

① Estimation

(1.1) Statistic: statistic is a RV function of observations X_1, \dots, X_n usually used to estimate some unknown parameter from the underlying probability distribution of the X_i 's.

Ex: $\bar{X} = \sum_{i=1}^n \frac{X_i}{n}$ (sample mean) $S^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{X})^2 = \frac{1}{n-1} \left(\sum_{i=1}^n X_i^2 - n\bar{X}^2 \right)$
(sample variance)

Let X_1, \dots, X_n be iid RV's and let $T(X) \equiv T(X_1, \dots, X_n)$ be a statistic based on the x_i 's (it's also RV, since if we had different samples, we'd expect to get different values for the statistic).

If we use $T(X)$ to estimate some unknown parameter θ , $T(X)$ is called a point estimator for θ .

Ex: \bar{X} is usually a point estimator for the mean $\mu = E[X]$
 S^2 is often a point estimator for the variance $\sigma^2 = \text{var}(X_i)$

Desired Properties for any point estimator

- Unbiasedness ($E[T(X)] = \theta$)
- low variance ("efficiency")

Think: $E[(T(X) - \theta)]^2 = \text{var}(T(X)) + [\text{Bias}(T(X))]^2$ $\left\{ \begin{array}{l} \lim_{n \rightarrow \infty} \text{var}(T(X)) = 0 \\ \lim_{n \rightarrow \infty} \text{Bias}(T(X)) = 0 \end{array} \right\} \begin{array}{l} T(X) \text{ is} \\ \text{consistent} \end{array}$

• \bar{X} is always unbiased for μ : $E[\bar{X}] = E\left[\frac{\sum_{i=1}^n X_i}{n}\right] = \frac{1}{n} \sum_{i=1}^n E[X_i] = \frac{1}{n} n\mu = \mu$

• S^2 is always unbiased for σ^2 : $E[S^2] = E\left[\frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1}\right] = \frac{1}{n-1} E\left[\sum_{i=1}^n (X_i^2 - 2X_i\bar{X} + \bar{X}^2)\right]$
 $= \frac{1}{n-1} E\left[\sum_{i=1}^n X_i^2 - 2\bar{X} \sum_{i=1}^n X_i + n\bar{X}^2\right] = \frac{1}{n-1} E\left[\sum_{i=1}^n X_i^2 - n\bar{X}^2\right]$

Recall
 $E[\bar{X}^2] = \text{var}(\bar{X}) + (E[\bar{X}])^2$
 $= \frac{\sigma^2}{n} + \mu^2$

(using $\sum_{i=1}^n X_i = n\bar{X}$) $= \frac{1}{n-1} \sum_{i=1}^n E[X_i^2] - nE[\bar{X}^2]$
 $= \frac{1}{n-1} \sum_{i=1}^n (\sigma^2 + \mu^2) - n\left(\frac{\sigma^2}{n} + \mu^2\right)$
 $= \frac{1}{n-1} (n\sigma^2 + n\mu^2 - \sigma^2 - n\mu^2) = \frac{1}{n-1} (n\sigma^2 - \sigma^2) = \sigma^2$

(1.2) Mean Squared Error: The MSE of an estimator $T(X)$ of θ is

$$MSE(T(X)) = E[(T(X) - \theta)^2]$$

"Easier" interpretation:

$$\begin{aligned} MSE(T(X)) &= E[T^2] - 2\theta E[T] + \theta^2 \\ &= \underbrace{E[T^2] - (E[T])^2}_{\text{Var}(T)} + \underbrace{(E[T])^2 - 2\theta E[T] + \theta^2}_{(E[T] - \theta)^2} \\ &= \text{Var}(T) + (E[T] - \theta)^2 \end{aligned}$$

And we define $\text{Bias}(T(X)) \equiv E[T(X)] - \theta$

Hence $MSE(T(X)) = \text{Var}(T(X)) + \text{Bias}^2(T(X))$

\Rightarrow MSE combines bias and variance of an estimator

(1.3) Next, we're gonna see two different methods of finding estimators (in Regression, we use Least Squares Est, but let's see 2 different ones):

(1.3.1) Method of Moments Estimator

Recall (from last lecture), the k^{th} Moment of a RV X is

$$E[X^k] = \begin{cases} \sum_x x^k f(x) & \text{if } X \text{ is discrete} \\ \int_{\mathbb{R}} x^k f(x) dx & \text{if } X \text{ is continuous} \end{cases}$$

Now suppose $X_1, \dots, X_n \stackrel{\text{iid}}{\sim}$ some pmf/pdf $f(x)$.

