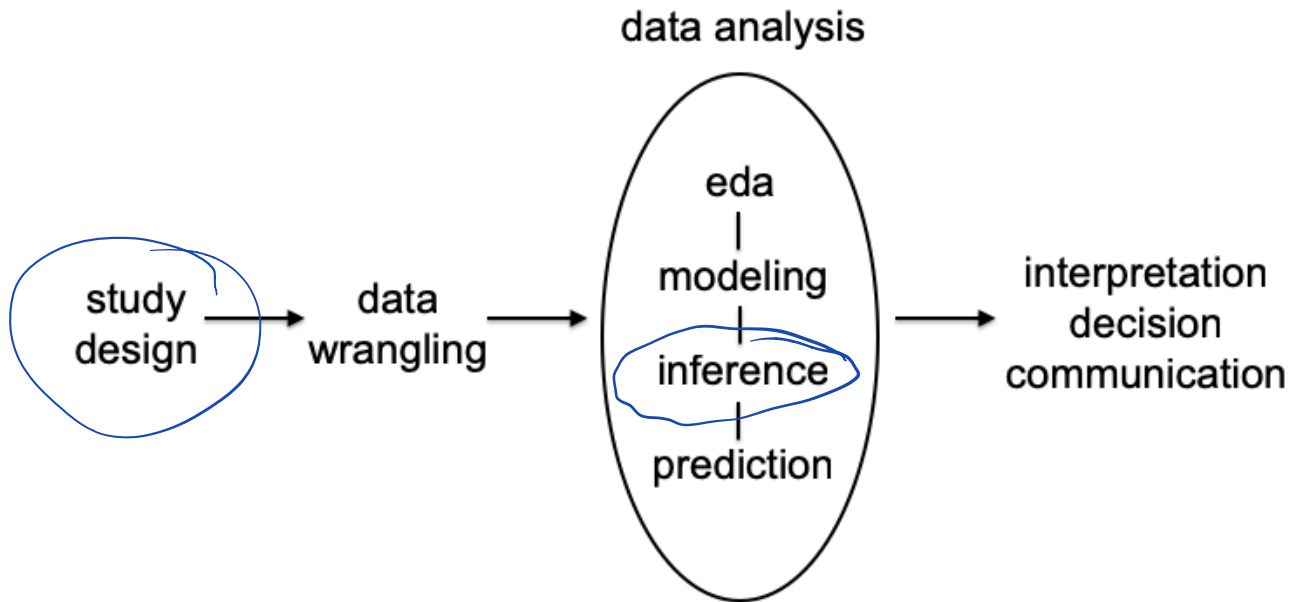
