

# Classical Statistical Inference

Xuening Bai (白雪宁)

Institute for Advanced Study (IASTU) &  
Department of Astronomy (DoA)



清華大學

Tsinghua University

Nov. 12, 2024

# Parametric vs non-parametric models

A **statistic model** is a set of distributions, among them there are:

A **parametric model** is one that can be parameterized by a finite number of parameters  $\theta$ .

Available models are restricted within **parameter space**  $\Theta$ .

There are also cases where only one or some parameters within  $\Theta$  needs to be estimated.

A **nonparametric model** one that can not be parameterized by a finite number of parameters.

Directly estimate the CDF/PDF without parameterization.

Estimate some **statistical functional** of the CDF without parameterization.

# Frequentist vs Bayesian inference

Frequentists consider model parameters  $\theta$  as an unknown constant. Given  $\theta$ , observed data are a realization (sampling) of some random variables.

This lecture.

Bayesians consider  $\theta$  as some random variables, whereas data are considered as known constants.

Next lecture

# Topics to address

Most classical inference problems can be identified into 3 categories:

(Point) estimation:

Provide the “best guess” of some quantity of interest.

Confidence interval (or *confidence set* in multi-D):

Given data, range of parameters that yield this data at certain probability.

Hypothesis testing:

Given a theory (*null hypothesis*), whether data provide sufficient evidence to *reject* this theory.

# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Fundamental concepts

Let  $X_1, \dots, X_n \sim F$  (the CDF) be an IID sample. A **point estimator** for some parameter  $\theta$  is usually of the form:

$$\hat{\theta}_n = g(X_1, \dots, X_n)$$

The distribution of  $\hat{\theta}_n$  is called **sampling distribution**.

The **bias** is defined as:  $\text{bias}(\hat{\theta}_n) = E(\hat{\theta}_n) - \theta$

An estimator is said to be **unbiased** if its bias is 0 for any  $n$ .

An estimator is said to be **consistent** if  $\hat{\theta}_n$  converges to  $\theta$  in probability.

The standard deviation of  $\hat{\theta}_n$  is called the **standard error**.

$$\text{se}_{\hat{\theta}_n} = \sqrt{\text{Var}(\hat{\theta}_n)}$$

# Example:

Suppose  $X_1, \dots, X_n$  are IIDs satisfying a normal distribution  $X \sim \mathcal{N}(\mu, \sigma)$ .

From the previous lecture:

Sample mean:  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$  unbiased estimate of  $E(X)$

Variance of mean:  $\text{Var}(\bar{X}) = \frac{\text{Var}(X)}{n}$

We can thus give an unbiased estimate of  $\mu$ :  $\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n X_i$

Similarly, we have an unbiased estimate of  $\sigma^2$ :  $\hat{\sigma}_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \hat{\mu}_n)^2$

Standard error of  $\hat{\mu}_n$ :  $\text{se}_{\hat{\mu}_n} = \frac{\sigma}{\sqrt{n}}$ , and can be estimated as  $\hat{\text{se}}_{\hat{\mu}_n} = \sqrt{\frac{\hat{\sigma}_n^2}{n}}$

# Fundamental concepts

The efficacy of a point estimate can be assessed by the **mean squared error**:

$$\text{MSE} = E(\hat{\theta}_n - \theta)^2$$

It is straightforward to show

$$\text{MSE} = \text{bias}^2(\hat{\theta}_n) + \text{Var}(\hat{\theta}_n)$$

A **(1- $\alpha$ ) confidence interval** for a parameter  $\theta$  is an interval  $C_n=(a,b)$ , where  $a, b$  are functions of  $X_1, \dots, X_n$  such that

$$P(\theta \in C_n) \geq 1 - \alpha$$

It is important to note that  $C_n$  is random but  $\theta$  is fixed (frequentists' view).

If the sample distribution is **asymptotically normal**, then confidence interval can be approximately determined from  $\text{se}_{\hat{\theta}_n}$  based on the normal distribution.



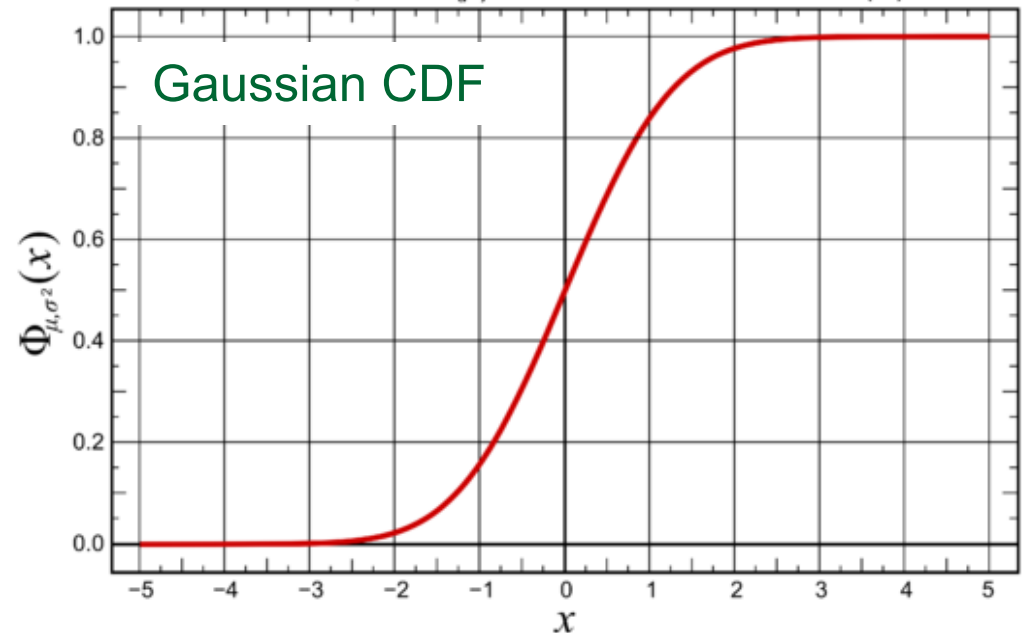
# Confidence interval for a Gaussian distribution

It often occurs that the sample distribution of  $\hat{\theta}_n$  is a Gaussian:

$$\hat{\theta}_n \sim \mathcal{N}(\theta, \hat{\text{se}}_{\hat{\theta}_n}^2)$$

Let  $\Phi$  be the CDF of a standard normal distribution. We define :

$$z_{\alpha/2} \equiv \Phi^{-1}(1 - \alpha/2)$$



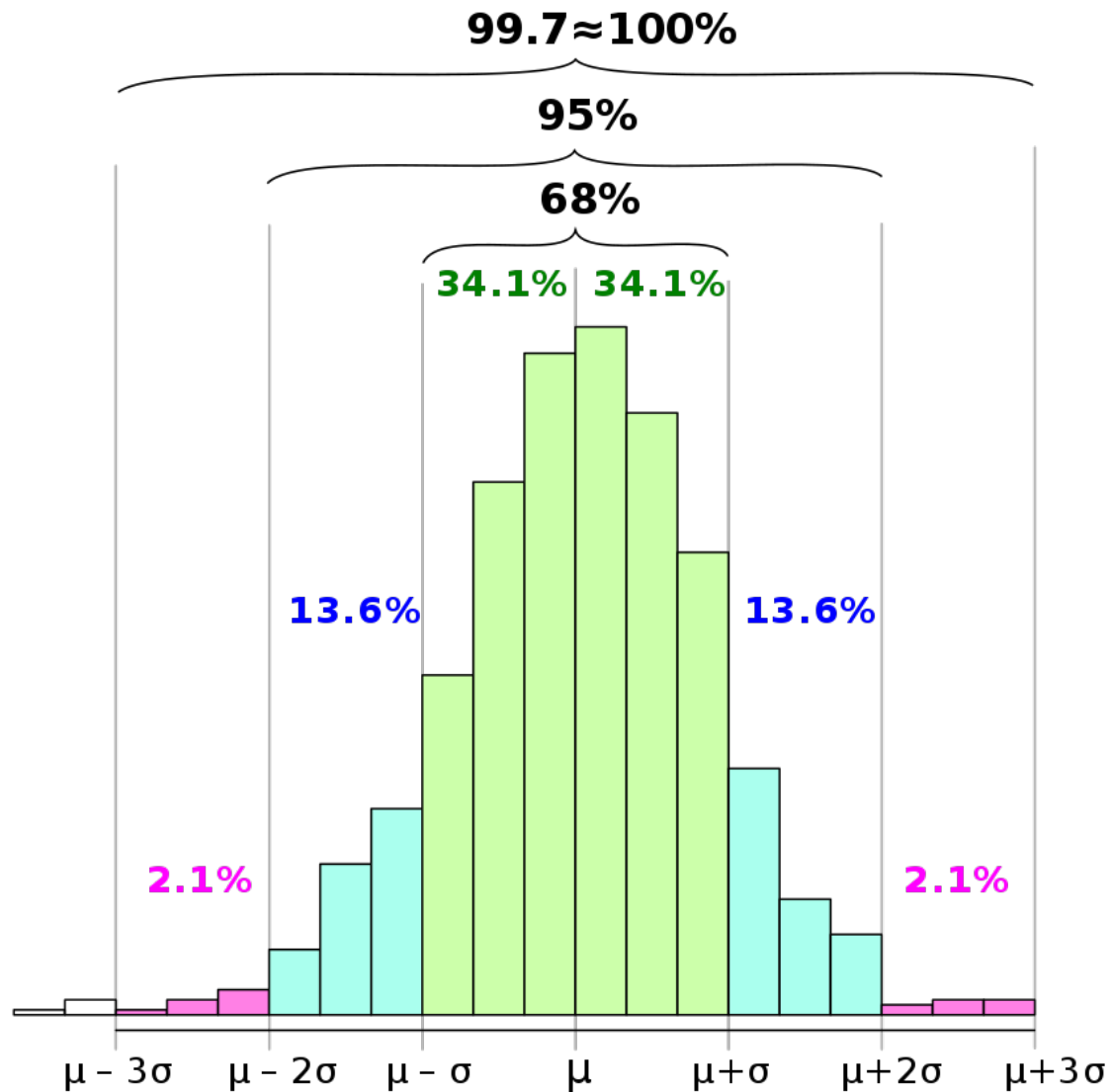
Then  $P(-z_{\alpha/2} < Z < z_{\alpha/2}) = 1 - \alpha$  for  $Z \sim \mathcal{N}(0,1)$ .

