

level, test statistic and rejection region

A test ϕ has level α if

$\Pr_{\theta_0}(\phi) \leq \alpha$, $\forall \theta_0 \in \Theta_0$

In general, a test has the form

$\phi: \Omega \rightarrow \{0, 1\}$ -> test statistic

for some statistic T_n and threshold $c \in \mathbb{R}$

& rejection region is $R = \{T_n > c\}$

One-sided v.s. two-sided tests:

When $\theta \in \Theta \subset \mathbb{R}$ and H_0 is of the form

$H_0: \theta \leq \theta_0 \Leftrightarrow \phi = \{0, 1\}$

$\{\theta: \theta > \theta_0\}$ or $H_0: \theta \geq \theta_0$ one-sided test

p-value

Definition: The asymptotic p-value of a test ϕ_θ is the smallest (asymptotic) level α at which ϕ_θ rejects H_0 .

It is random as it depends on the sample.

Golden rule: p-value $\leq \alpha \Leftrightarrow H_0$ is rejected by ϕ_θ at the (asymptotic) level α .

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

Concave and Convex

Definition:

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

$$A''(\theta) \geq 0, \forall \theta \in \Theta$$

It is said to be strictly convex if the inequality is strict, $A''(\theta) > 0$

More over, A is said to be (strictly) convex if $-A$ is (strictly) convex.

Multivariate convex function

For a multivariate function $A: \Theta \subset \mathbb{R}^d \rightarrow \mathbb{R}$, $d \geq 2$, define the gradient vector:

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

Hessian matrix:

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

A is convex $\Leftrightarrow \nabla^2 A(\theta) \geq 0, \forall \theta \in \Theta$

A is strictly convex $\Leftrightarrow \nabla^2 A(\theta) > 0, \forall \theta \in \Theta$

Optimizing Conditions

Strictly concave functions are easy to maximize: if they have a maximum then it is unique. It is the unique solution to

$$A'(\theta) = 0 \quad / \quad \text{high-convex}$$

high-convex

Convex optimization: close form solution

Multivariate Gaussian distribution

A Gaussian vector $x \in \mathbb{R}^d$ is completely determined by its expected value $E[x] \in \mathbb{R}^d$ and covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$

The pdf is given by

$$f(x) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{-\frac{1}{2} (x-E[x])^\top \Sigma^{-1} (x-E[x])} \quad (\in \mathbb{R}^d)$$

Multivariate CLT

The CLT may be generalized to averages of random vectors (also vectors of averages).

Let $X_1, \dots, X_n \in \mathbb{R}^d$ be independent copies of a random vector X such that

$$E[X] = \mu, \text{ Cov}(X) = \Sigma$$

$$\text{If } X \sim N(\mu, \Sigma), \text{ then } X \sim N(\mu, \Sigma)$$

Multivariate Delta Method

Let $(\hat{\theta}_n)$ be a sequence of random vectors in \mathbb{R}^d such that

$$\sqrt{n}(\hat{\theta}_n - \theta) \xrightarrow{D} N(0, \Sigma)$$

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 θ , then

$$\sqrt{n}(g(\hat{\theta}_n) - g(\theta)) \xrightarrow{D} N_g(0, \Sigma \otimes \Sigma)$$

where $\Sigma \otimes \Sigma = \sum_{i,j} \Sigma_{ij} \Sigma_{ij}^\top \in \mathbb{R}^{d \times d}$

Let $\theta^* = g^{-1}(\theta)$, then $\hat{\theta}_n \xrightarrow{P} \theta^*$

$\hat{\theta}_n \xrightarrow{P} \theta^* \Leftrightarrow \sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{D} N(0, \Sigma \otimes \Sigma)$

Asymptotic normality of the MLE

Theorem:

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

1. The parameter is identifiable.

2. For all $\theta \in \Theta$, the support of P_θ does not depend on θ .

3. θ^* is not on the boundary of Θ .

4. $I(\theta)$ is invertible in a neighborhood of θ^* .

5. A few more technical conditions.

Then $\hat{\theta}_n$ satisfies:

$$\hat{\theta}_n \xrightarrow{P} \theta^* \text{ and } P$$

$$\sqrt{n}(\hat{\theta}_n - \theta^*) \xrightarrow{D} N(0, I(\theta^*)^{-1}) \text{ with } P$$

Method of Moments

Moments:

Let x_1, \dots, x_n be an iid sample associated with a statistical model $(E, P_{\theta})_{\theta \in \Theta}$. Assume that $E \in \mathbb{R}$ and $\theta \in \mathbb{R}^d$, for some $d \geq 1$.