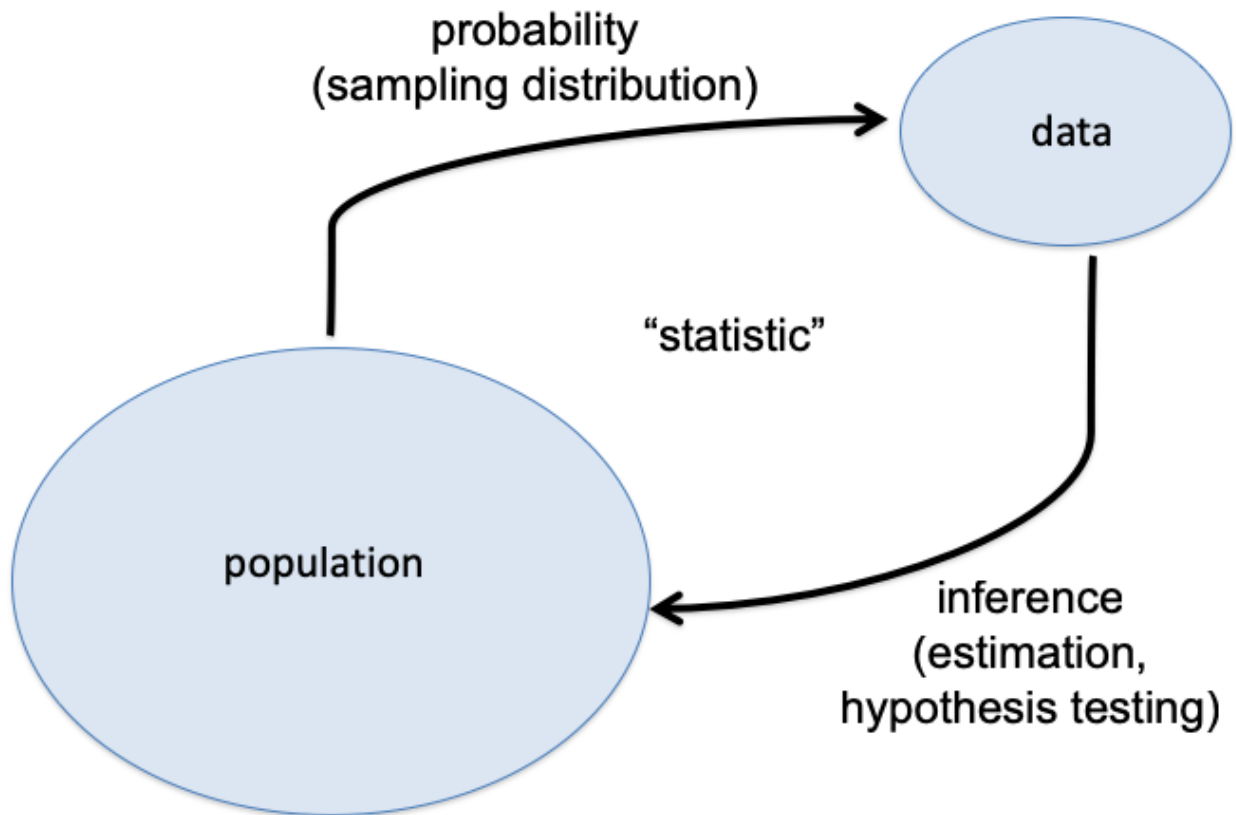