The confidence interval for  $\hat{\theta}_n$  is then given by

$$C_n = (\hat{\theta}_n - z_{\alpha/2} \hat{\text{se}}_{\hat{\theta}_n}, \hat{\theta}_n + z_{\alpha/2} \hat{\text{se}}_{\hat{\theta}_n})$$

For 95% confidence level,  $\alpha=0.05$ , we have  $z_{\alpha/2} = 1.96 \approx 2$ .

# Confidence interval for a Gaussian distribution



$1\sigma$ : 68%

$2\sigma$ : 95%

$3\sigma$ : 99.7%

$5\sigma$ : 99.999943%

# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Non-parametric estimate of the CDF

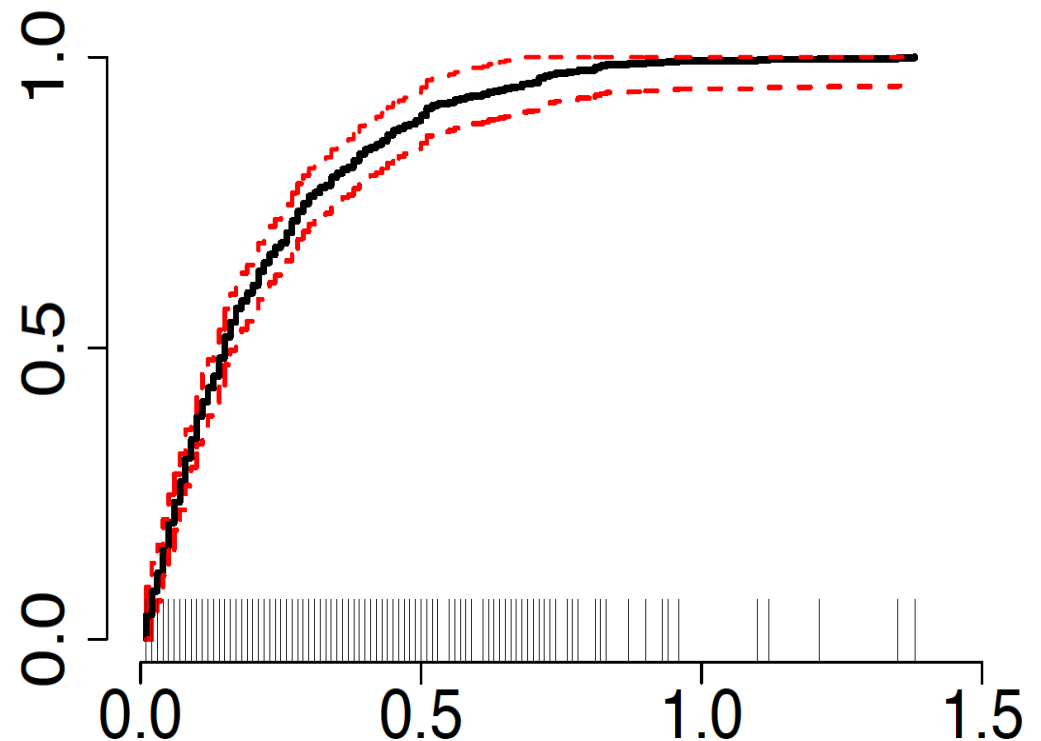
Let  $X_1, X_2, \dots, X_n \sim F$  be an IID sample, where  $F(x)$  is the CDF, and they take the values  $x_1, x_2, \dots, x_n$ .

We can estimate  $F$  via the so-called [empirical distribution function](#) (EDF):

$$\hat{F}_n(x) = \frac{\sum_{i=1}^n I(x_i \leq x)}{n}$$

where

$$I(x_i \leq x) = \begin{cases} 1 & \text{if } x_i \leq x \\ 0 & \text{if } x_i > x \end{cases}$$



# Non-parametric estimate of the CDF

One can show that

$$E(\hat{F}_n(x)) = F(x) , \quad \text{Var}(\hat{F}_n(x)) = \frac{F(x)(1 - F(x))}{n}$$

and that the EDF *almost surely* converges to the CDF as  $n \rightarrow \infty$ .

There is the [Dvoretzky-Kiefer-Wolfowitz \(DKW\) inequality](#), which allows one to find the nonparametric [confidence interval](#):

$$\begin{aligned} L(x) &= \max\{\hat{F}_n(x) - \epsilon_n, 0\} \\ U(x) &= \max\{\hat{F}_n(x) + \epsilon_n, 1\} \end{aligned} \quad \text{where} \quad \epsilon_n = \sqrt{\frac{1}{2n} \log \left( \frac{2}{\alpha} \right)}$$

so that

$$P(L(x) \leq F(x) \leq U(x)) \geq 1 - \alpha \quad (\text{for all } x)$$

# Plug-in estimator

With the EDF, one can estimate any statistical function of  $F$  as follows.

For  $\theta = T(F)$ , the **plug-in estimator** is defined as  $\hat{\theta} = T(\hat{F})$ .

Examples:

The mean:  $\mu = T(F) = \int x dF(x) \Rightarrow \hat{\mu} = \int x d\hat{F}_n(x) = \bar{x}_n$

The variance:

$$\sigma^2 = T(F) = \text{Var}_X = \int x^2 dF(x) - \left( \int x dF(x) \right)^2$$
$$\hat{\sigma}^2 = \int x^2 d\hat{F}(x) - \left( \int x d\hat{F}(x) \right)^2 = \frac{1}{n} \sum_{i=1}^n x_i^2 - \bar{x}_n^2 = \boxed{\frac{1}{n}} \sum_{i=1}^n (x_i - \bar{x}_n)^2$$

# Bootstrap

A **resampling** method for **estimating standard errors and confidence** intervals.

We have a data set  $\{x_i\}$  coming from some distribution with its CDF being  $F(x)$ .

We do not know  $F(x)$ , but we can use the EDF,  $\hat{F}_n(x)$  that is drawn from  $n$  IID samples, as an approximation. Note, the PDF of this function is essentially

$$\hat{f}_n(x) = \frac{1}{N} \sum_{i=1}^n \delta(x - x_i)$$

Suppose we are interested in some statistic  $T(F)$ , and with a finite sample it is given by  $T_n = g(X_1, \dots, X_n)$ , e.g., based on a plug-in estimator.

Further question: what is the variance and confidence interval of our estimate?

# Bootstrap

Key insight from bootstrap:

Drawing an observation from  $\hat{F}_n(x)$  is equivalent to drawing one point at random from the original data set.

draw random  
sample

compute  
statistic

Real world:  $F \implies X_1, \dots, X_n \implies T_n = g(X_1, \dots, X_n)$

Bootstrap:  $\hat{F}_n \implies X_1^*, \dots, X_n^* \implies T_n^* = g(X_1^*, \dots, X_n^*)$

In bootstrap, this can be done for large number of times, which allows one to compute the variance and confidence intervals.