Population moments: Let $m_\theta = E_\theta(X^1)$, $\text{cov}_\theta = E_\theta(X^1 X^2)$

Empirical moments: Let $\hat{m}_n = \bar{x}_n = \frac{1}{n} \sum_i x_i$, $\hat{\text{cov}}_n = \frac{1}{n} \sum_i (x_i - \bar{x}_n)(x_i - \bar{x}_n)^\top$

From CLT,

$$\frac{\hat{m}_n - m_\theta}{\sqrt{n}} \xrightarrow{D} 0$$

more compactly, we say that the whole vector converges:

$$(x_1, \dots, x_n) \xrightarrow{D} (m_\theta, \text{cov}_\theta, \dots, \text{cov}_\theta)$$

Unit 3: Methods of Estimation

Lecture 8: Distance measures between distributions

Total variation distance

• Definition:

$L(E, P_{\theta_0})$ is a statistical model associated with a sample of i.i.d. r.v. X_1, \dots, X_n . Assume that there exist $\theta \in \Theta$ that $\lambda \sim P_\theta$ is the true parameter distribution, and given X_1, \dots, X_n , find an estimator $\hat{\theta} = \hat{\theta}(X_1, \dots, X_n)$ s.t. P_θ is close to $P_{\hat{\theta}}$ for true parameter θ .

This means: $|P_\theta(X) - P_{\hat{\theta}}(X)|$ is small for all $\theta \in \Theta$

• Definition:

Total variation distance between two probability measures P_θ and $P_{\hat{\theta}}$:

$$TV(P_\theta, P_{\hat{\theta}}) = \max_{A \in \mathcal{B}(\mathbb{R})} |P_\theta(A) - P_{\hat{\theta}}(A)|$$

• Discrete sample space:

$$TV(P_\theta, P_{\hat{\theta}}) = \sum_{x \in \mathcal{X}} |P_\theta(x) - P_{\hat{\theta}}(x)|$$

• Continuous sample space:

$$TV(P_\theta, P_{\hat{\theta}}) = \int_{\mathbb{R}} |P_\theta(x) - P_{\hat{\theta}}(x)| dx$$

• Properties of total variation:

-Symmetric

$$TV(P_\theta, P_{\hat{\theta}}) = TV(P_{\hat{\theta}}, P_\theta)$$

-Positive

$$0 \leq TV(P_\theta, P_{\hat{\theta}}) \leq 1$$

-Definite

$$\text{If } TV(P_\theta, P_{\hat{\theta}}) = 0, \text{ then } P_\theta = P_{\hat{\theta}}$$

Triangle inequality

$$TV(P_\theta, P_{\hat{\theta}}) \leq TV(P_\theta, P_{\hat{\theta}'}) + TV(P_{\hat{\theta}'}, P_{\hat{\theta}})$$

Balakrishnan-Laike's (KL) divergence

• Definition:

The KL divergence between two probability measures P_θ and $P_{\hat{\theta}}$ is defined by

$$KL(P_\theta, P_{\hat{\theta}}) = \int_{\mathbb{R}} \log \frac{P_\theta(x)}{P_{\hat{\theta}}(x)} dP_{\hat{\theta}}(x) \quad \text{if } E[P_{\hat{\theta}}] \text{ is discrete}$$

$\int_{\mathbb{R}} \log \frac{P_\theta(x)}{P_{\hat{\theta}}(x)} dP_{\hat{\theta}}(x) \quad \text{if } E[P_{\hat{\theta}}] \text{ is continuous}$

• Properties of KL divergence:

-Symmetric: $KL(P_\theta, P_{\hat{\theta}}) = KL(P_{\hat{\theta}}, P_\theta)$ in general

-Positive: $KL(P_\theta, P_{\hat{\theta}}) \geq 0$

-Definite: If $KL(P_\theta, P_{\hat{\theta}}) = 0$, then $P_\theta = P_{\hat{\theta}}$

• Principle separation: $KL(P_\theta, P_{\hat{\theta}}) \leq KL(P_\theta, P_0) + KL(P_0, P_{\hat{\theta}})$ in general

Lecture 9: Introduction to MLE

Maximum Likelihood Estimation

• Definition:

$$KL(P_\theta, P_{\hat{\theta}}) = E_{\hat{\theta}}[\log \frac{P_\theta(X)}{P_{\hat{\theta}}(X)}] = E_{\hat{\theta}}[\log p_\theta(X)] - E_{\hat{\theta}}[\log p_{\hat{\theta}}(X)]$$