This Week's Topic: Frequentist Statistical Inference



Central Dogma of (Frequentist) Statistical Inference



Observing/data according to a
Collecting/ probabilistic model
Generating

Example: Simple Random Sample

Units are uniformly and independently sampled from a population:

- Political survey (survey samples)
- Flipping a coin n times

Example: Randomized controlled study

Units are randomly sampled and then randomized to two or more treatment groups

- Clinical trial

Example: Fair Coin?

Suppose I claim that a specific coin is fair, i.e., that it lands on heads or tails with equal probability.

I flip it 20 times and it lands on heads 16 times.

\hat{p} \sim p

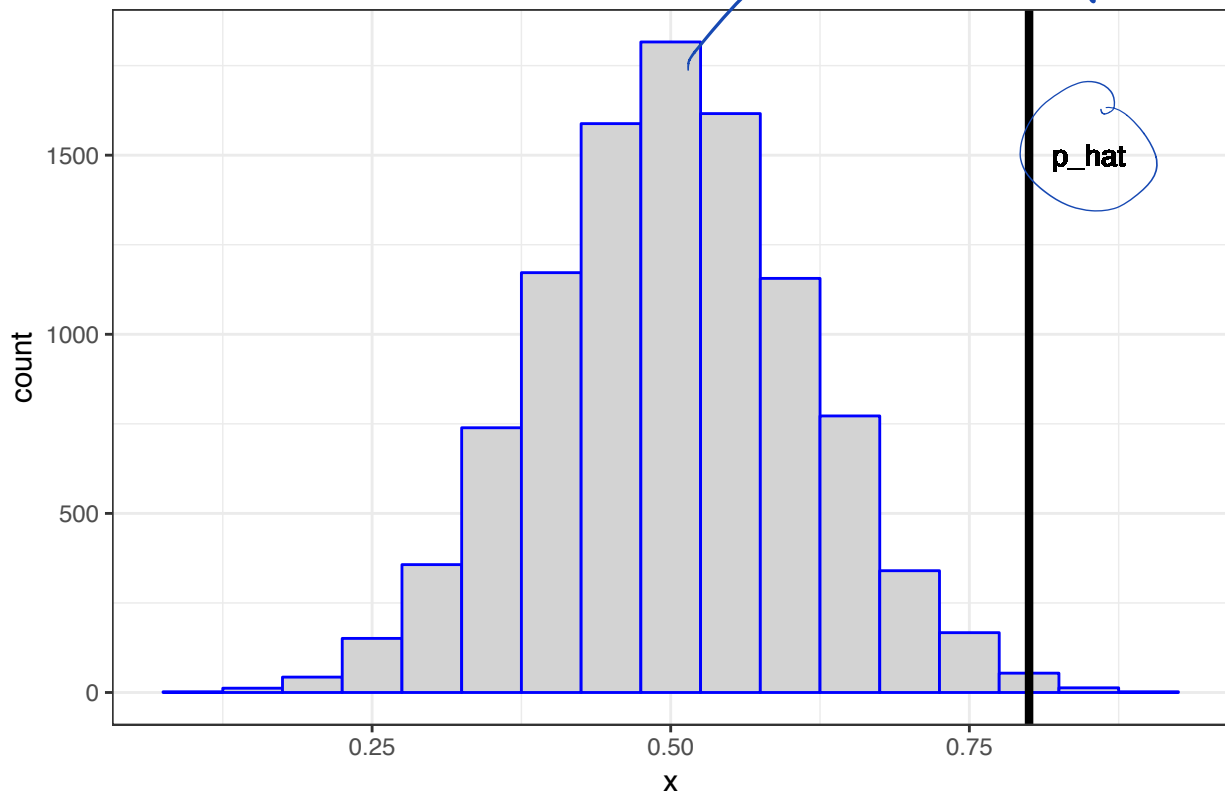
1. My data is $x = 16$ heads out of $n = 20$ flips.
2. My data generation model is $X \sim \text{Binomial}(20, p)$.
3. I form the statistic $\hat{p} = 16/20$ as an estimate of p .

Let's simulate 10,000 times what my estimate would look like if $p = 0.5$ and I repeated the 20 coin flips over and over.

```
> x <- replicate(n=1e4, expr=rbinom(1, size=20, prob=0.5))
> sim_p_hat <- x/20
> my_p_hat <- 16/20
```

What can I do with this information?

Histogram of Sampling Distribution



```
> sum(abs(sim_p_hat-0.5) >= abs(my_p_hat-0.5))/1e4
[1] 0.0083
```

Parameter a number that describes a population.

- Usually fixed
- We don't know its value
- appears in the probability model of how we collected data

Statistic a number calculated from a sample of data

A statistic is used to estimate a parameter

Sampling distribution of a statistic is the probability dist'n of the statistic under repeated realizations of the data from the assumed data generating probability dist'n

Goals of Inference

- ① Form estimates of the parameters
- ② Quantity uncertainty about the estimates
- ③ Test hypotheses on the parameters

Let's go through these goals
in simple scenario:

Data generating process is

$$X_1, X_2, \dots, X_n \stackrel{i.i.d.}{\sim} \text{Normal}(\mu, \sigma^2)$$

μ is unknown

σ^2 is known

$$\text{① } \sum_{i=1}^n X_i \sim \text{Normal}(n\mu, n\sigma^2)$$

② If $Z \sim \text{Normal}$ then
 $a + bZ \sim \text{Normal} \Rightarrow$

$$\bar{X} = \frac{\sum_{i=1}^n X_i}{n} \sim \text{Normal}(\mu, \sigma^2/n)$$

Exercise: verify this

③
$$\frac{\bar{X} - \mu}{\sqrt{\sigma^2/n}} \sim \text{Normal}(0, 1)$$

Point estimate of μ :

$$\hat{\mu} = \bar{X}$$

$$\hat{\mu} \sim \text{Normal}(\mu, \sigma^2/n)$$

↑
sampling distribution

$$\frac{\hat{\mu} - \mu}{\sqrt{\sigma^2/n}} \sim \text{Normal}(0, 1)$$

Pivotal Statistic a statistic whose sampling distribution does not depend on any unknown parameters

$$\frac{\hat{\mu} - \mu}{\sqrt{\sigma^2/n}} \text{ is pivotal}$$

Confidence Interval

Interval of the form

$$(\hat{\mu} - c_l, \hat{\mu} + c_u) \checkmark$$

$$c_l, c_u \geq 0$$

where

$$\Pr(\mu - c_l \leq \hat{\mu} \leq \mu + c_u)$$

forms the "level" or coverage probability of the interval

If $Z \sim \text{Normal}(0, 1)$ then

$$\Pr(-1.96 \leq Z \leq 1.96) = 0.95$$

$$\text{Let } Z = \frac{\hat{\mu} - \mu}{\sqrt{\sigma^2/n}}$$

$$\Pr\left(-1.96 \leq \frac{\hat{\mu} - \mu}{\sqrt{\sigma^2/n}} \leq 1.96\right) = 0.95$$

$$= \Pr\left(-1.96 \sqrt{\sigma^2/n} \leq \hat{\mu} - \mu \leq 1.96 \sqrt{\sigma^2/n}\right)$$

$$= \Pr\left(-1.96 \sqrt{\sigma^2/n} \leq \mu - \hat{\mu} \leq 1.96 \sqrt{\sigma^2/n}\right)$$

$$= \Pr\left(\hat{\mu} - 1.96 \sqrt{\sigma^2/n} \leq \mu \leq \hat{\mu} + 1.96 \sqrt{\sigma^2/n}\right)$$

$\Rightarrow C_L = C_U = 1.96 \sqrt{\sigma^2/n}$ gives
me a 95% confidence
interval

Confidence Interval Simulation

```
> mu <- 5
> n <- 20
> x <- replicate(10000, rnorm(n=n, mean=mu)) # 10000 studies
> m <- apply(x, 2, mean) # the estimate for each study
> ci <- cbind(m - 1.96/sqrt(n), m + 1.96/sqrt(n))
> head(ci)
      [,1]      [,2]
[1,] 4.797848 5.674386
[2,] 4.599996 5.476534
[3,] 4.472930 5.349468
[4,] 4.778946 5.655485
[5,] 4.778710 5.655248
[6,] 4.425023 5.301561

> cover <- (mu > ci[,1]) & (mu < ci[,2])
> mean(cover)
[1] 0.9512
```

Let z_α be the α -percentile
of $\text{Normal}(0,1)$.

If $Z \sim \text{Normal}(0,1)$ then

$$Pr(Z \leq z_\alpha) = \alpha$$

$(1-\alpha)$ level CI:

$$\left(\hat{\mu} - |z_{\alpha/2}| \frac{\sigma}{\sqrt{n}}, \hat{\mu} + |z_{\alpha/2}| \frac{\sigma}{\sqrt{n}} \right)$$

$$z_{\alpha/2} = -z_{1-\alpha/2}$$

When $\alpha = 0.05$ then

$$z_{\alpha/2} = -1.96, \quad z_{1-\alpha/2} = 1.96$$

$(1-\alpha)$ -level upper CI:

$$\left(-\infty, \hat{\mu} + |z_{\alpha}| \frac{b}{\sqrt{n}}\right)$$

$1-\alpha$
percentile

$(1-\alpha)$ -level lower CI:

$$\left(\hat{\mu} - |z_{\alpha}| \frac{b}{\sqrt{n}}, \infty\right)$$

α percentile

Hypothesis Testing

coin example I did a hypothesis test of $p = 0.5$ vs. $p \neq 0.5$

hypothesis test / significance test

is a formal procedure for

Comparing observed data with a hypothesis whose truth we want to assess

The results of a test are expressed in terms of how well the data and one of the hypotheses agree