# Bootstrap: example

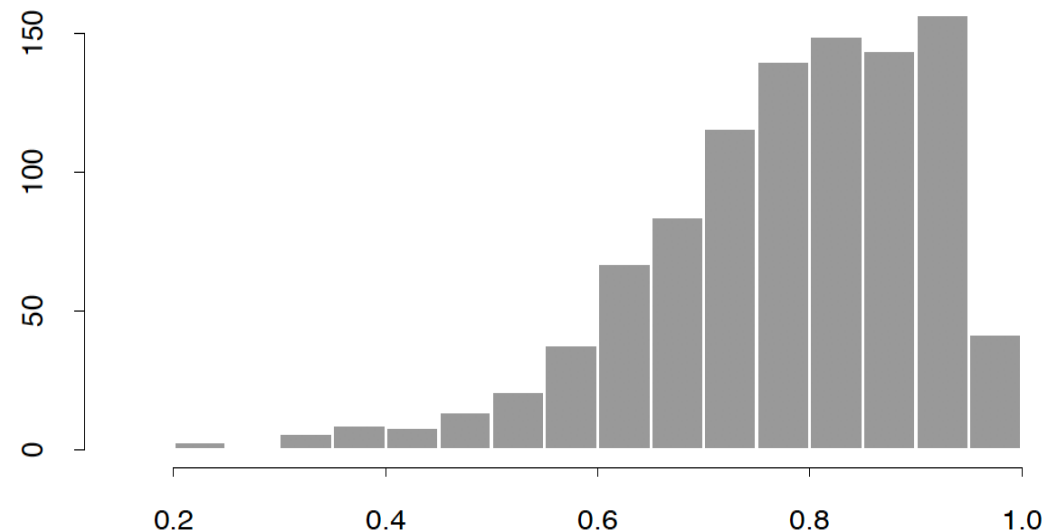
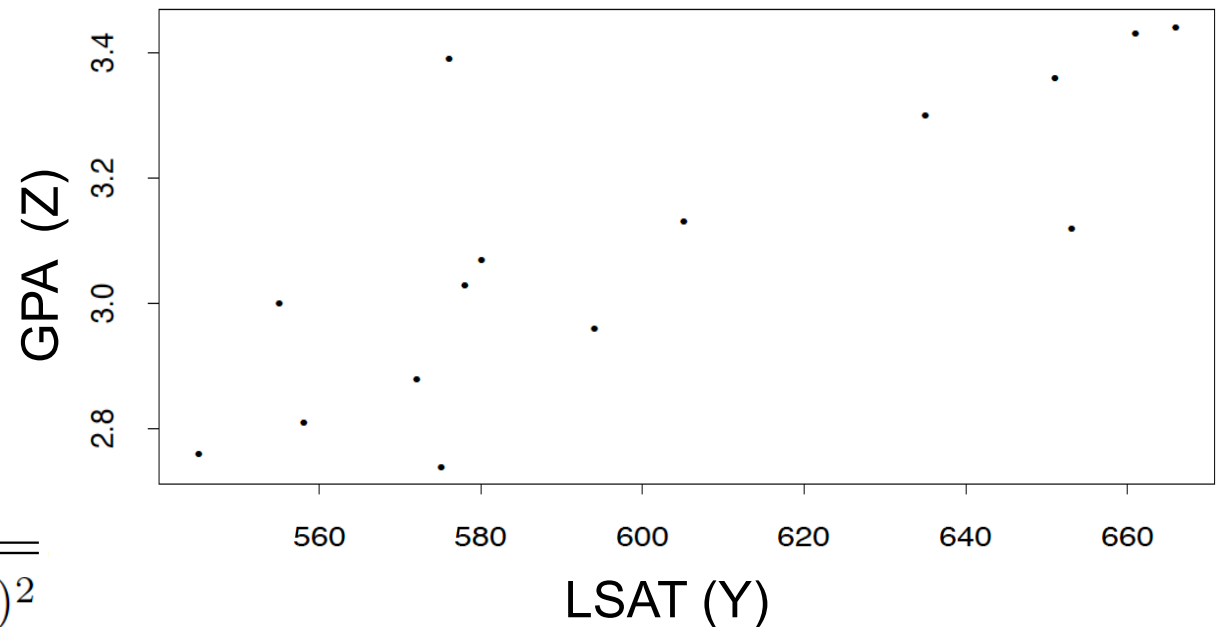
Given this data, what is the correlation between GPA and LSAT scores?

First computes the sample correlation coefficients:

$$\hat{\theta} = \frac{\sum_i (Y_i - \bar{Y})(Z_i - \bar{Z})}{\sqrt{\sum_i (Y_i - \bar{Y})^2 \sum_i (Z_i - \bar{Z})^2}}$$
$$= 0.776$$

Then conduct bootstrap for 1000 times, and each time one obtains a  $\hat{\theta}_B$ .

Confidence interval at 95% is found from this histogram to be (0.46, 0.96).



# Jackknife

Another resampling method in similar spirit to bootstrap.

Instead of drawing a data set of the same size as original, now we drop one or more observations at a time to compute the statistic of interest.

Let  $T_n = g(X_1, \dots, X_n)$  be a statistic, and  $T_{-i}$  be the statistic with  $i$ th observation removed. Further define:

$$\bar{T}_n = \frac{1}{n} \sum_{i=1}^n T_{-i}$$

Then the jackknife estimate of  $\text{Var}(T_n)$  is

$$V_{\text{jack}} = \frac{n-1}{n} \sum_{i=1}^n (T_{-i} - \bar{T}_n)^2$$

It is easier to implement than bootstrap, while the latter provides more reliable confidence intervals.

# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Prelude: method of moments

Express population moments as a function of parameters of interest.

$$\alpha_j = E(X^j) = g(\theta_1, \dots, \theta_k) , \quad j = 1, \dots, k$$

Approximate them with sample moments to and solve for the parameters.

$$\hat{\alpha}_j = \frac{1}{n} \sum_{i=1}^n X_i^j$$

Example: Normal distribution with  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ , then

$$\begin{aligned} \hat{\mu} &= \frac{1}{n} \sum_{i=1}^n X_i \\ \hat{\sigma}^2 + \hat{\mu}^2 &= \frac{1}{n} \sum_{i=1}^n X_i^2. \end{aligned}$$

Easy to calculate, but results are often not optimal.

Can be used as starting values for other methods that requires iteration.

# Maximum likelihood estimation (MLE)

The most common method for parameter estimation in a [parametric model](#).

Let  $X_1, \dots, X_n$  be IID with PDF  $f(x; \theta)$ , with observed values being  $x_1, \dots, x_n$ .

The [likelihood function](#) is defined as:

$$\mathcal{L}_n(\theta) = \prod_{i=1}^n f(x_i; \theta)$$

Often times, we use the [log-likelihood function](#):

$$l_n(\theta) \equiv \log \mathcal{L}_n(\theta) = \sum_{i=1}^n \log f(x_i; \theta)$$

The [maximum likelihood estimator, MLE](#), is the value of  $\theta$ , denoted by  $\hat{\theta}_n$ , that maximizes the (log) likelihood function.

# Example: Bernoulli distribution

Suppose  $X_1, \dots, X_n \sim \text{Bernoulli}(p)$ . You want to estimate  $p$ . The **probability function** is

$$f(x; p) = p^x (1 - p)^{1-x}, \quad (x = 0, 1) \quad (p \text{ is unknown parameter})$$

Given the data, the likelihood function is

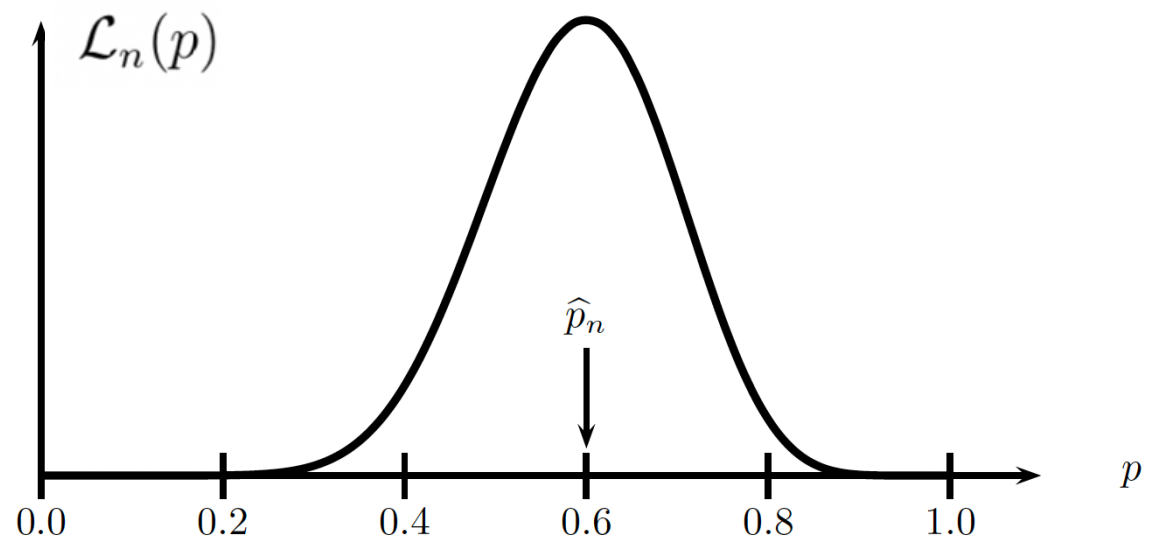
$$\mathcal{L}_n(p) = \prod_{i=1}^n p^{x_i} (1 - p)^{1-x_i} = p^m (1 - p)^{n-m} \quad \text{where } m = x_1 + \dots + x_n.$$

$$l_n(p) = m \log p + (n - m) \log(1 - p)$$

Maximizing this function:

$$\frac{\partial l_n}{\partial p} = \frac{m}{p} - \frac{n - m}{1 - p} = 0$$

$$\Rightarrow \hat{p}_n = \frac{m}{n}$$



# Example: normal distribution

Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$ . The likelihood function then reads (ignoring constants)

$$\mathcal{L}_n(\mu, \sigma) = \prod_i \frac{1}{\sigma} \exp \left[ -\frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$l_n(\mu, \sigma) = -n \log \sigma - \sum_{i=1}^n \frac{(x_i - \mu)^2}{2\sigma^2}$$

Taking derivatives on  $\mu$  and  $\sigma$ , one obtains:

$$\frac{\partial l_n}{\partial \mu} = \sum_{i=1}^n \frac{1}{\sigma^2} (x_i - \mu) = 0 \quad \Rightarrow \quad \hat{\mu}_n = \frac{1}{n} \sum_{i=1}^n x_i$$

$$\frac{\partial l_n}{\partial \sigma} = -\frac{n}{\sigma} + \sum_{i=1}^n \frac{(x_i - \mu)^2}{\sigma^3} \quad \Rightarrow \quad \hat{\sigma}_n^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \hat{\mu}_n)^2$$

# Properties of the MLE

Under certain regularity conditions:

- The MLE is **consistent**:

Converge to the true parameter in probability.