The MoM estimator for $E[X^k]$ is $\sum_{i=1}^n X_i^k / n$

(Does this make sense? Absolutely. Recall that the Law of Large Numbers implies that $\sum_{i=1}^n X_i^k / n \rightarrow E[X^k]$ (as n gets large - LLN))

(Baby) Examples:

• MoM estimator for μ (1st Moment): $\mu = E[X_1]$
 $\hat{\mu}_{\text{MOM}} = \sum_{i=1}^n \frac{X_i}{n} = \bar{X}$

• MoM estimator for σ^2 (variance is the second central moment, $E[(X - \mu)^2]$, but it can also be expressed as:

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

obs: MoM of $E[X^2]$ (2nd Moment) = $\sum_{i=1}^n X_i^2 / n$

$$\text{Hence } \hat{\sigma}_{\text{MOM}}^2 = \frac{1}{n} \sum_{i=1}^n X_i^2 - \left(\frac{1}{n} \sum_{i=1}^n X_i \right)^2 = \left(\frac{1}{n} \sum_{i=1}^n X_i^2 \right) - \bar{X}^2 = \frac{\sum_{i=1}^n X_i^2 - n\bar{X}^2}{n} = \frac{n-1}{n} s^2$$

1.3.2 Maximum Likelihood Estimators (MLE)

Consider an iid random sample X_1, \dots, X_n where each X_i has pdf/pmf $f(x)$. Suppose θ is some unknown parameter from X_i .

The likelihood function is $L(\theta) = \prod_{i=1}^n f(x_i)$

And the MLE of θ is the value of θ that maximizes $L(\theta)$ or in other words

$$\hat{\theta}^{MLE} = \operatorname{argmax}_{\theta} L(\theta)$$

→ In most common applications, optimizing the log-likelihood will make this problem a lot easier. Recall log is a one-to-one, monotonically increasing function.

Some properties (base a is replaced by "e" in MLE → natural log)

$$\begin{aligned} \textcircled{1} \log_a(x \cdot y) &= \log_a x + \log_a y & \textcircled{2} \log_a\left(\frac{x}{y}\right) &= \log_a x - \log_a y & \log_e^x = a \\ \textcircled{3} \log_a x^m &= m \log_a x & \textcircled{4} \log_a \sqrt[n]{x^m} &= \log_a x^{\frac{m}{n}} = \frac{m}{n} \log_a x & e^a = x \end{aligned}$$

Also, recall $\log_e 1 = 0$ and $\frac{d(\log x)}{dx} = \frac{1}{x}$

Example: $X_1, \dots, X_n \stackrel{iid}{\sim} \text{Bern}(p)$. Find \hat{p}^{MLE}

Recall pmf of Bernoulli(p) is $f(x) = \begin{cases} p^x(1-p)^{1-x} & \text{if } x=0 \text{ or } x=1 \\ 0 & \text{otherwise} \end{cases}$

$$\text{The likelihood: } L(p) = \prod_{i=1}^n f(x_i) = \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} = p^{\sum x_i} (1-p)^{n - \sum x_i}$$

$$\text{The log-likelihood: } l(p) = \sum_{i=1}^n x_i \log p + (n - \sum_{i=1}^n x_i) \log(1-p)$$

Take partial derivative with respect to p and set to 0.

$$\frac{\partial l(p)}{\partial p} = \frac{\sum x_i}{p} - \frac{n - \sum x_i}{1-p} \equiv 0$$

$$\Leftrightarrow \frac{\sum x_i}{p} = \frac{n - \sum x_i}{1-p} \Rightarrow \sum x_i - p \sum x_i - pn + p \sum x_i = 0$$

$$\hat{p}^{MLE} = \frac{\sum x_i}{n} = \bar{X} //$$

② Sampling Distribution: statistics (recall some examples \bar{X} & s^2) are RV's. So it's sometimes useful to figure out their distributions, which are called "sampling distribution".