Null hypothesis (H_0) is the statement being tested

Alternative hypothesis (H_1) is the opposite of the null, and it's the "interesting" state

H_0 and H_1 are defined in terms of parameter values (or probabilistic property)

Examples:

two-sided

$$H_0: \mu = 5 \quad H_1: \mu \neq 5$$

one-sided

$$H_0: \mu \leq 5 \quad H_1: \mu > 5$$

$$H_0: \mu \geq 5 \quad H_1: \mu < 5$$

test statistic a statistic designed to quantify evidence against the H_0 in favor of H_1 .

$$H_0: \mu = 5 \quad H_1: \mu \neq 5$$

$$|Z| = \left| \frac{\bar{X} - 5}{\sqrt{\sigma^2/n}} \right|$$

$$\hat{\mu} = \bar{X} \quad \frac{\hat{\mu} - \mu}{\sqrt{\sigma^2/n}} \sim N(0,1)$$

If $H_0: \mu = 5$ is true then

$$Z = \frac{\bar{X} - 5}{\sqrt{\sigma^2/n}} \sim \text{Normal}(0, 1)$$

is pivotal

Collect my n data points
and calculate my observed

statistic:

$$z = \frac{\bar{x} - 5}{\sqrt{\sigma^2/n}}$$

a number

Recall larger $|z|$ is, the
more evidence against H_0 in
favor of H_1 .

p-value is the probability of observing a test statistic "as or more extreme" than the observed statistic under the sampling distribution of the statistic when H_0 is true

Let $Z^* \sim \text{Normal}(0, 1)$

$$p\text{-value} = \Pr(|Z^*| \geq |z|)$$

↑ ↑
theoretical observed

Make a decision based on the p-value. Smaller it is, the more evidence against H_0 in favor of H_1 there is.

Let's say we call a test
"significant" if $p\text{-value} \leq \alpha$.

Type I error or false positive:

Call test significant, i.e.,
reject H_0 in favor of H_1
when H_0 is actually true

Type II error or false negative

Fail to call test significant
when H_1 is actually true

"Rate" is the probability of
these errors under a decision
rule.

If $p\text{-value} \leq \alpha$ is the rule
then false positive rate is α

Under two-sided test, $\mathbb{P} \leq \alpha$
Under a one-sided test

Exercise: convince yourself

P^* be the p -value under repeated studies. Sampling dist'n of P^* when H_0 is true is Uniform(0,1) for a two-sided test.

Show $\Pr(P^* \leq t; H_0 \text{ true}) = t$
for $t \in [0, 1]$.

$\Pr(P^* \leq t; H_0 \text{ true}) \leq t$
for one-sided

Power : Probability of a significant test when H_1 is true

$$\text{Power} = 1 - \text{false negative rate}$$

Power β calculate under a range of alternative parameter values

$$X: \Omega \longrightarrow \mathbb{R}$$

$$(\Omega, \mathcal{F}, P)$$

$$\mathcal{R} = \{X(\omega) : \omega \in \Omega\}$$

\mathcal{R} discrete or continuous

$f(x)$ p.m.f. or p.d.f.

$$\sum_{x \in \mathcal{R}} f(x) = 1 \quad \text{or} \quad \int_{x \in \mathcal{R}} f(x) dx = 1$$