- The MLE is **equivariant**:

If  $\hat{\theta}_n$  is the MLE of  $\theta$ , then  $g(\hat{\theta}_n)$  is the MLE of  $g(\theta)$ .

- The MLE is **asymptotically normal**:

This is to say,  $\frac{\hat{\theta}_n - \theta}{\widehat{\text{se}}_{\hat{\theta}_n}} \sim \mathcal{N}(0, 1)$  at large  $n$ .

- The MLE is **asymptotically optimal/efficient**:

Among all well-behaved estimators, the MLE achieves the smallest possible variance (known as Cramer-Rao bound), at least for large samples.



# Fisher information

First introduce the **score function**:

$$s(X; \theta) = \frac{\partial \log f(X; \theta)}{\partial \theta} \quad \text{It can be seen that } E[s(X; \theta)] = 0$$

The **expected Fisher information** (for IID) is defined as

$$I_n(\theta) = \text{Var}_\theta \left( \sum_{i=1}^n s(X_i; \theta) \right) = \sum_{i=1}^n \text{Var}_\theta [s(X_i; \theta)] = nI(\theta)$$

It can be shown that

$$I(\theta) = \text{Var}_\theta [s(X; \theta)] = E_\theta [s^2(X; \theta)]$$

$$I(\theta) = -E_\theta (-s'(X; \theta)) = - \int \left( \frac{\partial^2 \log f(x; \theta)}{\partial \theta^2} \right) f(x; \theta) dx$$

# MLE confidence intervals

The Fisher information can be used to estimate the **standard deviation** and prove **asymptotic normality** (based on the central limit theorem).

Let  $\hat{\theta}_n$  be the MLE, then asymptotically, its **standard error** is

$$\text{se}_{\hat{\theta}_n}^2 \approx \frac{1}{I_n(\theta)} = \frac{1}{nI(\theta)}$$

In practice, this **standard error** can be estimated as

$$\hat{\text{se}}_{\hat{\theta}_n}^2 \approx \frac{1}{I_n(\hat{\theta}_n)} = \frac{1}{nI(\hat{\theta}_n)}$$

The **asymptotic distribution of**  $\hat{\theta}_n$  is given by  $\hat{\theta}_n \sim \mathcal{N}(\theta, \text{se}_{\hat{\theta}_n}^2) \sim \mathcal{N}(\theta, \hat{\text{se}}_{\hat{\theta}_n}^2)$

This allows one to estimate the approximate confidence interval as

$$C_n \approx (\hat{\theta}_n - z_{\alpha/2} \hat{\text{se}}_{\hat{\theta}_n}, \hat{\theta}_n + z_{\alpha/2} \hat{\text{se}}_{\hat{\theta}_n})$$

# Example: Bernoulli distribution

Recall that the  
probability function:

$$f(x; p) = p^x (1 - p)^{1-x}, \quad (x = 0, 1)$$

$$\log f(x; p) = x \log p + (1 - x) \log(1 - p)$$

Now compute the  
score function:

$$s(X; p) = \frac{X}{p} - \frac{1 - X}{1 - p}, \quad -s'(X; p) = \frac{X}{p^2} + \frac{1 - X}{(1 - p)^2}$$

Fisher information:

$$I(p) = E_p(-s'(X; p)) = \frac{p}{p^2} + \frac{1 - p}{(1 - p)^2} = \frac{1}{p(p - 1)}$$

Therefore, the standard error is approximately

$$\widehat{\text{se}}_{\hat{p}_n}^2 = \frac{1}{nI(\hat{p}_n)} = \frac{\hat{p}_n(\hat{p}_n - 1)}{n}$$

# Multiparameter estimation

Let  $\theta = (\theta_1, \dots, \theta_k)$  and let  $\hat{\theta} = (\hat{\theta}_1, \dots, \hat{\theta}_k)$  be the MLE.

The log likelihood function is  $l_n = \sum_{i=1}^n \log f(x_i; \theta)$ , and define

$$H_{jj} \equiv \frac{\partial^2 l_n}{\partial \theta_j^2}, \quad H_{jk} \equiv \frac{\partial^2 l_n}{\partial \theta_j \partial \theta_k}$$

The **Fisher information matrix** is defined as

$$I_n(\theta) = - \begin{bmatrix} E_{\theta}(H_{11}) & E_{\theta}(H_{12}) & \dots & E_{\theta}(H_{1k}) \\ E_{\theta}(H_{21}) & E_{\theta}(H_{22}) & \dots & E_{\theta}(H_{2k}) \\ \vdots & \vdots & \vdots & \vdots \\ E_{\theta}(H_{k1}) & E_{\theta}(H_{k2}) & \dots & E_{\theta}(H_{kk}) \end{bmatrix}$$

Let  $J_n(\theta) = I_n^{-1}(\theta)$  be its inverse.

# Asymptotic normality

Under appropriate regularity conditions, we have

$$\hat{\theta} \sim \mathcal{N}(\theta, J_n(\theta))$$

Also, for individual components, we have

$$\hat{\theta}_j \sim \mathcal{N}(\theta_j, J_{n,jj}(\theta))$$

The errors in different parameters are not necessarily independent. In general, they are correlated, and the covariance is given by

$$\text{Cov}(\hat{\theta}_j, \hat{\theta}_k) = J_{n,jk}(\theta)$$

# Example: normal distribution

For a normal distribution with two parameters  $\mu, \sigma$ , we have

$$l_n(\mu, \sigma) = -n \log \sigma - \sum_{i=1}^n \frac{(X_i - \mu)^2}{2\sigma^2}$$
$$\frac{\partial l_n}{\partial \mu} = \sum_{i=1}^n \frac{1}{\sigma^2} (X_i - \mu), \quad \frac{\partial l_n}{\partial \sigma} = -\frac{n}{\sigma} + \sum_{i=1}^n \frac{(X_i - \mu)^2}{\sigma^3}$$

One can find that at MLE values, the [Fisher information matrix](#) reads:

$$\left. \frac{\partial^2 l_n}{\partial \mu^2} \right|_{\hat{\mu}, \hat{\sigma}} = -\frac{n}{\hat{\sigma}^2}, \quad \left. \frac{\partial^2 l_n}{\partial \sigma^2} \right|_{\hat{\mu}, \hat{\sigma}} = \frac{n}{\hat{\sigma}^2} - 3\frac{n}{\hat{\sigma}^2} = -\frac{2n}{\hat{\sigma}^2}, \quad \left. \frac{\partial^2 l_n}{\partial \mu \partial \sigma} \right|_{\hat{\mu}, \hat{\sigma}} = 0$$
$$I_n(\hat{\mu}, \hat{\sigma}) = \frac{n}{\hat{\sigma}^2} \begin{pmatrix} 1 & 0 \\ 0 & 2 \end{pmatrix} \Rightarrow J_n(\hat{\mu}, \hat{\sigma}) = \frac{\hat{\sigma}^2}{n} \begin{pmatrix} 1 & 0 \\ 0 & 1/2 \end{pmatrix}$$

Therefore, errors for  $\mu, \sigma$  are uncorrelated, and are given by

$$\widehat{\text{se}}(\hat{\mu}) = \frac{\hat{\sigma}}{\sqrt{n}}, \quad \widehat{\text{se}}(\hat{\sigma}) = \frac{\hat{\sigma}}{\sqrt{2n}}$$

# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Hypothesis testing: basic concepts

**Hypothesis testing** is a procedure for comparing observed data with a hypothesis whose plausibility is to be assessed.

Suppose we partition the parameter space into two disjoint sets  $\Theta_0$  and  $\Theta_1$ , and we wish to test:

$$H_0 : \theta \in \Theta_0 \text{ v.s. } H_1 : \theta \in \Theta_1$$

We call  $H_0$  the **null hypothesis**, and  $H_1$  the **alternative hypothesis**.

Three ingredients:

- **Significance level  $\alpha$** : (roughly) the probability of rejecting  $H_0$  when it is correct.
- **Test statistic  $T$** : calculated from data to measure its compatibility with  $H_0$ .
- **A rejection rule**: specify the range of  $T$  to reject  $H_0$ .



# Example: how to pose the problem?

Suppose you made an amazing discovery, e.g., detect a dark matter particle!

