

# Statistics Cheatsheet

Hongyi Guo\*

October 21, 2020

## 1 Probabilty & Distribution

**Definition 1** (Population). The entire group of individuals about which we want to learn something about.

**Definition 2** (Sample). Subset of the population from which the information is actually obtained.

**Definition 3** (Statistics). A numerical characteristic of the sample, a random variable.

**Remark 4.** Remarks on the difference between samples and population.

- (a) Population mean  $\mu = \frac{1}{n} \sum_{i=1}^n X_i$ .
- (b) Population variance  $\sigma^2 = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2$ .
- (c) Sample mean  $\bar{X} = \frac{1}{n} \sum_{i=1}^n X_i$ .
- (d) Sample variance  $s^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2$ .

**Definition 5** (SD, SE). SD is the standard deviation, while SE is the standard deviation of the sampling distribution, e.g.,  $SD(\bar{X}) = \frac{\sigma}{\sqrt{n}}$ ,  $SE(\bar{X}) = \frac{\hat{\sigma}}{\sqrt{n}} = \frac{s}{\sqrt{n}}$ .

**Definition 6** (Cdf, pmf, pdf). Let  $X$  be an r.v., then its *cumulative distribution function* (cdf) is defined by  $F_X(x)$ , where

$$F_X(x) = \mathbb{P}_X((-\infty, x]) = \mathbb{P}(X \leq x).$$

For a discrete r.v.  $X$ , the *probability mass function* (pmf) of  $X$  is given by  $p_X(x) = \mathbb{P}[X = x]$ , for  $x \in \mathcal{D}$ . For a continuous r.v.  $X$ , the *probability density function* is given by  $f_X(s) = \frac{d}{ds} F_X(s)$ .

---

\*Northwestern University; [hongyiguo2025@u.northwestern.edu](mailto:hongyiguo2025@u.northwestern.edu)

**Definition 7** (Support). The support of a r.v.  $X$  is the points  $x$  in the space of  $X$  that  $p_X(x) > 0$  (discrete r.v.) or  $f_X(x) > 0$  (continuous r.v.).

**Definition 8** (Point).

**Remark 9.** The cdf, pmf/pdf, mle of various distributions.

Name	Support	pmf/pdf	Mean	Variance	MLE	Fisher
Bernouli	$(0, 1)$	$f(x; p) = \begin{cases} p : & x = 1 \\ 1 - p : & x = 0 \end{cases}$	$p$	$p(1 - p)$	$\hat{p} = \bar{x}$	$\frac{1}{p(1-p)}$
Uniform	$[a, b]$	$f(x; p) = \begin{cases} \frac{1}{b-a} & a \leq x \leq b \\ 0 : & o.w. \end{cases}$	$\frac{a+b}{2}$	$\frac{(b-a)^2}{12}$	$\hat{b} = \max\{x_1, x_2, \dots\}$ $\hat{a} = \min\{x_1, x_2, \dots\}$	-
Normal	$(-\infty, \infty)$	$f(x) = \frac{1}{\sqrt{2\pi}\sigma} \exp \left\{ -\frac{1}{2} \left( \frac{x-\mu}{\sigma} \right)^2 \right\}$	$\mu$	$\sigma$	$\hat{\mu} = \bar{x}$ $\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$	$\begin{bmatrix} \frac{1}{\sigma^2} & 0 \\ 0 & \frac{2}{\sigma^2} \end{bmatrix}$

## 2 Finding Point Estimators and Test Statistics

### 2.1 Point estimator

**Definition 10** (Point estimates). Based on the data sample, come up with the best single guess  $\hat{\theta}$  for the unknown true parameter  $\theta$ .