Joint Distributions

Distribution of two or more random variables

Bivariate joint dist'n:

rv's X, Y

have pmf or pdf $f(x, y)$

discrete

$$f(x, y) = \Pr(X=x, Y=y)$$

$$= \Pr(\{\omega: X(\omega)=x\} \cap \{\omega: Y(\omega)=y\})$$

analogous for pdf's

$A_x \subseteq \mathbb{R}, A_y \subseteq \mathbb{R}$ then

$\Pr(X \in A_x, Y \in A_y)$ is:

$$\sum_{x \in A_x} \sum_{y \in A_y} f(x, y)$$

discrete

$$\int_{x \in A_x} \int_{y \in A_y} f(x, y) dx dy$$

continuous

Bivariate cdf:

$$F(a, b) = P_r(X \leq a, Y \leq b)$$

Marginal Dist'n

$$f(x) = \sum_{y \in R_y} f(x, y)$$

$$f(x) = \int_{y \in R_y} f(x, y) dy$$

Independence of X and Y

$$f(x, y) = f(x) f(y)$$

Conditional dist'n's

$X | Y = y$

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

$$\sum_{x \in \mathcal{R}_x} f(x|y) = 1$$

$$\int_{x \in \mathcal{R}_x} f(x|y) dx = 1$$

Law of
Total
Variance

$$E[X^k | Y=y] = \sum_{x \in \mathcal{R}_x} x^k f(x|y)$$

or

$$= \int_{x \in \mathcal{R}_x} x^k f(x|y) dx$$

General Joint Distributions

$$X_1, X_2, \dots, X_n$$

$$f(x) = f(x_1, x_2, \dots, x_n)$$

$$x = (x_1, x_2, \dots, x_n)$$

If the r.v.'s are independent

$$\text{then } f(x) = \prod_{i=1}^n f(x_i)$$

Likelihood

$X_1, X_2, \dots, X_n \sim \text{Bernoulli}(p)$

write their joint pmf as

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

Generically write θ as the parameter.

Joint pmf or pdf

$$f(x; \theta)$$

$$\prod_{i=1}^n f(x_i; \theta) \quad (\text{independence})$$

For observed data x_1, x_2, \dots, x_n

I can calculate $f(x; \theta)$

\Rightarrow this a function of θ

Likelihood

$$L(\theta; \mathbf{x}) = f(\mathbf{x}; \theta)$$

viewed as a function of θ

for observed $\mathbf{x} = (x_1, x_2, \dots, x_n)$

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

log Likelihood

$$\log(L(\theta; \mathbf{x})) = l(\theta; \mathbf{x})$$

independence \Rightarrow

$$l(\theta; \mathbf{x}) = \sum_{i=1}^n l(\theta, x_i)$$

Sufficient statistic

A sufficient statistic $T(\mathbf{x})$ is such that $\mathbf{x} | T(\mathbf{x})$ does not depend on θ

If $f(x; \theta) = g(T(x); \theta) h(x)$
then $T(x)$ is sufficient

$$L(\theta; x) = g(T(x), \theta) h(x) \\ \propto L(\theta; T(x))$$

Example $X_1, X_2, X_3, \dots, X_n \stackrel{iid}{\sim} N(\mu, \sigma^2)$

$T(x) = \bar{X}$ is sufficient for μ

$$\bar{X} \sim \text{Normal}(\mu, \sigma^2/n)$$

Likelihood Principle

Suppose x and y are two
data sets such that

$$L(\theta; x) \propto L(\theta; y)$$

$$\text{i.e. } L(\theta; x) = L(\theta; y) c(x, y)$$

Then inference on θ should be the same for x and y .