Before announcing this exciting result, you don't want to make a fool of yourself.

How would you pose the hypothesis test problem?

Null hypothesis:

This is consistent  
with dark matter.

This is inconsistent  
with dark matter.

Alternative hypothesis:

There is definitely  
NO dark matter.

There is definitely  
dark matter!

Read: data suggest the Null hypothesis is rejected at 5% significance level.

# Hypothesis testing: basic concepts

Let  $X$  be a random variable and  $\mathcal{X}$  be its range.

We test a hypothesis by finding an appropriate subset of outcomes  $W \subset \mathcal{X}$  called the **rejection region** so that:

$$X \in W \implies \text{reject } H_0$$

$$X \notin W \implies \text{retain (do not reject) } H_0$$

Usually, the rejection region  $W$  is of the form

$$W = \left\{ x : T(x) > c \right\}$$

Here,  $T(X)$  is called **test statistic**, and  $c$  is called a **critical value**.

The problem in hypothesis test is to **find an appropriate test statistic  $T$  and an appropriate critical value  $c$**  (equivalently, a rejection region).

# Example: source detection

Suppose your detector is subject to Poisson noise, with **background count rate being  $\lambda$  (known)**. Over time interval  $t$ , total count satisfies the Poisson distribution:

$$f_0(n, t) = \frac{(\lambda t)^n}{n!} e^{-\lambda t}$$

Let  $N$  be the number of photons gathered in an experiment. We want to ask whether there is a signal (which would give an excess).

**Null hypothesis:** there is no signal, so that data is consistent with  $f = f_0$ .

**Test statistic:**  $T(N) = N$  (this is a trivial case: there is no parameter)

The **rejection region** can be defined as:  $W : \{n : T(n) > c\}$

We need to specify a **significance level  $\alpha$** , say 5%. Thus, we require

$$P(X \in W) = P(N > c) = \alpha = 0.05$$

In other words, we should choose  $c$  such that  $\sum_{k=0}^c \frac{(\lambda t)^k}{k!} e^{-\lambda t} \approx 0.95$  .

# Type I/II errors and significance level

The conclusion drawn from hypothesis test can be false, leading to an error.

	Retain Null	Reject Null	
$H_0$ true	✓	type I error	false positive
$H_1$ true	type II error	✓	false negative

We are generally more concerned with type I error.

Define the **power function**:  $\rho_W(\theta) \equiv P(X \in W ; \theta)$

Given  $\theta$ , what is the probability of being rejected?

If  $H_0$  is true, that is,  $\theta \in \Theta_0$ , how likely is it rejected (type-I error)?

Define **significance level**:  $\alpha = \sup_{\theta \in \Theta_0} \rho_W(\theta)$ . (i.e., maximum type I error rate)

# Example: normal distribution

Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sigma$  is known.

We want to test  $H_0: \mu \leq 0$  versus  $H_1: \mu > 0$ , that is,  $\Theta_0 = (-\infty, 0]$ ,  $\Theta_1 = (0, \infty)$

Consider using  $T = \bar{X}$  as a test statistic, and

reject  $H_0$  if  $T > c$  for some  $c$  TBD.

We know that  $T$  should satisfy  $T \sim \mathcal{N}(\mu, \sigma^2/n)$ . That is,

$$Z \equiv \frac{\sqrt{n}(T - \mu)}{\sigma} \sim \mathcal{N}(0, 1)$$

Therefore,

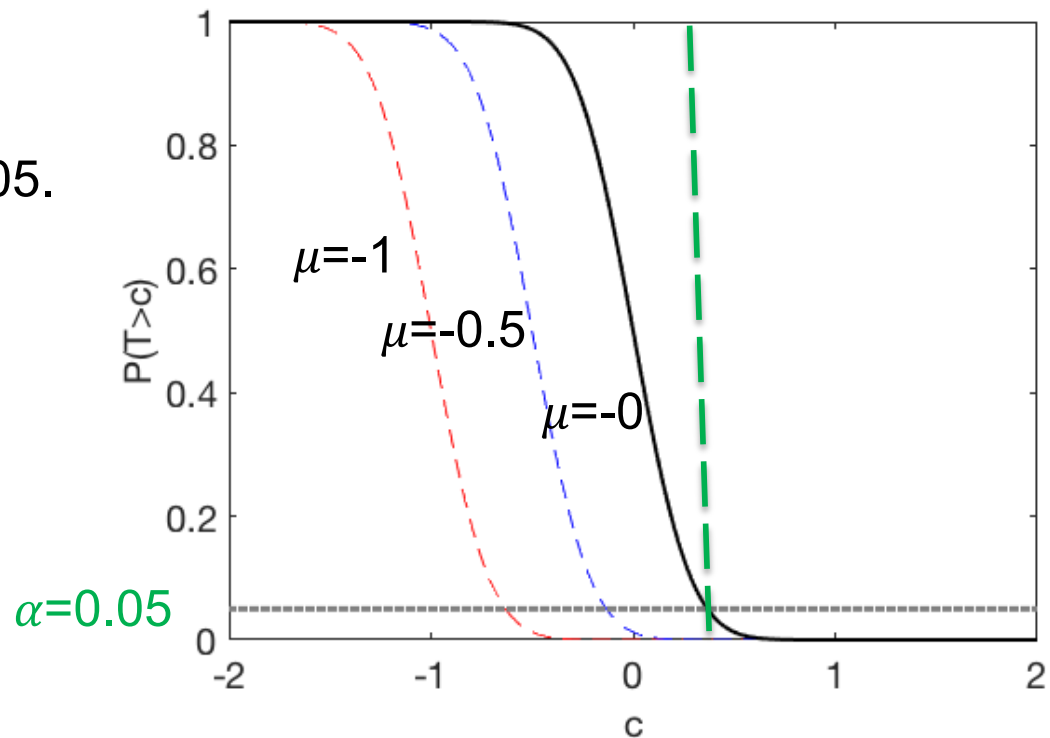
$$P(T > c; \mu) = P\left(Z > \frac{\sqrt{n}(c - \mu)}{\sigma}; \mu\right) = 1 - \Phi\left(\frac{\sqrt{n}(c - \mu)}{\sigma}\right)$$

# Example: normal distribution

We want significance level  $\alpha=0.05$ .

This should apply to all  $\mu \leq 0$ .

Clearly, the case that maximizes type-I error corresponds to  $\mu=0$ .



Therefore, we should determine  $c$  according to:

$$P(T > c; \mu = 0) = \alpha, \text{ or } \Phi\left(\frac{\sqrt{nc}}{\sigma}\right) = 1 - \alpha$$

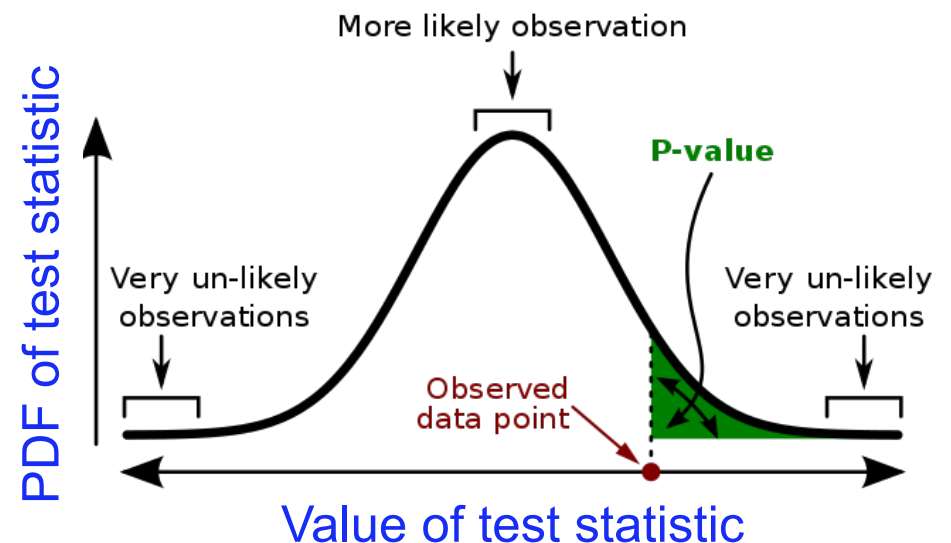
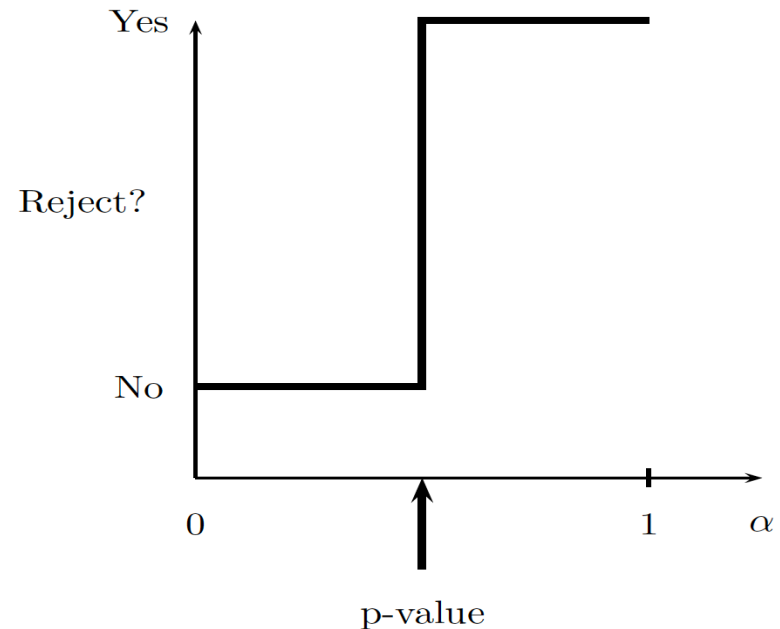
This yields: 
$$c = \frac{\sigma}{\sqrt{n}} \Phi^{-1}(1 - \alpha)$$