**Definition 11** (Bias, Var). As follows,

$$\text{Bias}(\hat{\theta}) = \mathbb{E}[\hat{\theta}] - \theta, \quad \text{Var}(\hat{\theta}) = \mathbb{E}[\hat{\theta} - \mathbb{E}[\hat{\theta}]]^2.$$

**Definition 12** (Consistent).  $\hat{\theta}$  is consistent if for any  $\epsilon > 0$ ,  $\lim_{n \rightarrow \infty} \mathbb{P}\{|\hat{\theta} - \theta| \geq \epsilon\} = 0$ .

### 2.2 Pivitol

**Theorem 13.** The CI of  $\bar{X} - \bar{Y}$  is

### 2.3 MLE

**Definition 14** (Likelihood). Likelihood function  $L(\theta; \mathbf{x}) \equiv f(\mathbf{x}; \theta)$  is a function of  $\theta$  for fixed  $\mathbf{x}$ . It's often simpler to maximize log-likelihood:  $\ell(\theta; \mathbf{x}) \equiv \log(L(\theta; \mathbf{x}))$ . When i.i.d.,

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

$$\ell(\theta; \mathbf{x}) = \log \left( \prod_{i=1}^n f(x_i; \theta) \right) = \sum_{i=1}^n \log(f(x_i; \theta)).$$

**Remark 15.** For normal distribution,

$$L(\boldsymbol{\theta}; \mathbf{x}) = f(\mathbf{x}; \boldsymbol{\theta}) = \prod_{i=1}^n f(x_i; \boldsymbol{\theta}) = \frac{1}{(2\pi)^{n/2} \sigma^n} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^n (x_i - \mu)^2 \right\}.$$

**Definition 16** (MLE Principle). Choose estimator of  $\boldsymbol{\theta}$  to maximize  $L(\boldsymbol{\theta}; \mathbf{x})$ :  $\hat{\boldsymbol{\theta}} \equiv \operatorname{argmax}_{\boldsymbol{\theta}} L(\boldsymbol{\theta}; \mathbf{x})$ .

**Remark 17.**  $S^2$  is an unbiased estimator of  $\sigma^2$ .  $S$  is a biased estimator of  $\sigma$ .  $\hat{\sigma}$  is an unbiased estimator of  $\sigma$ .  $\hat{\sigma}^2$  is a biased estimator of  $\sigma^2$ .

**Assumption 18** (Regularity assumptions). These assumptions are

- i) The pdf's are distinct for different  $\boldsymbol{\theta}$ ;
- ii) The pdf's have common support for all  $\boldsymbol{\theta}$ ;
- iii)  $\boldsymbol{\theta}_0$  is in the interior of  $\Omega$ .

**Theorem 19** (Rao-Cramer Lower Bound). Let  $X_1, \dots, X_n$  be i.i.d. with common pdf  $f(x; \theta)$  for  $\theta \in \Omega$ . Assume that the regularity conditions hold. Let  $Y = u(X_1, X_2, \dots, X_n)$  be a statistic with mean  $E(Y) = E[u(X_1, X_2, \dots, X_n)] = k(\theta)$ . Then

$$\operatorname{Var}(Y) \geq \frac{[k'(\theta)]^2}{nI(\theta)}.$$

## 2.4 Method of Moments

**Lemma 20** (Shannon's Lemma). Uniqueness.

## 3 Find confidence regions and critical regions

**Definition 21** (CR). For each  $\boldsymbol{\theta}_0 \in \Omega$ , let  $C(\boldsymbol{\theta}_0) \subset \mathbb{R}^n$  denote the critical region for a size- $\alpha$  test of  $H_0 : \boldsymbol{\theta} = \boldsymbol{\theta}_0$  and a suitable  $H_1$ , and for each  $\mathbf{x}$ ,

$$R(\mathbf{x}) \equiv \{\boldsymbol{\theta}_0 : \mathbf{x} \notin C(\boldsymbol{\theta}_0)\} \subset \mathbb{R}^p.$$

Then  $R(\mathbf{x})$  is a  $1 - \alpha$  confidence region for  $\boldsymbol{\theta}$ .