Ex: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \Rightarrow \bar{X} \sim N(\mu, \sigma^2/n)$

2.1 Chi-squared distribution (χ^2)

Def: if $Z_1, \dots, Z_K \stackrel{iid}{\sim} N(0,1)$ then $Y = \sum_{i=1}^K Z_i^2$ has a χ_K^2 dist

$\therefore Y \sim \chi^2(K)$ (chi-squared dist. with K degrees of freedom)

and pdf is $f_Y(y) = \frac{1}{2^{K/2} \Gamma(K/2)} y^{K/2-1} e^{-y/2}$, for $y > 0$

$E[Y] = K$, $\text{var}(Y) = 2K$
 gamma function

Property: χ^2 's add up. Let $Y_1, \dots, Y_n \stackrel{iid}{\sim} \chi^2(d_i) \forall i$
 then $\sum_{i=1}^n Y_i \sim \chi^2(\sum_{i=1}^n d_i)$

Usually χ^2 dist comes up when we try to estimate σ^2 , for example,
 if $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \rightarrow s^2 = \frac{\sum_{i=1}^n (X_i - \bar{X})^2}{n-1} \sim \frac{\sigma^2}{n-1} \chi^2(n-1)$

2.2 t distribution

Def: suppose that $Z \sim N(0,1)$, $Y \sim \chi^2(K)$ and Z and Y are indep.

$T \equiv \frac{Z}{\sqrt{Y/K}}$ has the student t dist. with K degrees of freedom

$\therefore T \sim t(K)$ with pdf $f_T(x) = \frac{\Gamma(\frac{K+1}{2})}{\sqrt{\pi K} \Gamma(\frac{K}{2})} \left(\frac{x^2}{K} + 1 \right)^{-\frac{K+1}{2}}$, $x \in \mathbb{R}$

$\rightarrow k=1$ gives Cauchy dist, with really fat tails

$\rightarrow t(K) \rightarrow N(0,1)$ as K becomes large

$E[T] = 0$ and $\text{var}(T) = \frac{K}{K-2}$ ($K > 2$)

Usually the t dist. is used for finding confidence intervals and conduct hypothesis tests for the mean μ .

2.3 F distribution

Def: Suppose that $X \sim \chi^2(n)$, $Y \sim \chi^2(m)$ and X, Y are indep.

Then $F \equiv \frac{mX}{nY}$ has the F dist with n and m degrees of freedom

$$\therefore F \sim F(n, m) \text{ with pdf } f_F(x) = \frac{\Gamma\left(\frac{n+m}{2}\right) \left(\frac{n}{m}\right)^{\frac{n}{2}} x^{\frac{n}{2}-1}}{\Gamma\left(\frac{n}{2}\right) \Gamma\left(\frac{m}{2}\right) \left(\frac{n}{m}x + 1\right)^{\frac{n+m}{2}}}, x > 0$$

F distribution is used to find confidence intervals and conduct hypothesis tests for the ratio of variances for 2 processes.

(Remember in Regression when we analyze the ANOVA Table, we get F-stats for testing the equality of treatment means, i.e.

$$F = \frac{MSTreatment}{MSError}, \text{ which comes from } \begin{cases} E(MSTr) = \sigma^2 + \beta_1^2 \sum_{i=1}^2 (x_i - \bar{x})^2 \\ E(MSE) = \sigma^2 \end{cases}$$

$$\hookrightarrow H_0: \beta_1 = 0 \text{ vs } H_1: \beta_1 \neq 0$$

2.4 Normal (or Gaussian) distribution

Def: $X \sim N(\mu, \sigma^2)$ if it has pdf $f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}}, x \in \mathbb{R}$

$$E[X] = \mu \text{ and } \text{var}(X) = \sigma^2$$

Thm: Additive property of normals: if $X_1, \dots, X_n \stackrel{\text{iid}}{\sim} N(\mu_i, \sigma_i^2), i=1, \dots, n$

$$\text{then } Y \equiv \sum_{i=1}^n a_i X_i + b \sim N\left(\sum_{i=1}^n a_i \mu_i + b, \sum_{i=1}^n a_i^2 \sigma_i^2\right)$$

(i.e. a linear combination of normals is itself normal)

\Rightarrow Fact: a normal dist. is completely characterized by its mean and variance (if we add up normals, the result is still normal, just need to compute new mean and variance).