↓ $\min_{\theta \in \Theta} KL(P_\theta, P_{\hat{\theta}})$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[\log p_\theta(X)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[\log p_\theta(X)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[p_\theta(X)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[p_\theta(X)] = \min_{\theta \in \Theta} \prod_{i=1}^n p_\theta(x_i)$

↓ Maximum Likelihood Principle

• Definition:

$$KL(P_\theta, P_{\hat{\theta}}) = E_{\hat{\theta}}[\log p_\theta(X)] = E_{\hat{\theta}}[\log p_{\hat{\theta}}(X)]$$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[\log p_{\hat{\theta}}(X)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[\log p_{\hat{\theta}}(X)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_{\hat{\theta}}(X)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_{\hat{\theta}}(X)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(X)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(X)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] + \dots + E_{\hat{\theta}}[-\log p_\theta(x_n)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)] = \min_{\theta \in \Theta} E_{\hat{\theta}}[-\log p_\theta(x_1)]$

↓ $\min_{\theta \$

Small sample size:
 - Can not really apply Student's t-test
 - We need to find the (asymptotic) distribution of quantiles of the form
 $\frac{X_n}{\sigma_n}$
 when $X_1, \dots, X_n \sim N(\mu, \sigma^2)$
 It turns out that this distribution does not depend on μ or σ , so we can compute its quantiles.

The t^* distribution

Definition:
 For a positive integer d , the t^* (pronounced "Kai-squared") distribution with d degrees of freedom is the law of the random variable T^* where $Z \sim N(0, 1)$, $V \sim \chi^2_d$ and $Z \perp V$ (Z is independent of V).

$$Z^2 + \frac{V}{d} \sim \chi^2_{d+2} \text{ where } Z = \frac{X_n - \mu}{\sigma_n} \sim N(0, 1)$$

Examples:
 - If $Z \sim N(0, 1)$, then $|Z| \sim t^*$

$$Z \sim \text{Exp}(1)$$

Properties for $V \sim \chi^2_n$

$$E(V) = E(Z^2) = 1 = E(\chi^2_1) = d$$

$$\text{Var}(V) = \text{Var}(Z^2) = 2 = \text{Var}(\chi^2_2) = 2d$$

Important example: the sample variance

Sample variance is given by

$$S^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$$

Cochran's theorem: for $X_1, \dots, X_n \sim N(\mu, \sigma^2)$, S^2 is the sample variance:

$$S^2 \sim \chi^2_{n-1}$$

and $\frac{S^2}{\sigma^2} \sim \chi^2_{n-1}$

* We often prefer the unbiased estimator of σ^2 : $\hat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$

Likelihood ratio test: A test based on the log-likelihood

Consider an iid sample X_1, \dots, X_n with statistical model $(E, P_{\theta(x)})$, where $\theta \in \Theta(d, n)$.

Suppose the null hypothesis has the form

$$H_0: (\theta_1, \dots, \theta_d) = (\theta_1^{(0)}, \dots, \theta_d^{(0)})$$

for some fixed and given numbers $\theta_1^{(0)}, \dots, \theta_d^{(0)}$

Let

$$\theta_0 = \arg\max_{\theta \in \Theta} L(\theta) \quad (\text{MLE})$$

and

$$\hat{\theta}_n = \arg\max_{\theta \in \Theta} L_n(\theta) \quad (\text{"constrained MLE"})$$

Test statistic:

$$T_n = 2(L(\hat{\theta}_n) - L(\theta_0))$$

Wilks' theorem:

Assume H_0 is true and MLE technical conditions are satisfied, then

$$T_n \xrightarrow{d} \chi^2_{k-d}$$

Likelihood ratio test with asymptotic level α (0.05):

$$\Phi = \Phi(T_n > \chi^2_{k-d}) \rightarrow (1-\alpha) \text{-quantile of } \chi^2_{k-d}$$

Goodness of fit test (discrete, χ^2)

Let $X_1, \dots, X_n \sim P_p$, for some unknown p e.g., and $p \in \Delta$ is fixed

Test: $H_0: p = p^*$ vs. $H_1: p \neq p^*$, with asymptotic level α (0.05)

Likelihood of model:

$$L(X_1, \dots, X_n; p) = p_1^{X_1} p_2^{X_2} \cdots p_n^{X_n}$$

where $p_j = p_j(p_1, \dots, p_n, X_j)$

Let \hat{p} be the MLE:

$$\hat{p}_j = \frac{X_j}{n}, \quad j=1, \dots, k$$

→ minimizes $\log L(X_1, \dots, X_n; p)$ under the constraint

χ^2 -test:

If H_0 is true, then $\hat{\chi}^2 = \sum_{j=1}^k (p_j - \hat{p}_j)^2 / \hat{p}_j$ is asymptotically normal and the following holds

Theorem (under H_0):

$$\sqrt{n} \frac{\hat{\chi}^2 - (k-1)}{\sqrt{(k-1)}} \xrightarrow{d} \chi^2_{k-1}$$

• χ^2 -test with asymptotic level α :

$$\Phi = \Phi(\hat{\chi}^2 > \chi^2_{k-1}) \rightarrow (1-\alpha) \text{-quantile of } \chi^2_{k-1}$$

• Asymptotic p-value of this test:

$$p\text{-value} = P(\hat{\chi}^2 > \hat{\chi}^2)$$

• Pointwise convergence:

Uniform convergence

$$\lim_{n \rightarrow \infty} p(\hat{\chi}^2 > \hat{\chi}^2) = 1 - \alpha$$

For every $\epsilon > 0$, there exists an n_0 s.t.

such that $p(\hat{\chi}^2 > \hat{\chi}^2) \leq \epsilon$ for all $n \geq n_0$

Goodness of fit test (continuous, Kolmogorov-Smirnov):

Test statistic: $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| \rightarrow$ p-value, use table

Quantile-Quantile ($Q-Q$) plot:

A visual way of performing GOF tests

Not a formal test, but quick and easy way to check if a distribution is plausible

Main idea: check visually if the plot F_n is close to that of F or equivalently if the plot of F_n is close to that of F

• Define: $F_n(x_m) = x_m$ → the largest observation

• Check if points $(F_n(x_1), F_n(x_2)), (F_n(x_2), F_n(x_3)), \dots, (F_n(x_n), F_n(x_1))$ are near $y=x$

graph

Left tail: $F_n(x) < x$

Right tail: $F_n(x) > x$

Converges to $y=x$

• F -test with asymptotic level α :

$$\Phi = \Phi(\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| > \chi^2_{k-1})$$

• Given p_1, \dots, p_n , let $F_n(x) = p_1 I_{(-\infty, x]} + \dots + p_n I_{(-\infty, x]}$

• Hence: $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| = \sup_{x \in \mathbb{R}} |p_1 I_{(-\infty, x]} + \dots + p_n I_{(-\infty, x)} - F(x)|$

• Bayes's formula

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

• Bayes's theorem: $P(A|B) = \frac{P(B|A)P(A)}{P(B|A)P(A) + P(B|\bar{A})P(\bar{A})}$

• Given p_1, \dots, p_n , let $P(A) = p_1 + \dots + p_n$

• Hence: $\sup_{x \in \mathbb{R}} |F_n(x) - F(x)| = \sup_{x \in \mathbb{R}} |\frac{p_1 I_{(-\infty, x]} + \dots + p_n I_{(-\infty, x)}}{p_1 + \dots + p_n} - F(x)|$

• Conjugate prior

(posterior is family with prior)

Student's T distribution

Definition:

For a positive integer d , the Student's T distribution with d degrees of freedom (denoted by T^*) is the law of the random variable T^* where $Z \sim N(0, 1)$, $V \sim \chi^2_d$ and $Z \perp V$ (Z is independent of V).

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

Let $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ where both μ and σ^2 are unknown

We want to test:

$$H_0: \mu = \mu_0 \quad \text{vs.} \quad H_1: \mu \neq \mu_0$$

Test statistic:

$$T_n = \frac{\bar{X}_n - \mu_0}{\hat{\sigma}_n} \sim t_{n-1}$$

Since $\bar{X}_n \sim N(\mu, \sigma^2/n)$, and $\hat{\sigma}_n^2 \sim \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ are independent (Cochran's theorem)

$T_n \sim t_{n-1}$ (Student's T distribution with $n-1$ degrees of freedom)

Student's T test with (non asymptotic) level α (0.05):

$$\Phi = \Phi(|T_n| > t_{n-1, 1-\alpha/2}) \rightarrow (1-\alpha) \text{-quantile of } t_{n-1}$$

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

as above

We want to test:

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

Test statistic:

$$T_n = \frac{\bar{X}_n - \mu_0}{\hat{\sigma}_n} \sim t_{n-1} \quad (\text{under } H_0)$$

Student's T test with (non asymptotic) level α (0.05):

$$\Phi = \Phi(T_n > t_{n-1, 1-\alpha}) \rightarrow (1-\alpha) \text{-quantile of } t_{n-1}$$

Advantage & disadvantage of Student's T test

- Advantage: Non-asymptotic, it can be run on small samples + can always use large sample size
- Disadvantage: assumption that sample is Gaussian