**Remark 22** (CR). A few remarks on CR.

1. CI is for scalar  $\theta$ . CR is for  $p > 1$  dimensional  $\boldsymbol{\theta}$ .

2. When adopting pivotal quantities to find CRs, in the scalar case, one-sided intervals are unique, but two-sided are not. Taking the shortest length interval is equivalent to requiring that  $\theta$  has constant likelihood on the boundary. In  $p > 1$  case, taking a minimum-volume CR is equivalent to requiring that  $\theta$  has constant likelihood on the boundary.

### 3.1 Asymptotic distribution

**Definition 23** (Score function). The score function is defined as

$$\mathbf{S}(\theta) \equiv \nabla \log f(x; \theta) \equiv \left[ \frac{\partial f(x; \theta) / \partial \theta_1}{f(x; \theta)}, \frac{\partial f(x; \theta) / \partial \theta_2}{f(x; \theta)}, \dots, \frac{\partial f(x; \theta) / \partial \theta_p}{f(x; \theta)} \right]^\top$$

**Definition 24** (Fisher Information, Hessian Matrix). The fisher information is defined as

$$\mathbf{I}(\theta) \equiv -\mathbb{E}_\theta \left[ \frac{\partial^2 \log(f(X; \theta))}{\partial \theta^2} \right] = \mathbb{E}_\theta \left[ \left( \frac{\partial \log f(X; \theta)}{\partial \theta} \right) \left( \frac{\partial \log f(X; \theta)}{\partial \theta} \right)^\top \right] = \text{Cov}_\theta[\nabla \log f(X; \theta)] \geq 0.$$

The Hessian matrix is defined as

$$\mathbf{H}(\theta) \equiv \mathbb{E}_\theta \left[ \frac{\partial^2 \log(f(X; \theta))}{\partial \theta^2} \right] = -\mathbf{I}(\theta).$$

**Remark 25.** Fisher information relates to how accurately we can identify  $\theta$ .  $\mathbf{I}(\theta)$  is inversely proportional to the variance of the MLE.

**Theorem 26.** Let  $X_1, \dots, X_n$  be iid. with pdf  $f(x; \theta)$  for  $\theta \in \Omega$ . Assume the regularity conditions hold. Then

1. The likelihood function  $\frac{\partial}{\partial \theta} l(\theta) = \mathbf{0}$  has a solution  $\hat{\theta}_n$  s.t.  $\hat{\theta}_n \xrightarrow{P} \theta$ .
2. For any sequence which satisfies (1),

$$\hat{\theta}_n \xrightarrow{D} N_p \left( \theta, \frac{\mathbf{I}^{-1}(\theta)}{n} \right).$$

**Theorem 27.** Let  $\mathbf{g}$  be a transformation  $\mathbf{g}(\theta) = (g_1(\theta), \dots, g_k(\theta))^\top$  s.t.  $1 \leq k \leq p$  and that the  $k \times p$  matrix of a partial derivatives  $\mathbf{B} = \left[ \frac{\partial g_i}{\partial \theta_j} \right]$ ,  $i = 1, \dots, k$ ,  $j = 1, \dots, p$  has continuous elements and does not vanish in a neighborhood of  $\theta$ . Let  $\hat{\eta} = \mathbf{g}(\hat{\theta})$ . Then  $\hat{\eta}$  is the mle of  $\eta = \mathbf{g}(\theta)$ , and

$$\hat{\eta} \xrightarrow{D} N_k \left( \eta, \frac{\mathbf{B} \mathbf{I}^{-1}(\theta) \mathbf{B}^\top}{n} \right).$$

Hence,  $\mathbf{I}(\eta) = [\mathbf{B} \mathbf{I}^{-1}(\theta) \mathbf{B}^\top]$ .

**Remark 28.** If expectation is tractable, take it. Substitute  $\hat{\mathbf{I}} = \mathbf{I}(\hat{\boldsymbol{\theta}})$  for  $\mathbf{I}(\boldsymbol{\theta}_0)$  if needed. Otherwise, calculate observed Fisher info matrix.

