) \equiv h(X_i^\top \beta)$$

where  $h$  is defined as

$$h = (b')^{-1} \circ g^{-1} = (g \circ b')^{-1}.$$

If  $g$  is the *canonical link function*,  $h$  is the **identity**  $g = (b')^{-1}$ .

**Log-likelihood** The log-likelihood is given by

$$\begin{aligned} \ell_n(\mathbf{Y}, \mathbb{X}, \beta) &= \sum_i \frac{Y_i \theta_i - b(\theta_i)}{\phi} + \text{constant} \\ &= \sum_i \frac{Y_i h(X_i^\top \beta) - b(h(X_i^\top \beta))}{\phi} + \text{constant} \end{aligned}$$

When we use the *canonical link function*, we obtain the expression

$$\ell_n(\mathbf{Y}, \mathbb{X}, \beta) = \sum_i \frac{Y_i X_i^\top \beta - b(X_i^\top \beta)}{\phi} + \text{constant}$$

**Strict concavity** The log-likelihood  $\ell(\theta)$  is **strictly concave** (if  $\text{rank}(\mathbb{X}) = p$ ) using the canonical function when  $\phi > 0$ . As a consequence, the maximum likelihood estimator is unique.

On the other hand, if another parametrization is used, the likelihood function may not be strictly concaving leading to *several local maxima*.

# Recommended Resources

- Probability and Statistics (DeGroot and Schervish)
- Mathematical Statistics and Data Analysis (Rice)
- Fundamentals of Statistics [Lecture Slides] (<http://www.edx.org>)

Please share this cheatsheet with friends!