Lecture 14: Wald's Test, Likelihood Ratio test and Implicit Hypothesis Test

Wald's Test: A test based on the MLE

Consider an iid sample X_1, \dots, X_n with statistical model $(E, P_{\theta(x)})$, where $\theta \in \Theta(d, n)$

Suppose the null hypothesis has the form

$$H_0: (\theta_1, \dots, \theta_d) = (\theta_1^{(0)}, \dots, \theta_d^{(0)})$$

for some fixed and given numbers $\theta_1^{(0)}, \dots, \theta_d^{(0)}$

Let

$$\hat{\theta} = \arg\max_{\theta \in \Theta} L(\theta) \quad (\text{MLE})$$

and

$$\hat{\theta}_n = \arg\max_{\theta \in \Theta} L_n(\theta) \quad (\text{"constrained MLE"})$$

Test statistic:

$$T_n = 2(L(\hat{\theta}_n) - L(\hat{\theta}))$$

Wald's test with asymptotic level α (0.05):

$$\Phi = \Phi(T_n > \chi^2_{d-1}) \rightarrow (1-\alpha) \text{-quantile of } \chi^2_{d-1}$$

* Wald's test is also valid if H_0 has the form " $\theta \in \Theta$ " or " $\theta \in \Theta'$ " (but less powerful)

• Wald's test with asymptotic level α (0.05):

$$\Phi = \Phi(T_n > \chi^2_{d-1}) \rightarrow (1-\alpha) \text{-quantile of } \chi^2_{d-1}$$

• Advantage: Non-asymptotic, it can be run on small samples + can always use large sample size

• Disadvantage: assumption that sample is Gaussian

Lecture 15: Goodness of Fit Test for Discrete Distributions

Goodness of fit tests - motivation

e.g. Let X be a r.v. Given n iid copies of X we want to answer

- Does X have distribution $N(0, 1)$?
- Does X have distribution $\text{Bin}(0, 1)$?
- Does X have PAF $p_1 = 0.2, p_2 = 0.3, p_3 = 0.5$?

These are all goodness of fit tests: we want to know if the hypothesized distribution is a good fit to the data

* Key characteristic no parameters

Null hypothesis would be \sim parametric

parametric alternative could be more complicated

• Discrete distribution - Probability simplex

Let $\theta = (a_1, \dots, a_k)$ be a finite set, and (p_1, \dots, p_k) be the family of all probability distributions on θ :

• Probability simplex in \mathbb{R}^k , is the set of all vectors $p = (p_1, \dots, p_k)$ such that

$$p_1 + \dots + p_k = 1, \quad p_i \geq 0 \quad \forall i = 1, \dots, k$$

• Equivalently:

$$a = \{p \mid p = (p_1, \dots, p_k) \in \mathbb{R}^k, \quad p_1, \dots, p_k \geq 0, \quad p_1 + \dots + p_k = 1\}$$

For $p \in \Delta$ and $X \sim p$,

$$p_X(x) = p_j, \quad j=1, \dots, k$$

• Categorical statistical model:

$$(a_1, \dots, a_k), (p_1, \dots, p_k)$$

• Categorical likelihood: a_1, \dots, a_k , outcomes X_1, \dots, X_n , number of occurrences n_1, \dots, n_k

$$L(a_1, \dots, a_k; p_1, \dots, p_k) = \prod_{j=1}^k p_j^{n_j}$$

• Categorical random variable p :

$$P(p=a_1, \dots, a_k) = \prod_{j=1}^k p_j^{n_j}$$

• Categorical distribution:

$$(a_1, \dots, a_k), (p_1, \dots, p_k)$$

• Quantiles and p-values:

For some large integer M :

Simulate M iid copies of T_1, \dots, T_M of T_n

Estimate $(1-\alpha)$ -quantile $q_{1-\alpha}$ of T by taking the sample $(1-\alpha)$ -quantile $q_{1-\alpha} = q_{M, 1-\alpha}$ of T_1, \dots, T_M

Test with approximate level α :

$$\Phi = \Phi(T_n > q_{1-\alpha})$$

Approximate p-value of this test:

$$p\text{-value} = \frac{\#\{T_i > T_n\}}{M}$$

* other goodness of fit tests

We want to know distance of two functions: F and F_0

Holmogorov-Smirnov: $d(F, F_0) = \sup_{x \in \mathbb{R}} |F(x) - F_0(x)|$

Cramér-von Mises: $d(F, F_0) = \int_{-\infty}^{\infty} (F(x) - F_0(x))^2 dF(x)$

Anderson-Darling: $d(F, F_0) = \frac{\int_{-\infty}^{\infty} (x - F(x))^2 dF(x)}{\int_{-\infty}^{\infty} (x - F_0(x))^2 dF_0(x)}$

* Motivation for Composite Goodness of Fit Tests

What if I want to test: "Does X have Gaussian distribution?" but I don't have parameters?