# p-value

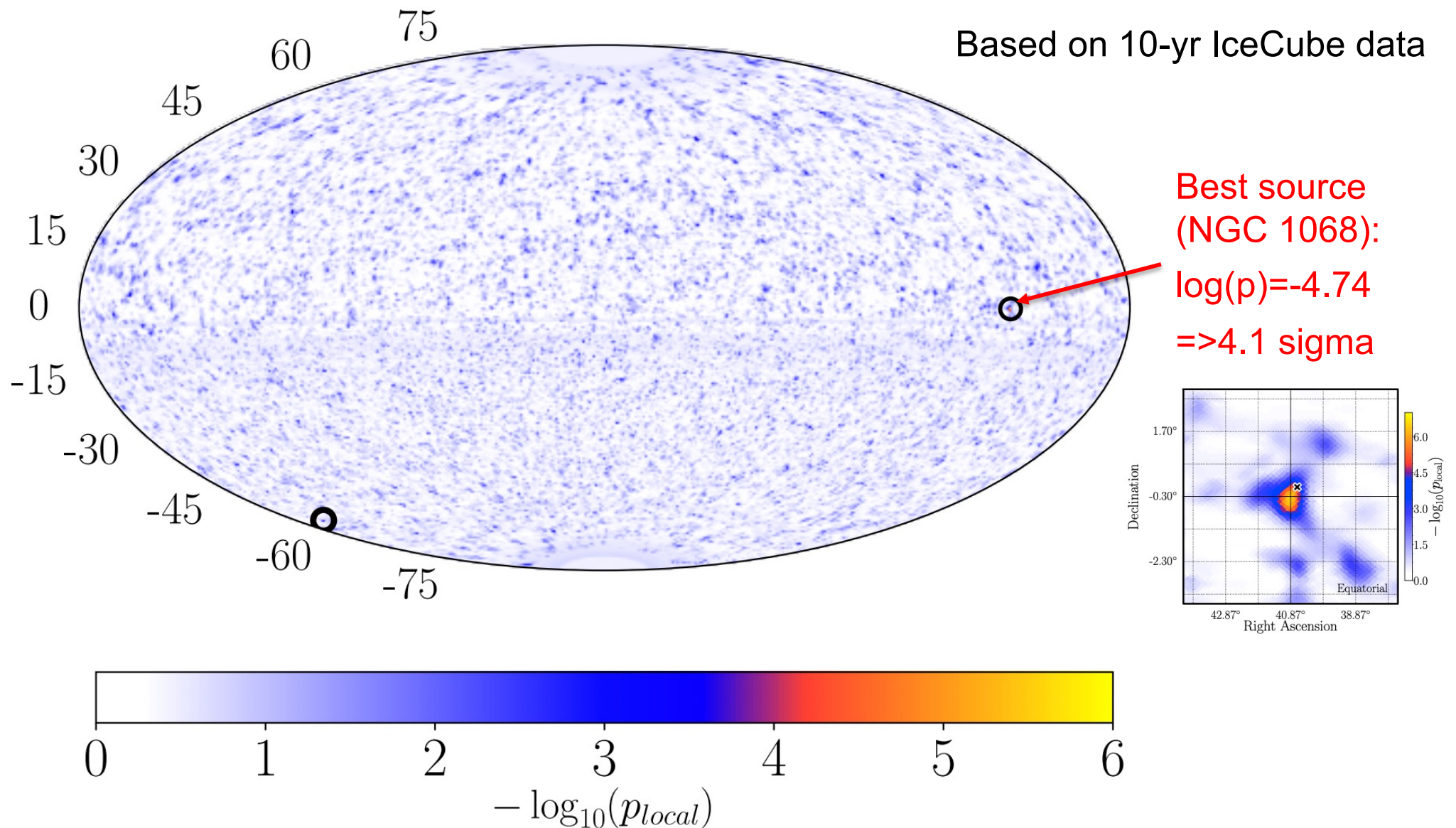
Given data, whether we reject the hypothesis depends on the significance level  $\alpha$ .

We can further ask that given the data, at what level that the null hypothesis can be rejected.

This level is called the **p-value**: the probability of observing the test statistic at current or more extreme values, assuming the null hypothesis is true.



# Example: astro neutrino sources





# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Likelihood ratio test

Consider testing:  $H_0 : \theta \in \Theta_0$  versus  $H_1 : \theta \in \Theta/\Theta_0$

The **likelihood ratio statistic** is defined as

$$\lambda = 2 \log \left( \frac{\sup_{\theta \in \Theta} \mathcal{L}(\theta)}{\sup_{\theta \in \Theta_0} \mathcal{L}(\theta)} \right) = 2 \log \left( \frac{\mathcal{L}(\hat{\theta})}{\mathcal{L}(\hat{\theta}_0)} \right)$$

parameter that  
achieves maximum  
likelihood

Clearly, one should have  $\lambda \geq 0$ .

The log likelihood ratio should be close to 0 if  $H_0$  is true, and larger otherwise.

The rejection region should be determined as

$$W_0 = \{\mathbf{x} : \lambda(\mathbf{x}) > \lambda_0\}$$

data  
values

critical value set by  
significance level

# Example

Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sigma$  is known.

Now test  $H_0: \mu = \mu_0$  versus  $H_1: \mu \neq \mu_0$ .

With  $\sigma$  readily known, the likelihood function reads:

$$\mathcal{L}_n(\mu) \propto \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right)$$

The likelihood ratio for this problem is thus given by

$$\lambda = 2 \log \frac{\mathcal{L}_n(\bar{x})}{\mathcal{L}_n(\mu_0)} = \frac{1}{\sigma^2} \left( \sum_{i=1}^n (x_i - \mu_0)^2 - \sum_{i=1}^n (x_i - \bar{x})^2 \right) = \frac{n}{\sigma^2} (\bar{x} - \mu_0)^2$$

Effectively, one can simply choose  $\bar{X}$  as test statistic, which satisfies a Gaussian distribution, with rejection region given by

$$W_0 = \{\mathbf{x} : \lambda(\mathbf{x}) > \lambda_0\} = \{\mathbf{x} : |\bar{x} - \mu_0| > c\}$$

with  $c$  determined by  $P(|\bar{X} - \mu_0| > c) = \alpha = 0.05$

# Example

Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sigma$  is unknown.

Now test  $H_0: \mu = \mu_0$  versus  $H_1: \mu \neq \mu_0$ .

With  $\sigma$  unknown, the likelihood function reads:

$$\mathcal{L}_n(\mu, \sigma) \propto \frac{1}{\sigma^n} \exp \left( -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right) = \frac{1}{\sigma^n} \exp \left( -\frac{1}{2\sigma^2} \left[ \sum_{i=1}^n (x_i - \bar{x})^2 + n(\bar{x} - \mu)^2 \right] \right)$$

Recall that the maximum likelihood estimate of parameters are

$$\hat{\mu} = \bar{x}, \quad \hat{\sigma}^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2 \quad \text{or when } \mu = \mu_0 \text{ is known: } \hat{\sigma}_0^2 = \frac{1}{n} \sum_{i=1}^n (x_i - \mu_0)^2$$

For full parameter space:

$$\mathcal{L}_n(\hat{\mu}, \hat{\sigma}) \propto \frac{1}{\hat{\sigma}^n} e^{-n/2} \propto \left[ \frac{n}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{\frac{n}{2}}$$

For null hypothesis, similarly:

$$\mathcal{L}_n(\mu_0, \hat{\sigma}_0) \propto \left[ \frac{n}{\sum_{i=1}^n (x_i - \mu_0)^2} \right]^{\frac{n}{2}}$$

# Example

Let  $X_1, \dots, X_n \sim \mathcal{N}(\mu, \sigma^2)$  where  $\sigma$  is unknown.

Now test  $H_0: \mu = \mu_0$  versus  $H_1: \mu \neq \mu_0$ .