Maximum Likelihood Estimation

Estimate θ as the value that maximizes $L(\theta; x)$

$$\begin{aligned}\hat{\theta}_{MLE} &= \operatorname{argmax}_{\theta} L(\theta; x) \\ &= \operatorname{argmax}_{\theta} l(\theta; x) \checkmark \\ &= \operatorname{argmax}_{\theta} L(\theta; T(x))\end{aligned}$$

Example:

$X \sim \text{Binomial}(n, p)$

$$\begin{aligned}L(p; x) &= \binom{n}{x} p^x (1-p)^{n-x} \\ &\propto p^x (1-p)^{n-x}\end{aligned}$$

$$l(p; x) \propto x \log(p) + (n-x) \log(1-p)$$

- $\frac{d}{dp} l(p; x)$ set it to 0

- solve for p .

$$\Rightarrow \hat{p} = \frac{x}{n}$$

It happens to be the case that
true parameter

MLE $\frac{\hat{p} - p}{\sqrt{\frac{\hat{p}(1-\hat{p})}{n}}} \sim \text{Normal}(0, 1)$
standard error with \hat{p} plugged

Approx 95% CI :

$$\left(\hat{p} - 1.96 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}}, \hat{p} + 1.96 \sqrt{\frac{\hat{p}(1-\hat{p})}{n}} \right)$$

$$\text{Var}(X) = np(1-p)$$

$$\begin{aligned}\text{Var}\left(\frac{X}{n}\right) &= \frac{1}{n^2} \text{Var}(X) \\ &= \frac{np(1-p)}{n^2} \\ &= \frac{p(1-p)}{n}\end{aligned}$$

Properties of MLEs

Under certain "regularity conditions", we have the following:

Consistent: $\hat{\theta}_n$ MLE n observations

$$\hat{\theta}_n \xrightarrow{P} \theta$$

For all $\varepsilon > 0$,

$$\Pr(|\hat{\theta}_n - \theta| > \varepsilon) \rightarrow 0$$

as $n \rightarrow \infty$

• Equivariant :

If $\hat{\theta}$ is MLE for θ

then $g(\hat{\theta})$ is MLE $g(\theta)$

$$\left(\begin{array}{l} \hat{p} \text{ MLE } p \\ \frac{\hat{p}(1-\hat{p})}{n} \text{ MLE } \frac{p(1-p)}{n} \end{array} \right)$$

• Asymptotically Normal distributed

• Asymptotically Efficient

Take any estimate $\tilde{\theta}_n$

$$\frac{\text{Var}(\tilde{\theta}_n)}{\text{Var}(\hat{\theta}_n)} \rightarrow \text{to a quantity} \leq 1.$$

Fisher Information

Let $X_1, X_2, \dots, X_n \stackrel{iid}{\sim} F_\theta$

$$I_n(\theta) = \text{Var} \left(\frac{d}{d\theta} \log f(X; \theta) \right)$$

$$= \sum_{i=1}^n \text{Var} \left(\frac{d}{d\theta} \log f(X_i; \theta) \right)$$

$$= -E \left[\frac{d^2}{d\theta^2} \log f(X; \theta) \right]$$

$$= -\sum_{i=1}^n E \left[\frac{d^2}{d\theta^2} \log f(X_i; \theta) \right]$$

$$\text{Var}(\hat{\theta}_n) \approx \frac{1}{I_n(\theta)}$$

Standard Error

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

$$se(\hat{\theta}_n) \approx \frac{1}{\sqrt{I_n(\theta)}}$$

$$\hat{se}(\hat{\theta}_n) = \frac{1}{\sqrt{I_n(\hat{\theta}_n)}}$$

MLE CLT :

$$\frac{\hat{\theta}_n - \theta}{se(\hat{\theta}_n)} \xrightarrow{D} \text{Normal}(0,1)$$

$$\frac{\hat{\theta}_n - \theta}{\hat{se}(\hat{\theta}_n)} \xrightarrow{D} \text{Normal}(0,1)$$

$$Z = \frac{\hat{\theta}_n - \theta}{\hat{se}(\hat{\theta}_n)} \text{ is approx pivoted Normal}(0,1)$$

Repeat our special case
where $\hat{\mu} \leftarrow \hat{\theta}$

$$\sqrt{\frac{\sigma^2}{n}} \leftarrow \hat{se}(\hat{\theta})$$