Plug in: $\sup_{\theta \in \Theta} |F_\theta(x) - F(x)|$

$\theta = \theta_0$ → $\hat{\theta}$ → $\hat{F}_{\hat{\theta}}$ → $\hat{F}_{\hat{\theta}}(x)$

Bonker's theorem is no longer valid

Jeffreys prior: $\pi(\theta) \propto \det(\partial \ln L(\theta)/\partial \theta)$ → depends on choice of prior

Another choice: $\pi(\theta) \propto \exp(-\phi(\theta))$ → maximum a posteriori, provided ϕ is superexponential

$\phi(\theta) = \int_{-\infty}^{\infty} (x - \theta)^2 dF(x)$

= conjugate prior

Jeffreys prior satisfies a representation-reversal principle

If π is a reparameterization of $\pi(\theta)$, $\pi(\theta')$ for some one-to-one map ϕ , then the pdf $\pi(\theta')$ of θ' satisfies

$$\pi(\theta') = \pi(\theta) \det(\partial \phi/\partial \theta)$$

further information matrix of statistical model

• Examples: Bernoulli: $\pi(p) \propto p^{-1} (1-p)^{-1}$ → Jeffreys prior

Beta: $\pi(\theta_1, \theta_2) \propto \theta_1^{\theta_1-1} \theta_2^{\theta_2-1}$ → Jeffreys prior

Gaussian: $\pi(\theta) \propto 1/\theta$ → Jeffreys prior

• In general, one can still define a posterior distribution using Bayes's formula

Bayesian prior: $\pi(\theta) \propto f(\theta)$ → Jeffreys prior

Bayesian confidence regions:

* Definition: Bayesian confidence region with level α ($0 < \alpha < 1$) is a random subset

depends on \mathcal{L} & π of parameter space Θ , which depends on sample x_1, \dots, x_n , such that

$P(\theta \in \text{confidence region}) = 1 - \alpha$

* Bayesian confidence region and confidence interval are two distinct notions

Unit 8: Linear Regression

Lecture 19: Linear Regression I

Modeling Assumption

- $(x_i, y_i), i=1, \dots, n$ are iid from some unknown joint distribution P
- P can be described entirely by assuming all exists
- Estimate a joint PDF density
- The marginal density of X is $f_X(x) = \int f_{(X,Y)}(x,y) dy$ and the conditional density $f_{Y|X}(y|x) = \frac{f_{(X,Y)}(x,y)}{f_X(x)}$

Partial modeling

We can also describe the distribution only partially, e.g. using the expectation of Y : $E[Y]$

The conditional expectation of Y given $X=x$: $E[Y|X=x]$ the function:

$$x \mapsto f_{\theta}(x) = E[Y|X=x] = \int y f_{Y|X}(y|x) dy$$

is called regression function

Other possibilities: conditional median,

$$\text{min such that } \int |y - \hat{y}| f_{Y|X}(y|x) dy = \frac{1}{2}$$

conditional quantiles

conditional variance (no information about locations)

Linear regression

Focus on modeling the regression function $f_{\theta}(x) = E[Y|X=x]$

* too many possible regression functions / too parametric

Useful to restrict to simple functions that are described by a few parameters

Simplified: $f_{\theta}(x) = \alpha + \beta x$ linear/slope function

Under this assumption, we discuss linear regression

Bonferroni's test

Test whether a group of explanatory variables is significant in the linear regression

$H_0: \beta_1 = \dots = \beta_k = 0$ vs. $H_1: \beta_j \neq 0$, $j=1, \dots, k$

Bonferroni's test: $R_{\text{test}} = \sqrt{\frac{k}{n}} R_{\text{OLS}}$, where $k=|\beta|$

This test has non-asymptotic level at most α

Remarks:

* Linear regression exhibits correlation NOT causality

* Normality of the noise: One can use goodness-of-fit tests to test whether the residuals $\hat{e}_i = Y_i - \hat{Y}_i$ are Gaussian

* Deterministic design: If X is not deterministic, all the above can be understood conditionally on X , if the noise is assumed to be Gaussian, conditionally on X .