Ex: $X \sim N(3, 4)$ and $Y \sim N(4, 6)$, X, Y are indep. Find dist of $2X - 3Y$

$$E[2X - 3Y] = 2E[X] - 3E[Y] = -6 \quad \text{var}(2X - 3Y) = 4\text{var}(X) + 9\text{var}(Y) = 70$$

$$\Rightarrow 2X - 3Y \sim N(-6, 70)$$

2.4.1 Standard Normal

Def: The Normal $(0,1)$ distribution is called standard normal, often denoted by Z (tables available for cdf).

Any normal can be standardized by applying the transformation

$$Z = (X - \mu) / \sigma$$

pdf of $N(0,1)$ (usually represented by ϕ)

$$\phi(z) = \frac{1}{\sqrt{2\pi}} e^{-z^2/2}, \quad z \in \mathbb{R}$$

cdf of $N(0,1)$ (usually represented by Φ)

$$\Phi(z) = \int_{-\infty}^z \phi(t) dt, \quad z \in \mathbb{R}$$

Remarks

1. $P(Z \leq a) = \Phi(a)$
2. $P(Z \geq b) = 1 - \Phi(b)$
3. $P(a \leq Z \leq b) = \Phi(b) - \Phi(a)$
4. $\Phi(0) = 0.5$
5. $\Phi(-b) = P(Z \leq -b) = P(Z \geq b) = 1 - \Phi(b)$
6. $P(-b \leq Z \leq b) = \Phi(b) - \Phi(-b) = 2\Phi(b) - 1$

Then

$P(\mu - k\sigma \leq X \leq \mu + k\sigma) = P(-k \leq Z \leq k) = \Phi(k) - \Phi(-k) = 2\Phi(k) - 1$
 \hookrightarrow probability that any normal RV is within k std. dev of μ doesn't depend on μ or σ^2 .

Ex: for $k=1$ standard deviation, $2\Phi(1) - 1 = 2(0.8413) - 1 = 0.6826$.

for $k=2$, $2\Phi(2) - 1 = 0.95$

for $k=3$, $2\Phi(3) - 1 = 0.997$

Example: $X \sim N(21, 4)$. Find $P(19 < X < 22.5)$

$$P(19 < X < 22.5) = P\left(\frac{19-21}{\sqrt{4}} < Z < \frac{22.5-21}{\sqrt{4}}\right)$$

$$= P(-1 < Z < 0.75)$$

$$= \Phi(0.75) - \Phi(-1)$$

$$= \Phi(0.75) - [1 - \Phi(1)]$$

$$= 0.7734 - [1 - 0.8413] = 0.6147$$

2.4.2 Some theorems / corollaries:

→ Corollary: $X_1, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2) \Rightarrow \bar{X} \sim N(\mu, \frac{\sigma^2}{n})$
 (dist of sample mean has normal dist and as the number of observations increases, $\text{var}(\bar{X})$ gets smaller).

→ Central Limit Theorem: suppose X_1, \dots, X_n iid with $E[X_i] = \mu$ and $\text{var}(X_i) = \sigma^2$. Then as $n \rightarrow \infty$

$$\frac{\sum_{i=1}^n X_i - n\mu}{\sigma\sqrt{n}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \xrightarrow{d} N(0, 1)$$

Remarks ① If n is large (usually $n \geq 30$ will do the job) $\rightarrow \bar{X} \approx N(\mu, \frac{\sigma^2}{n})$

② The X_i 's don't need to be normal for the CLT to work. In fact, CLT works even on discrete distributions.

③ You can almost always use the CLT if the observations are iid.