Therefore, the likelihood ratio is

$$\frac{\mathcal{L}_n(\hat{\mu}, \hat{\sigma})}{\mathcal{L}_n(\mu_0, \hat{\sigma}_0)} = \left[ \frac{\sum_{i=1}^n (x_i - \mu_0)^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right]^{\frac{n}{2}} = \left( 1 + \frac{T^2}{n-1} \right)^{\frac{n}{2}}$$

where  $T = \frac{\sqrt{n(n-1)}(\bar{x} - \mu_0)}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2}}$  satisfying the student's  $t$  distribution.  
(with  $n-1$  degrees of freedom)

Effectively, we can choose  $T$  as the test statistic, with rejection region given by

$$W_0 = \{\mathbf{x} : |T| > c\}$$

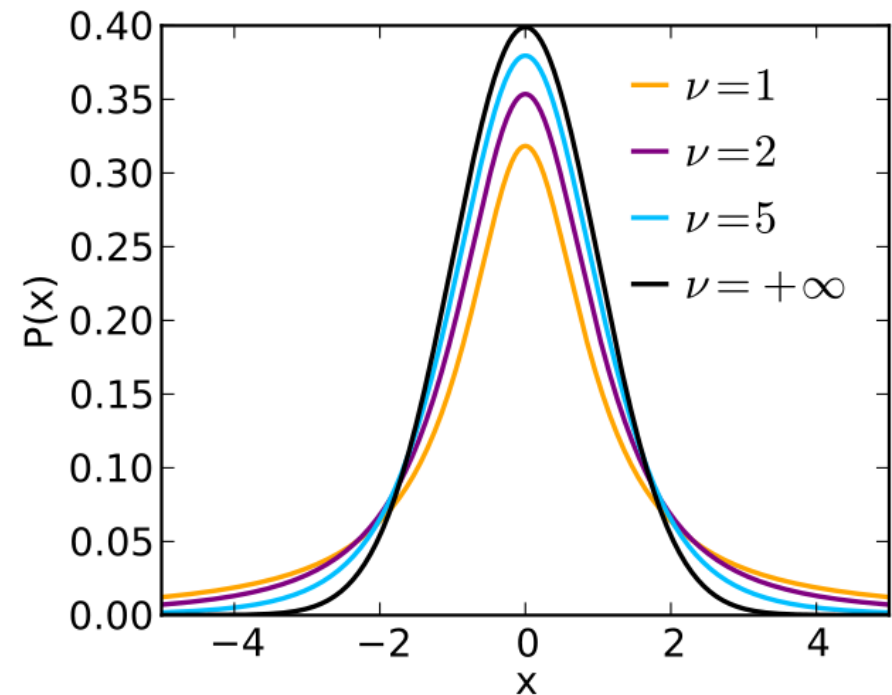
with  $c$  determined by  $P(|T| > c) = \alpha = 0.05$ .

# Recall: Student's t distribution

The PDF of student's t distribution with  $\nu$  degrees of freedom reads

$$f(x) = \frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\pi\nu}\Gamma(\nu/2)} \left(1 + \frac{x^2}{\nu}\right)^{-\frac{\nu+1}{2}}$$

It approaches standard normal distribution for large  $\nu$ .



If  $X_1, \dots, X_n$  are independent, standard normal random variables, define

$$\bar{X}_n = \frac{1}{n}(X_1 + \dots + X_n), \quad S_n^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X}_n)^2 \quad (\text{unbiased estimates of mean and variance})$$

Then the variable  $T \equiv \frac{\sqrt{n}}{S_n}(\bar{X}_n - \mu)$  satisfies student's t distribution with  $n-1$  degrees of freedom.

# General cases with Gaussian distribution

One-sample test:

Situation	Test statistic	Distribution	Dof	
Test on $\mu$ with $\sigma=\sigma_0$ known	$\bar{X}$	$\mathcal{N}(\mu, \sigma_0^2/n)$		a.k.a. Z-test
Test on $\mu$ with $\sigma$ unknown	$T = \frac{\sqrt{n(n-1)}(\bar{X} - \mu_0)}{\sqrt{\sum_{i=1}^n (X_i - \bar{X})^2}}$	t	n-1	
Test on $\sigma$ with $\mu=\mu_0$ known	$\chi^2 = \sum_{i=1}^n \frac{(X_i - \mu_0)^2}{\sigma_0^2}$	$\chi^2$	n	
Test on $\sigma$ with $\mu$ unknown	$\chi^2 = \sum_{i=1}^n \frac{(X_i - \bar{X})^2}{\sigma_0^2}$	$\chi^2$	n-1	

# Applications:

- Quality check: is the weight of the chocolate bars the same as claimed?

One-sample test with  $t$  statistic (one-sample  $t$ -test).

- How effective is a certain method to reduce weight?

One-sample test with  $t$  statistic (one-sample  $t$ -test).

- The city bus is supposed to run every 10 minutes. How can we tell if the bus is sufficiently on time (i.e., variance less than say 2 minutes)?

One-sample test with  $\chi^2$  statistic.

- Are men on average taller than women?

Two-sample test with  $t$  statistic (two-sample  $t$ -test).

- Goodness of fit: is there a linear relation between  $X$  and  $Y$ ?

Read: fit  $y=ax+b$ , does  $a=0$ ? => Does linear fit significantly improve the residuals?

Two-sample test with  $F$  statistic (two-sample  $F$ -test).



# Two-sample tests

Let  $X_1, \dots, X_{n_1} \sim \mathcal{N}(\mu_1, \sigma_1^2)$ ,  $Y_1, \dots, Y_{n_2} \sim \mathcal{N}(\mu_2, \sigma_2^2)$  where  $\mu_{1,2}$  are unknown.

Now test  $H_0: \sigma_1^2 = \sigma_2^2$  versus  $H_1: \sigma_1^2 \neq \sigma_2^2$ .

One can go through similar (but tedious) procedures to derive the likelihood ratio and arrive at the test statistic to be

$$F = \frac{\sum_{i=1}^{n_1} (x_i - \bar{x})^2 / (n_1 - 1)}{\sum_{j=1}^{n_2} (y_j - \bar{y})^2 / (n_2 - 1)}$$

which satisfies the  $F$  distribution with  $(n_1-1, n_2-1)$  degrees of freedom.

The rejection region can be set by requiring:  $P(F < c_1) = P(F > c_2) = \alpha/2$

The values of  $c_1, c_2$  can be obtained by consulting standard tables.

This test is known as the F-test, which is important in analysis of variance (ANOVA).

# Recall: Fisher's F distribution

The Fisher's F distribution is a two-parameter family whose PDF reads

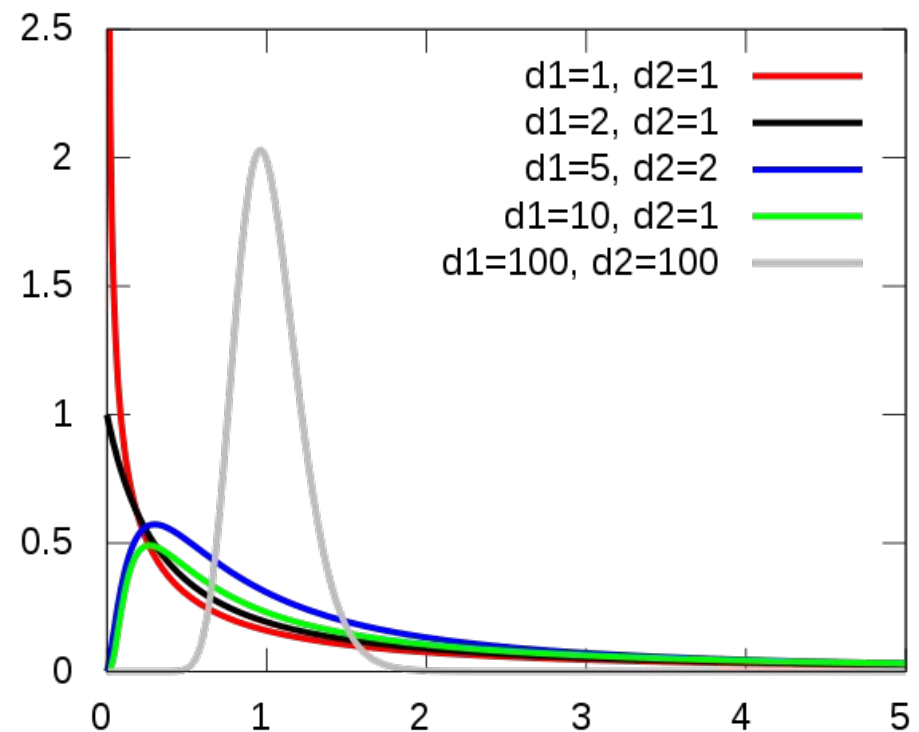
$$f(x) = \underbrace{B\left(\frac{d_1}{2}, \frac{d_2}{2}\right)}_{\text{beta function}}^{-1} \left(\frac{d_1}{d_2}\right)^{\frac{d_1}{2}} x^{\frac{d_1}{2}-1} \left(1 + \frac{d_1}{d_2}x\right)^{-\frac{d_1+d_2}{2}}$$

This is usually denoted by  $F(d_1, d_2)$ .