Probabilistic analysis

- Let X and Y be two real 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 $(a^*, b^*) = \arg\min_{(a,b)} E[(Y - a - bx)^2]$

Setting partial derivatives to 0 gives:

$$\begin{aligned} b^* &= \frac{\partial E[(Y-a-bx)^2]}{\partial b} \\ a^* &= E[Y] - b^* E[X] = E[Y] - \frac{\partial E[(Y-a-bx)^2]}{\partial a} \end{aligned}$$

Noise:

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

$$\hat{e} = \begin{cases} E[\hat{e}] = 0 \\ \text{cov}(\hat{e}, \hat{e}) > 0 \end{cases}$$

Least squares

Definition: The least squares estimator (LSE) of (a^*, b^*) is the minimizer of sum of squared errors:

$$\hat{a}^*, \hat{b}^* \in \arg\min_{(a,b)} \sum_{i=1}^n (Y_i - a - bx_i)^2$$

\hat{a}^*, \hat{b}^* is given by

$$\begin{cases} \hat{b}^* = \frac{\sum_i (Y_i - \bar{Y})(x_i - \bar{x})}{\sum_i (x_i - \bar{x})^2} \\ \hat{a}^* = \bar{Y} - \hat{b}^* \bar{x} \end{cases}$$

Residuals: $\hat{e}_i = Y_i - (\hat{a}^* + \hat{b}^* x_i)$

Unit 7: Generalized linear models

Lecture 21: Introduction to Generalized Linear Models, Exponential Families

Linear model

A Gaussian linear model assumes:

$$\begin{aligned} Y &\sim N(\mu, \sigma^2) \\ \text{Response function} \\ E[Y|X=x] &= \mu = \beta_0 + \beta_1 x \end{aligned}$$

Two model components: (μ & σ^2 related in this case)

- Random component: the response variable Y is continuous and $Y|X=x$ is Gaussian with mean $\mu(x)$
- Regression function: $\mu(x) = \beta_0 + \beta_1 x$ (linear)

Generalization

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

- Random component: $Y|X=x$ ~ some distribution (e.g. Bernoulli, exponential, Poisson)
- Regression function: $\mu(x) = g^{-1}(x)$ link function \rightarrow regression function

Exponential family

A family of Distribution (P_θ) , $\theta \in \Theta$, is said to be a k -parameter exponential family on \mathbb{R}^k , if there exist real valued functions:

$$\begin{aligned} f(y, \theta) &\text{ and } g_i(\theta) \text{ and } h_i(y) \\ f(y, \theta) &\text{ such that the density function/pdf of } P_\theta \text{ is} \\ f(y, \theta) = \exp \left(\frac{y \cdot \theta - h(\theta)}{g(\theta)} \right) &= \exp(y \cdot \theta - h(\theta)) \\ g_i(\theta) = \frac{\partial}{\partial \theta_i} \ln(f(y, \theta)) &= \frac{\partial}{\partial \theta_i} \ln \left(\exp(y \cdot \theta - h(\theta)) \right) \\ g_i(\theta) = \theta_i &\Rightarrow g_i(\theta) = \theta_i \end{aligned}$$

Back to β

- Given a link function g , note the following relationship between point θ :

$$\theta = (g^{-1})'(\mu) \quad \theta = (g^{-1})'(\hat{y})$$

where θ is defined as:

$$\theta = (g^{-1})' \circ g^{-1} \circ (g(\cdot))^{-1}$$

Remark: if g is canonical link function, θ is the identity $(g^{-1})' = (g^{-1})^{-1}$

Log-likelihood

The log-likelihood is given by:

$$\ln(Y|X,\theta) = \frac{\sum_i (Y_i - h(\theta))}{g(\theta)} + \text{constant}$$

Note that when using canonical link function, we obtain simpler expression

$$\ln(Y|X,\theta) = \frac{\sum_i (Y_i - \hat{y}_i)}{g(\theta)} + \text{constant}$$

Strictly concavity:

The log-likelihood is strictly concave using the canonical function when $\theta > 0$

As a consequence, μ is unique

If another parameterization is used the likelihood function may not be strictly concave, leading to several local maxima

Concluding remarks

Maximum likelihood for Bernoulli Y and the logit link is called logistic regression

In general, there is no closed form for the ml and we have to use optimization algorithms

The asymptotic normality of MLE also applies to GLMs

Lecture 20: Linear Regression 2

Multivariate case

$$Y = X^T \beta^* + \epsilon, \quad \epsilon \sim N(0, \sigma^2 I_n)$$

Vector of explanatory variables: $X \in \mathbb{R}^{n \times p}$ \rightarrow log. assume its first coordinate Response/dependent variable: $Y \in \mathbb{R}^n$

$\beta^* = (\beta^*, \beta^*)^T$, $\beta^* = \arg\min_{\beta} \text{var}(\epsilon)$

