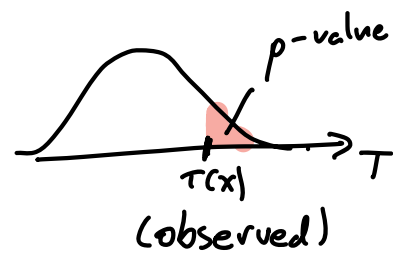# $p$-Values, Confidence Regions

## Outline

1) $p$-Values

2) Confidence regions

3) (Mis-)interpreting tests

# $p$-Values

<u>Informal definition:</u> Suppose $\phi(X)$ rejects for large values of $T(X)$.

$$p(x) = {}^{``}\mathbb{P}_{H_o}\left(T(X) \geq T(x)\right){}^{"}$$

$$= \sup_{\theta \in \Theta_o} \mathbb{P}_\theta\left(T(X) \geq T(x)\right)$$



$p$-value

$T(x)$ (observed)

$T$

<u>Ex</u>   $X \sim N(\theta, 1)$    $H_o: \theta = 0$   vs.   $H_1: \theta \neq 0$

Two-sided test rejects for large $T(X) = |X|$

$$\left(\iff \phi_\alpha(x) = \mathbb{1}\{|X| > z_{\alpha/2}\}\right)$$

The two-sided $p$-value is $p(X)$ where

$$p(x) = \mathbb{P}_o\left(|X| > |x|\right)$$

$$= 2\left(1 - \Phi(|x|)\right)$$

# Formal definition : $\mathcal{P}, \textcircled{$\Theta$}_o, \textcircled{H}$.

Assume we have a test $\phi_\alpha$ for each significance level, $\sup\limits_{\theta \in \textcircled{$\Theta$}_o} \mathbb{E}_\theta \phi_\alpha(X) \le \alpha$

(non-randomized case : $\phi_\alpha = 1\{x \in R_\alpha\}$ )

Assume tests are monotone in $\alpha$ :

if $\alpha_1 \le \alpha_2$ then $\phi_{\alpha_1}(x) \le \phi_{\alpha_2}(x)$

(non-randomized : $R_{\alpha_1} \subseteq R_{\alpha_2}$)

Then $p(x) = \inf\{\alpha : \phi_\alpha(x) = 1\}$

$(= \inf\{\alpha : x \in R_\alpha\})$

<span style="color:red">(possible to define randomized p-value but not worth it)</span>

Note $p(x) \le \alpha \iff \phi_{\tilde\alpha}(x) = 1 \quad \forall \tilde\alpha > \alpha$

For $\theta \in \textcircled{H}_o$, $\mathbb{P}_\theta(p(x) \le \alpha) \le \inf\limits_{\tilde\alpha > \alpha} \underbrace{\mathbb{P}_\theta(\phi_{\tilde\alpha}(X) = 1)}_{\color{red}{\le \tilde\alpha}} \le \alpha$

$\implies$ p-value <u>stochastically dominates</u> $u[0,1]$

If $\phi_\alpha$ rejects for large $T(X)$,

reduces to original definition.

Note the p-value depends on

- the model & null hyp.,
- the data, AND
- the choice of test

Ex    $X \sim N_d(\theta, I_d)$    $H_o : \theta = 0$ vs $H_1 : \theta \neq 0$

We can use    $T_1(x) = \|X\|^2$    ($\chi^2$ test)

or    $T_2(x) = \|X\|_\infty$
$= \max_i |X_i|$    (max test)

Very different p-values / power if d large
(choice reflects belief about whether $\theta$ is sparse)

# Confidence Sets

[Accept/reject decision only so interesting:
- usually we care how big $\theta$ is
- tiny p-value doesn't imply big $\theta$
  (big p-value doesn't imply small $\theta$ either)]

<u>Def</u> $\mathscr{P} = \{P_\theta : \theta \in \Theta\}$

$C(X)$ is a <u>$1-\alpha$ confidence set</u> for $g(\theta)$ if

$$P_\theta \left( C(X) \ni g(\theta) \right) \geq 1-\alpha, \quad \forall \theta \in \Theta$$

subject ↑ verb ↑ object ↑

We say $C(X)$ <u>covers</u> $g(\theta)$ if $C(X) \ni g(\theta)$

$P_\theta(C(X) \ni g(\theta))$ is <u>coverage probability</u>

$\inf_\theta P_\theta(C \ni g(\theta))$ is <u>conf. level</u>

<u>Notes</u>
- $C(X)$ is random, not $g(\theta)$
- Often misinterpreted as Bayesian guarantee
- Say "$C(X)$ has a 95% chance of covering"
  NOT "$g(\theta)$ has a 95% chance of being in $C$"
  NEVER "95% chance $g(\theta) \in [0.5, 1.5]$" (e.g.)

# Duality of Testing & Confidence Sets

Suppose we have a level-$\alpha$ test $\phi(x; a)$
  of $H_0: g(\theta) = a$    vs.    $H_1: g(\theta) \neq a$,   $\forall a \in g(\Theta)$

We can use it to make a confidence set for $g(\theta)$:

Let    $C(X) = \{a : \phi(x; a) < 1\}$

$\qquad\qquad = $ "all non-rejected values of $\theta$"

Then    $\mathbb{P}_\theta \left( C(X) \not\ni g(\theta) \right) = \mathbb{P}_\theta \left( \phi(x; g(\theta)) = 1 \right)$

$$\leq \alpha \qquad\qquad \forall \theta$$

Alternatively, suppose   $C(X)$ is a $1 - \alpha$ confidence set for $g(\theta)$.