Ex: $X_1, \dots, X_{100} \stackrel{iid}{\sim} \exp(1/1000)$. Find $P(950 \leq \bar{X} \leq 1050)$

First: since X_i 's are $\exp(\frac{1}{1000})$: $E[X_i] = \frac{1}{\lambda}$ $\text{var}(X_i) = \frac{1}{\lambda^2}$

By CLT: $E[\bar{X}] = E[X_i] = \frac{1}{\lambda} = 1000$

$$\text{Var}(\bar{X}) = \text{var}(X_i)/n = 10000$$

$$\text{Hence } P(950 \leq \bar{X} \leq 1050) = P\left(\frac{950 - 1000}{\sqrt{10000}} \leq Z \leq \frac{1050 - 1000}{\sqrt{10000}}\right)$$

$$\approx P\left(-\frac{1}{2} \leq Z \leq \frac{1}{2}\right) = 2\Phi(0.5) - 1 = 0.383$$

So the above probability can be approximated by Normal dist.
 (The exact answer can be obtained by using Erlang distribution).

Summary

$$t_n = \frac{Z}{\sqrt{\frac{\chi_n^2}{n}}} \quad \text{indep}$$

$$F_{m,n} = \frac{\chi_m^2}{\chi_n^2} \quad \text{indep}$$

$$Z \sim N(0,1)$$

 $z_i\text{'s} \stackrel{iid}{\sim} N(0,1)$

$$\chi_n^2 = \sum_{i=1}^n z_i^2$$

③ Confidence Intervals

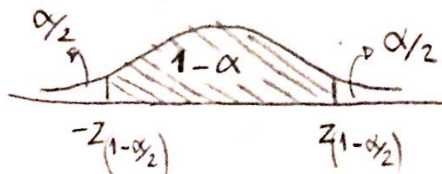
③.1 Quantiles: the $(1-\alpha)$ quantile of a RV X is the value of x_α such that $P(X > x_\alpha) = 1 - F(x_\alpha) = \alpha$. Note that $x_\alpha = F^{-1}(1-\alpha)$ where $F^{-1}(\cdot)$ is the inverse cdf of X .

③.2 Confidence interval: instead of estimating a parameter by a point estimator give a (random) interval that contains the unknown parameter with a certain probability: there's a $(1-\alpha)$ chance that the parameter lies between two lower and upper limits.

• Two-sided CI: most times, we're interested in two-sided CI

$\Rightarrow X_1, \dots, X_n \sim N(\mu, \sigma^2)$, σ^2 known

$$Z = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}}$$



$$-z_{1-\alpha/2} \leq \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \leq z_{1-\alpha/2}$$

$(1-\alpha)$ CI is then:

$$\bar{X} \pm \underbrace{z_{1-\alpha/2} \frac{\sigma}{\sqrt{n}}}_{\text{half-width}}$$

$\Rightarrow X_1, \dots, X_n \sim N(\mu, \sigma^2)$, σ^2 unknown
($n \leq 30$)

$$t = \frac{\bar{X} - \mu}{s/\sqrt{n}} = \frac{\bar{X} - \mu}{\sigma/\sqrt{n}} \cdot \frac{\sigma}{s} = Z \cdot \frac{1}{\sqrt{\frac{s^2}{\sigma^2}}} = Z \cdot \frac{1}{\sqrt{\frac{(n-1)s^2}{\sigma^2} \cdot \frac{1}{n-1}}} = \frac{Z}{\sqrt{\frac{\chi^2_{n-1}}{n-1}}} = t_{n-1}$$

$(1-\alpha)$ CI is

$$\bar{X} \pm t_{n-1, 1-\alpha/2} \frac{s}{\sqrt{n}}$$

• One-sided CI (just interested in one bound)

$\Rightarrow 100(1-\alpha)\%$ upper CI for μ

$$\mu \leq \bar{X} + z_\alpha \sqrt{\sigma^2/n}$$

$\Rightarrow 100(1-\alpha)\%$ lower CI for μ

$$\mu \geq \bar{X} - z_\alpha \sqrt{\sigma^2/n}$$

(use α quantile, not $\frac{\alpha}{2}$)