$\{\epsilon_i\}_{i=1, \dots, n}$: noise vectors satisfying $\text{cov}(\epsilon_i, \epsilon_j) = 0$

Definition: the least squares estimator (LSE) of β^* is the minimizer of squared error

$$\hat{\beta} = \arg\min_{\beta} \sum_{i=1}^n (Y_i - \beta^* - X_i \beta)^2$$

LSE in matrix form:

$$5. \quad \hat{\beta} = (X^T X)^{-1} X^T Y$$

$X \in \mathbb{R}^{n \times p}$ matrix whose rows are X_1^T, \dots, X_n^T (design matrix)

$\epsilon \in \mathbb{R}^n$ (unobserved noise)

6. $\beta^* = X^T X \beta^* + X^T \epsilon$ β^* unknown

7. The LSE $\hat{\beta}$ satisfies:

$$\hat{\beta} = \arg\min_{\beta} \|Y - X\beta\|^2$$

8. Closed form solution:

Assume $\text{rank}(X) = p$, analytical computation of the LSE:

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

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

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

10. Test with non-asymptotic level α : $R_{\text{test}} = \sqrt{\frac{n}{k+1}} > \chi_{\alpha, k+1}^2$ $\hat{\beta}$ $\sim \text{N}(0, \sigma^2 X(X^T X)^{-1} X^T)$

Statistical inference

Properties of LSE

1. LSE = MLE

2. Distribution of $\hat{\beta}$: $\hat{\beta} \sim N(\beta^*, \sigma^2 (X^T X)^{-1})$

3. Quadratic risk of $\hat{\beta}$: $E[(\hat{\beta} - \beta^*)^T (\hat{\beta} - \beta^*)] = \sigma^2 \text{tr}(X^T X)$

4. Predicted errors: $\hat{Y} = X\hat{\beta} \sim N(X\beta^*, \sigma^2 I_n)$

5. Unbiased estimator of σ^2 : $\hat{\sigma}^2 = \frac{1}{n-p} \sum_{i=1}^n (Y_i - \hat{Y}_i)^2$ $\hat{\sigma}^2 \sim \text{Chi-squared}_{n-p}$ (Central limit theorem)

Significance tests:

6. Test whether the j -th explanatory variable is significant in the linear regression: $H_0: \beta_j = 0$ vs. $H_1: \beta_j \neq 0$

7. If β_j is the j -th diagonal coefficient of $(X^T X)^{-1}$:

8. Let $T_j = \frac{\hat{\beta}_j}{\hat{\sigma}} \sim t_{n-p}$ (t-test) \rightarrow (1 - α) quantile of t-distr.

9. Variance prof: We have: $\frac{\partial^2 L}{\partial \beta_j^2} = \frac{\partial^2 \ell}{\partial \beta_j^2} = -\frac{1}{\hat{\sigma}^2} (X^T X)_{jj} = -\frac{1}{\hat{\sigma}^2} (X^T X)_{jj} + \left(\frac{\partial^2 \ell}{\partial \beta_j^2} \right)_{\text{true}}$ from previous result

together with second identity, this yields: $\hat{\sigma}^2 = \frac{\partial^2 \ell}{\partial \beta_j^2} + \frac{\partial^2 \ell}{\partial \beta_j^2} \text{ true}$

which leads to: $\text{var}(\hat{\beta}_j) = \frac{\partial^2 \ell}{\partial \beta_j^2} \cdot \hat{\sigma}^2$

Lecture 22: Log-link function and the canonical link function

Link function

1. g is the parameter of interest, needs to appear somewhere in likelihood function, to enable usage of maximum likelihood

2. A link function g relates the linear predictor $X\beta$ to the main parameter, $\mu = g^{-1}(X\beta)$

3. g is required to be monotonically increasing and differentiable

$g'(x) = g''(x) \neq 0$

4. Examples:

① Linear model: $g(x) = \text{identity}$

② Poisson data: suppose $Y|X \sim \text{Poisson}(\mu)$

$\mu = \lambda$

$\log(\mu) = \lambda$

→ In general, a link function for count data should map $(0, \infty)$ to \mathbb{R} , the log-link is a natural one

③ Bernoulli/Binomial data:

o $g(x) = \text{logit}$

g should map $(0, 1)$ to \mathbb{R}

2 choices: $\log(\mu) = \lambda$ or $\log(\mu) = \lambda^2$

public: $\mu = \lambda \exp(\lambda)$ natural logit inverse

o $g(x) = \text{log}$

o $g(x) = \text{logit}$

o $g(x) = \text{log}$

o $g(x) = \text{log}$