We can use $C$ to construct a test $\phi(x)$ of

$\qquad H_0: g(\theta) = a$    vs.    $H_1: g(\theta) \neq a$

$\qquad \phi(x) = 1\{a \notin C(X)\}$

$\qquad$ For $\theta$ s.t. $g(\theta) = a$:

$\qquad \mathbb{E}_\theta \, \phi(x) = \mathbb{P}_\theta \left( C(X) \not\ni g(\theta) \right) \leq \alpha$

This is called <u>inverting the test</u>

# Confidence Intervals / Bounds

If $C(x) = [C_1(x), C_2(x)]$ we say $C(x)$ is a <u>confidence interval</u> (c<u>I</u>)

$C(x) = [C_1(x), \infty)$: <u>lower conf. bd.</u> (L<u>CB</u>)

$C(x) = (-\infty, C_2(x)]$: <u>upper conf bd.</u> (u<u>CB</u>)

We usually get LCB / UCB by inverting a one-sided test in appropriate direction

Called <u>uniformly most accurate</u> (u<u>MA</u>) if test UMP

Get CI by inverting a two-sided test

Called UMAU if test is UMPU

**Ex**   $X \sim Exp(\theta) = \frac{1}{\theta} e^{-x/\theta}$     $x > 0, \theta > 0$

CDF $\mathbb{P}_\theta(X \leq x) = 1 - e^{-x/\theta}$

**LCB:** Invert test for $H_0: \theta \leq \theta_0$

Solve $\alpha = \mathbb{P}_{\theta_0}(X > c(\theta_0)) = e^{-c(\theta_0)/\theta_0}$

$c(\theta_0) = \theta_0 \log(1/\alpha)$  $(> 0)$

$X \leq c(\theta_0) \implies \theta_0 \geq \frac{X}{-\log \alpha}$

$C(X) = \left[\frac{X}{-\log \alpha}, \infty\right)$

**UCB:** Similar, $C(X) = \left(-\infty, \frac{X}{-\log(1-\alpha)}\right]$

**Equal-tailed CI:**

Invert equal-tailed test of $H_0: \theta = \theta_0$

$\underbrace{\phi_\alpha^{ET}(x)}_{\substack{\text{equal-tailed} \\ H_0: \theta = \theta_0}} = \underbrace{\phi_{\alpha/2}^{\geq \theta_0}(x)}_{H_0: \theta \geq \theta_0} + \underbrace{\phi_{\alpha/2}^{\leq \theta_0}(x)}_{H_0: \theta \leq \theta_0}$

$\implies C(X) = \left[\frac{X}{-\log \frac{\alpha}{2}}, \infty\right) \cap \left(-\infty, \frac{X}{-\log(1-\alpha/2)}\right]$

$= \left[\frac{X}{-\log \frac{\alpha}{2}}, \frac{X}{-\log(1-\frac{\alpha}{2})}\right]$

Similar for UMPU 2-sided test

# (Mis-) Interpreting Hypothesis Tests

Hypothesis tests ubiquitous in science

Common misinterpretations:

1) $p < 0.05$ therefore "there is an effect"
   or "the effect size = the estimate"

2) $p > 0.05$ therefore "there is no effect"

3) $p = 10^{-6}$ therefore "the effect is huge"

4) $p = 10^{-6}$ therefore "the data are signif."
   and everything about our model
   is correct in most naive interp.

5) Effect CI for men is $[0.2, 3.2]$,
   for women is $[-0.2, 2.8]$ therefore
   "there is an effect for men and not
   for women."

Dichotomous test doesn't eliminate uncertainty
   (CIs usually less misleading to novices)

# How to interpret testing

Learning about the world from data
is not easy or automatic!

Hypothesis tests let us ask specific
questions about specific data sets
under specific modeling assumptions,
using specific testing method.

All of these choices bear on the interpretation.

Top-tier medical journals let people
publish claims, reporting p-values
without saying what model was used or
what test was employed

Pretty bad when you think about it!

Hyp. tests can be a good companion to
critical thinking, never a substitute

"All models are wrong, some are useful" but
need experience and theory to understand
when assumptions do or don't cause real trouble

# Conceptual Objections

**Q1:** Why should I test $H_0: \theta = 0$? No $\theta$ is ever exactly $0$.

**A1:** a) Test $H_0: |\theta| \leq \delta$ if you want

If $s.e.(\hat{\theta}) >> \delta$, not much difference.

b) Most two-sided tests justify <u>directional inference</u>:

"If $T > c_a$ declare $\theta > 0$, if $T < c$, declare $\theta < 0$" with $\mathbb{P}(\text{false claim}) \leq \alpha$

c) Harder to answer in non-parametric problems,
e.g. $H_0: P = Q$ vs $H_1: P \neq Q$ for perm. test, but alternative frameworks like Bayes force very strong assumptions on us.

**Q2:** People only like frequentist results like p-values, CIs because they mistake them for Bayesian results.

95% chance $C(x) \ni \theta$ is misinterpreted as a claim about $p(\theta | X)$.

**A2:** True, but subjective Bayesian results often misinterpreted as "<u>the</u> posterior dist. of $\theta$" when really should be "<u>my</u> posterior opinion about $\theta$"