# Interval estimation

# Motivation

Suppose we have data $(X_1, Y_1), \ldots, (X_n, Y_n)$ with

$$Y_i \sim Bernoulli(p_i)$$

$$\log\left(\frac{p_i}{1 - p_i}\right) = \beta^T X_i$$

So far, we have discussed:

+ Finding point estimates $\widehat{\beta}$

+ Testing hypotheses about the true (but unknown) parameters $\beta$

> What are the limitations of point estimates and hypothesis tests for inference about $\beta$?

# Confidence interval

```
...
##               Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.6415063  0.1213233   21.77   <2e-16 ***
## WBC         -0.2892904  0.0134349  -21.53   <2e-16 ***
## PLT         -0.0065615  0.0005932  -11.06   <2e-16 ***
## ---
...
```

> How would I calculate a 95% confidence interval for $\beta_1$ (the change in the log odds of dengue for a one-unit increase in WBC, holding PLT fixed)?

$$\hat{\beta}_1 \quad \pm \quad Z_{\frac{\alpha}{2}} \; SE(\hat{\beta}_1) \qquad \leftarrow 1-\alpha \text{ Wald CI}$$

$$95\% \; CI: \qquad -0.289 \quad \pm 1.96 \, (0.0134)$$

$$(-0.315, \; -0.262)$$

# Confidence interval

$$P\left(\beta_1 \in (-0.315, -0.262)\right)$$

$$= \begin{cases} 1 & \text{(if interval contains } \beta_1\text{)} \\ 0 & \text{(if interval does not contain } \beta_1\text{)} \end{cases}$$

```
...
##             Estimate Std. Error z value Pr(>|z|)
## (Intercept)  2.6415063  0.1213233   21.77   <2e-16 ***
## WBC         -0.2892904  0.0134349  -21.53   <2e-16 ***
## PLT         -0.0065615  0.0005932  -11.06   <2e-16 ***
## ---
...
```

95% confidence interval for $\beta_1$: (-0.315, -0.262)

How do I interpret this confidence interval?

95% confident: if we take many samples and we calculate an interval from each sample, 95% of those intervals should contain the true (unknown) parameter.

Let L be the lower endpoint, U be the upper endpoint (random variables that are functions of the sample):

$$P\left(\beta_1 \in (L, U)\right) = 0.95 \qquad \text{(for 95\% interval)}$$

$$\hat{\theta} \sim N(\theta, Var(\hat{\theta})) \implies \frac{\hat{\theta} - \theta}{SE(\hat{\theta})} \sim N(0,1)$$

# Deriving the coverage probability

$$P\left( -z_{\frac{\alpha}{2}} \leq \frac{\hat{\theta} - \theta}{SE(\hat{\theta})} \leq z_{\frac{\alpha}{2}} \right) = 1-\alpha$$

$$P\left( -z_{\frac{\alpha}{2}} SE(\hat{\theta}) \leq \hat{\theta} - \theta \leq z_{\frac{\alpha}{2}} SE(\hat{\theta}) \right) = 1-\alpha$$

$$\implies P\left( \hat{\theta} - z_{\frac{\alpha}{2}} SE(\hat{\theta}) \leq \theta \leq \hat{\theta} + z_{\frac{\alpha}{2}} SE(\hat{\theta}) \right) = 1-\alpha$$

endpoints of $1-\alpha$ Wald CI: $\hat{\theta} \pm z_{\frac{\alpha}{2}} SE(\hat{\theta})$

(Also true: $P\left( \theta - z_{\frac{\alpha}{2}} SE(\hat{\theta}) \leq \hat{\theta} \leq \theta + z_{\frac{\alpha}{2}} SE(\hat{\theta}) \right) = 1-\alpha$

but we can't calculate these endpoints! )