This statistics arises from the ratio of two independent reduced chi squared:

$$X_1 \sim \chi_{d_1}^2, X_2 \sim \chi_{d_2}^2$$

$$\text{Then } \frac{X_1/d_1}{X_2/d_2} \sim F(d_1, d_2)$$



# General cases with Gaussian distribution

Two-sample test:

Situation	Test statistic	Distribution	Dof
Compare $\sigma_1^2/\sigma_2^2$ with $\Delta_0$ , with $\mu_{1,2}$ unknown	$F = \frac{\sum_{i=1}^{n_1} (X_i - \bar{X})^2}{\sum_{j=1}^{n_2} (Y_j - \bar{Y})^2} \frac{1}{\Delta_0} \frac{n_2 - 1}{n_1 - 1}$	$F$	$(n_1-1, n_2-1)$
Compare $\mu_1, \mu_2$ knowing $\sigma_1 = \sigma_2$	$T = \frac{\sqrt{n_1 n_2 (n_1 + n_2 - 2) / (n_1 + n_2)} (\bar{X} - \bar{Y})}{\sqrt{\sum_{i=1}^{n_1} (X_i - \bar{X})^2 + \sum_{j=1}^{n_2} (Y_j - \bar{Y})^2}}$	$t$	$n_1 + n_2 - 2$
Compare $\mu_1, \mu_2$ with $\sigma_1 \neq \sigma_2$	$T = \frac{\bar{X} - \bar{Y}}{\sqrt{S_1^2/n_1 + S_2^2/n_2}}$ $S_{1,2}^2 = \frac{1}{n_{1,2}} \sum_1^{n_{1,2}} (X_i - \bar{X})^2 \text{ or } (Y_j - \bar{Y})^2$	$t$ (approximately)	round( $m^*$ )
<p>where <math>m^* = \left( \frac{S_1^2}{n_1} + \frac{S_2^2}{n_2} \right)^2 \left[ \frac{1}{n_1 - 1} \left( \frac{S_1^2}{n_1} \right)^2 + \frac{1}{n_2 - 1} \left( \frac{S_2^2}{n_2} \right)^2 \right]^{-1}</math></p>			

# Outline

## ■ Parameter estimation and confidence interval

- Fundamental concepts
- Non-parametric estimations
- Bootstrap and jackknife
- Maximum likelihood method

## ■ Hypothesis testing

- Basic concepts
- Likelihood ratio test and applications to Gaussian samples
- Comparing distributions

# Comparing two distributions

Given two sets of data, we want to ask, are they drawn from the same distribution function?

- Are the arrival directions of ultra-high-energy cosmic-rays consistent from coming uniformly in the sky?
- Does the galaxy luminosity function evolve with redshift?
- Does my source property change between two measurements?
- How does the architecture of the solar system compare with other known exoplanetary systems?

**Null hypothesis:** two distributions are the same.

Note: data can be either continuous or discrete/binned, and we either compare data to a known distribution, or to another data set.

# Kolmogorov-Smirnov (K-S) test

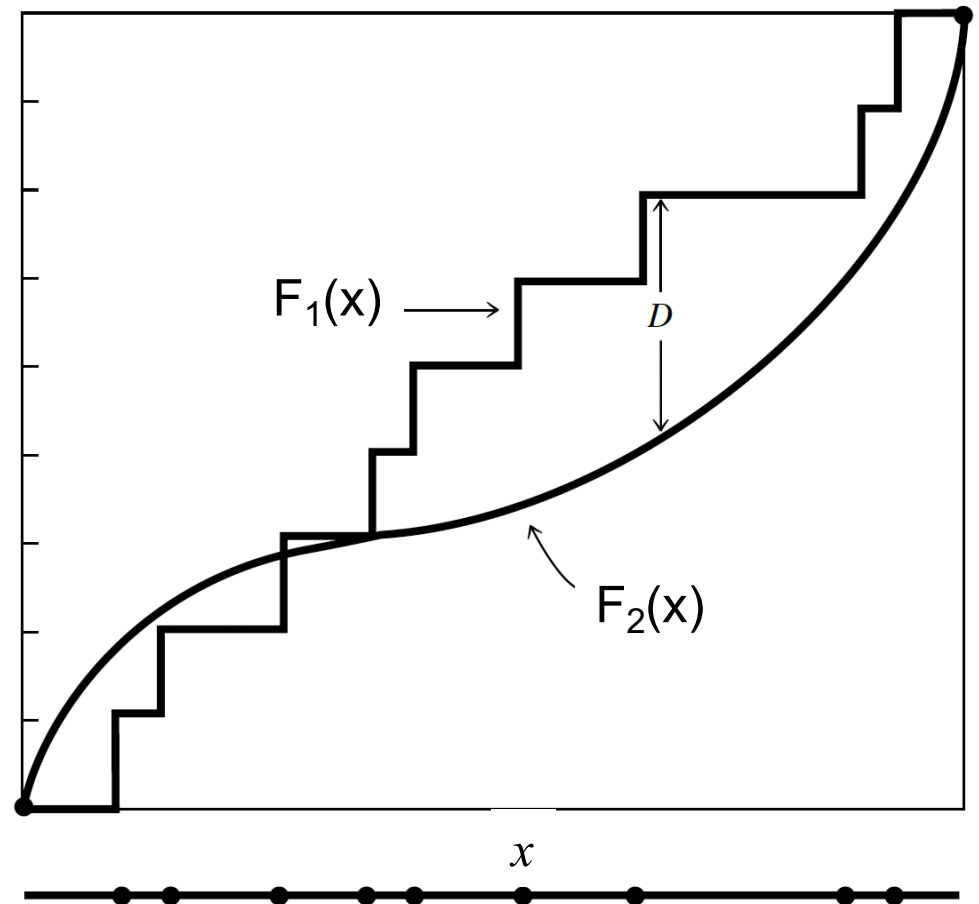
The most popular non-parametric test for comparing two distributions.

The test is based on [the K-S statistic](#), which measures the maximum distance between two CDFs (EDFs):

$$D = \max_{-\infty < x < \infty} |F_1(x) - F_2(x)|$$

[One-sample K-S test](#):  $F_2$  is a reference distribution;

[Two-sample K-S test](#):  $F_2$  is the CDF/EDF from another sample.



# Kolmogorov-Smirnov (K-S) test

The probability to obtain a sample distribution with D larger than measured is

$$Q_{KS}(\lambda) = 2 \sum_{j=1}^{\infty} (-1)^{j-1} e^{-2j^2 \lambda^2} \quad [\text{Note: } Q_{KS}(0)=1, Q_{KS}(\infty)=0]$$

where

$$\lambda = \left( 0.12 + \sqrt{N_e} + \frac{0.11}{\sqrt{N_e}} \right) D \quad \text{with} \quad N_e = \frac{N_1 N_2}{N_1 + N_2}$$

For one-sample distribution, take  $N_2 = \infty$ , so that  $N_e = N_1$ .

It is not the best choice for Gaussian distribution (not sensitive to the tails).

There are ways to generalize it to multi-D, though not very straightforward.

# Other methods

Kuiper test:

$$D^* = \max[F_1(x) - F_2(x)] + \max[F_2(x) - F_1(x)]$$

It is invariant under cyclic transformation.

Mann-Whitney U test and Wilcoxon signed-rank test:

Non-parametric analog of t-test of the mean for Gaussian distributions, involving using ranks (sorting the data).

Anderson-Darling test and Shapiro-Wilk test:

More sensitive to differences in the tail of the distribution, and are used primarily for testing whether data is drawn from a normal distribution.



# Pearson's $\chi^2$ test (and G-test)

Statistical test applied to sets of **categorical data** (taking finite # of values) to evaluate how likely the observed differences between the sets arise by chance.

The **test statistic** is defined as

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

# of observations of type i

expected # of type i

Superseded by  
(likelihood ratio):

$$G = 2 \sum_{i=1}^n O_i \ln \left( \frac{O_i}{E_i} \right)$$

Expected to approximately satisfy  $\chi^2$  distribution. The # of **degree of freedom** is:

Total # of data points – total # of constraints (e.g., parameters)

Three types of comparison: **goodness of fit**, **homogeneity**, and **independence**.

Require large sample size (~1000) to be exact.

# Example: fairness of dice

A 6-sided dice is thrown 60 times. Given the result, is it fair?

$i$	$O_i$	$E_i$	$O_i - E_i$	$(O_i - E_i)^2$	$\frac{(O_i - E_i)^2}{E_i}$
1	5	10	-5	25	2.5
2	8	10	-2	4	0.4
3	9	10	-1	1	0.1
4	8	10	-2	4	0.4
5	10	10	0	0	0
6	20	10	10	100	10
Sum					13.4

Null hypothesis: it is fair.

Value of test statistic: 13.4

Dof: 6-1=5.

We find a p-value of 0.02.

Therefore, fairness is rejected at  $\alpha=0.05$  level, but retained at  $\alpha=0.01$  level.

# Summary

- Parameter estimation and confidence interval
  - Concepts: sample distribution, standard error, confidence set.
  - Non-parametric estimations: EDF and plug-in estimator.
  - Bootstrap and jackknife: important resampling methods.
  - Maximum likelihood method: asymptotic normality, Fisher information, asymptotically optimal.
- Hypothesis testing
  - Basic concepts: null/alternative hypothesis, type 1/2 error, significance level, p-value.
  - Likelihood ratio test and applications to Gaussian samples: contains several cases of  $\chi^2$  test,  $t$ -test,  $F$ -test, etc.
  - Comparing distributions: KS test, Pearson's  $\chi^2$ /G -test, etc.