**Theorem 29** (Find approximate CRs or hyp tests). Individual normal CI on  $\theta_j$ :

$$\hat{\theta}_j \xrightarrow{D} N\left(\theta_j, \frac{[\mathbf{I}^{-1}(\boldsymbol{\theta})]_{j,j}}{n}\right).$$

Then,  $SD(\hat{\theta}_j) = \sqrt{\frac{[\mathbf{I}^{-1}(\boldsymbol{\theta})]_{j,j}}{n}}$ ,  $SE(\hat{\theta}_j) = \sqrt{\frac{[\hat{\mathbf{I}}^{-1}]_{j,j}}{n}}$ , the approx.  $1 - \alpha$  CI is  $\theta_j \in \hat{\theta}_j \pm z_{\alpha/2} SE(\hat{\theta}_j)$ .

**Theorem 30** (Find joint  $\chi^2$  CR on  $\boldsymbol{\theta}$ ). With the construction  $[\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}]^\top [n\mathbf{I}(\boldsymbol{\theta})][\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}] \xrightarrow{D} \chi_p^2$ , an approx.  $1 - \alpha$  joint CR is

$$\left\{ \boldsymbol{\theta} : [\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}]^\top [n\hat{\mathbf{I}}][\hat{\boldsymbol{\theta}} - \boldsymbol{\theta}] \leq \chi_{p,\alpha}^2 \right\}.$$

The CR is an ellipsoid centered at  $\hat{\boldsymbol{\theta}}$ .

**Remark 31.** A few remarks:

1. If  $\mathbf{X} \sim N_k(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ , then  $(\mathbf{X} - \boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{X} - \boldsymbol{\mu}) \sim \chi_k^2$ .
2. Assume  $\boldsymbol{\Sigma}$  has orthonormal eigenvectors  $\mathbf{v}_1, \mathbf{v}_2, \dots$  and eigenvalues  $\lambda_1, \lambda_2, \dots$ , and let  $\mathbf{V} = [\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_k]$  and  $\mathbf{D} = \text{diag}\{\lambda_1, \lambda_2, \dots, \lambda_k\}$ . Then  $\boldsymbol{\Sigma} = \mathbf{V}\mathbf{D}\mathbf{V}^\top = [\mathbf{V}\mathbf{D}^{1/2}][\mathbf{V}\mathbf{D}^{1/2}]^\top = \mathbf{A}\mathbf{A}^\top$ .

### 3.2 Critical Regions from Asymptotic Dist.

**Theorem 32.** Assume that

1.  $X_i : i = 1, 2, \dots, n$  i.i.d.,  $X \sim f(x; \boldsymbol{\theta})$ ,  $n \rightarrow \infty$ , same regularity conditions as for asymptotic distribution of MLE.
2.  $\omega_0 = \{\boldsymbol{\theta} \in \Omega : g_i(\boldsymbol{\theta}) = a_i, i = 1, 2, \dots, q\}$  for some set of  $q \leq p$  smooth independent functions  $g_i(\cdot)$  and constants  $a_i$ , and  $\omega_0$  is in the interior of  $\Omega$ . ( $\omega_0$  is a  $p - q$  dimensional manifold)
3.  $\hat{\boldsymbol{\theta}}_0$  and  $\hat{\boldsymbol{\theta}}$  in the LRT are consistent MLE solutions.

Then, when  $H_0$  is true:

$$-2 \log \Lambda(\mathbf{X}) \xrightarrow{D} \chi_q^2.$$

We reject  $H_0$  if  $-2 \log \Lambda(\mathbf{X}) > \chi_{q,\alpha}^2$ , i.e., reject  $H_0$  if  $\Lambda(\mathbf{X}) < c$  with  $c = \exp \left\{ \frac{-\chi_{q,\alpha}^2}{2} \right\}$ .

**Remark 33.** Test with asymptotic distribution can better control  $\alpha$ -risk.