<u>Note</u>: for any particular value $\theta_0$ of $\theta$

$\theta_0 \in \left( \hat{\theta} - z_{\alpha/2} SE(\hat{\theta}), \hat{\theta} + z_{\alpha/2} SE(\hat{\theta}) \right)$ if and only if

$\left| \frac{\hat{\theta} - \theta_0}{SE(\hat{\theta})} \right| \leq z_{\alpha/2}$   i.e. fail to reject

$H_0: \theta = \theta_0$   $H_A: \theta \neq \theta_0$

<u>Summarize</u> :

$$\theta_0 \in \left( \hat{\theta} - z_{\frac{\alpha}{2}} SE(\hat{\theta}) , \hat{\theta} + z_{\frac{\alpha}{2}} SE(\hat{\theta}) \right)$$

if <u>and only if</u> fail to reject $H_0 : \theta = \theta_0$

(vs. $H_A : \theta \neq \theta_0$ )

To test a hypothesis $H_0 : \theta = \theta_0$ vs. $H_A : \theta = \theta_0$ (at level $\alpha$):

1) Construct a $1-\alpha$ CI for $\theta$

2) Check if $\theta_0 \in$ CI

(But, we don't get a p-value)

To create a $1-\alpha$ CI for $\theta$ :

1) Test $H_0 : \theta = \theta_0$ vs. $H_A : \theta \neq \theta_0$ for all $\theta_0$

(at level $\alpha$)

2) $1-\alpha$ CI = { all $\theta_0$ for which we fail to reject }

$$X_1, \ldots, X_n \overset{iid}{\sim} \quad \text{with} \quad \text{mean } \mu \; \& \; \text{variance } \sigma^2$$

$$\text{Var}(\bar{X}_n) \;=\; \frac{\sigma^2}{n} \;=\; \frac{\text{Var}(X_i)}{n}$$

$$\text{SE}(\bar{X}) \;=\; \frac{\sigma}{\sqrt{n}}$$

# Formal definition

Let $\Theta \in \boxed{H}$ be a parameter of interest, and $X_1, \dots, X_n$ a sample. Let $C(X_1, \dots, X_n) \subseteq \boxed{H}$ be a set constructed from $X_1, \dots, X_n$ ($\Rightarrow C(X_1, \dots, X_n)$ is a random set).

$C(X_1, \dots, X_n)$ is a $1-\alpha$ <u>confidence set</u> for $\Theta$ if

$$\inf_{\Theta \in \boxed{H}} P_\Theta \left( \Theta \in C(X_1, \dots, X_n) \right) = 1-\alpha$$

$$\left( \forall \ \Theta, \ P_\Theta \left( \Theta \in C(X_1, \dots, X_n) \right) \geq 1-\alpha \right)$$

# Inverting a test

Theorem: Let $\theta \in \boxed{\mathcal{H}}$ be a parameter of interest.

For each value of $\theta_0 \in \boxed{\mathcal{H}}$, consider testing $H_0: \theta = \theta_0$ vs. $H_A: \theta \neq \theta_0$, and let $R(\theta_0)$ be the rejection region for a level $\alpha$ test of these hypotheses.

Let $C(X_1, \ldots, X_n) = \{ \theta_0 : (X_1, \ldots, X_n) \notin R(\theta_0) \}$

Then $C(X_1, \ldots, X_n)$ is a $1-\alpha$ confidence set for $\theta$

# Example

Suppose $X_1, \ldots, X_n \overset{iid}{\sim} Uniform[0, \theta]$. We want to test

$$H_0 : \theta = \theta_0 \quad H_A : \theta \neq \theta_0$$

Find the LRT statistic for this test.

# Example

Suppose $X_1, \ldots, X_n \overset{iid}{\sim} Uniform[0, \theta]$. Inverting the LRT gives us a confidence interval of the form

$$C(X_1, \ldots, X_n) = \left\{ \theta : X_{(n)} \leq \theta \leq X_{(n)} k' \right\}$$

Find a value $k'$ such that the test is size $\alpha$.