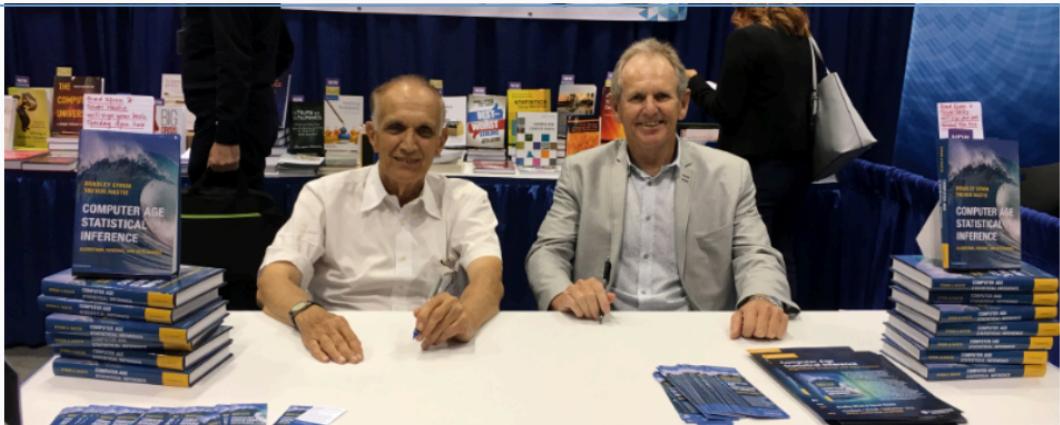


Mathematical Statistics I

STA2212H S LEC9101

Week 2

January 20 2021



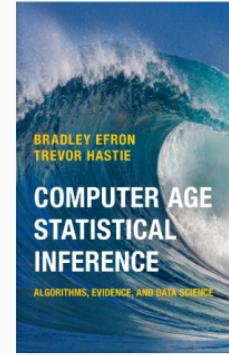


Volume 88, Issue S1

Special Issue: "Data Science versus Classical Inference: Prediction, Estimation, and Attribution", honouring Prof. Brad Efron's International Prize in Statistics in 2019

Pages: S1-S224

December 2020



The Computer Age Statistical Inference book makes the distinction between the two levels of statistics, the algorithmic level and the inferential level, which is somewhat an artificial distinction but a pretty good one. It says that the first level is doing something and the second level is understanding what you did in the first level. The algorithmic level always gets more action, in particular in these days of these big prediction algorithms like deep learning. You'd think that's the only thing going on. It isn't the only thing going on. The deeper understanding of the kind of thing that Fisher and these people – Neyman, Hotelling – did for early 20th-century statistics, putting it on a solid intellectual ground so you can understand what's at stake, is terribly important.

Recap

- likelihood notation $\ell'(\theta) = \mathcal{U}(\theta), s(x, \theta)$

notes on likelihood

- score function, maximum likelihood estimate, observed and expected Fisher information $\hat{j}(\hat{\theta}) \hat{I}(\hat{\theta}) \underline{I}(\theta)$

- asymptotic normality of maximum likelihood estimators

$$\sqrt{n}(\hat{\theta} - \theta) I_1^{1/2}(\hat{\theta}) \xrightarrow{d} N(0, 1)$$

- estimating the asymptotic variance

$$\hat{\ell}'(\hat{\theta}) = j(\hat{\theta}), I_n(\hat{\theta})$$

$$\tau = g(\theta)$$

- the delta method

see notes p.6

- profile likelihood

- sufficient statistics $t(x)$

$$\ell(\theta; t(x)) + \dots$$

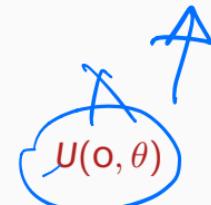
- Newton-Raphson method for computing $\hat{\theta}$

- irregular models

- Quasi-Newton \leftarrow

- EM Algorithm

1st order
??
? proof?



Friday

1. Quasi-Newton
2. Hypothesis testing AoS 10.1
3. Significance testing SM 7.3.1; AoS 10.2
4. Tests based on likelihood AoS 10.6

- **January 25 3.00 – 4.00 Aleeza Gerstein** Data Science and Applied Research Series
- “Turning qualitative observation to quantitative measurement through statistical computing” [Link](#)



Notes on optimization: Tibshirani, Pena, Kolter CO 10-725 CMU

- Goal: $\max_{\theta} \ell(\theta; \mathbf{x})$

- Solve: $\ell'(\hat{\theta}; \mathbf{x}) = 0$ mle

- Iterate: $\hat{\theta}^{t+1} = \hat{\theta}^t + \frac{\ell'(\hat{\theta}^t)}{-\ell''(\hat{\theta}^t)}$

- Rewrite:

- Quasi-Newton: $\begin{bmatrix} \hat{\mathbf{J}} \end{bmatrix}^{-1} (\hat{\theta}^{t+1} - \hat{\theta}^t) = \ell'(\hat{\theta}^t)$

-

- $(-\Delta \ell)^{-1} [\nabla \ell]$

$$\frac{\partial}{\partial \theta} \ell(\theta; \mathbf{x}) \quad \hat{\theta}(\mathbf{x})$$

$$= \hat{\theta}^t + [-\ell''(\hat{\theta}^t)]^{-1} \ell'(\hat{\theta}^t)$$

$p \times p \qquad p \times 1$

replace $\hat{\mathbf{J}}(\hat{\theta}^t)^{-1}$ by
by an approx=

? optim(par, fn, gr = NULL, ...,

method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN", "Brent"),
lower = -Inf, upper = Inf, control = list(), hessian = FALSE)

doesn't use 2nd deriv.

Notes on optimization: Tibshirani, Pena, Kolter CO 10-725 CMU

- Goal: $\max_{\theta} \ell(\theta; \mathbf{x})$
- Solve: $\ell'(\theta; \mathbf{x}) = 0$
- Iterate: $\hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + \{j(\hat{\theta}^{(t)})\}^{-1} \ell'(\hat{\theta}^{(t)})$
- Rewrite: $j(\hat{\theta}^{(t)})(\hat{\theta}^{(t+1)} - \hat{\theta}^{(t)}) = \ell'(\hat{\theta}^{(t)})$
- Quasi-Newton:
 - approximate $j(\hat{\theta}^{(t)})$ with something easy to invert \leftarrow
 - use information from $j(\hat{\theta}^{(t)})$ to compute $j(\hat{\theta}^{(t+1)}) \leftarrow$
- optimization notes add a step size to the iteration $\hat{\theta}^{(t+1)} = \hat{\theta}^{(t)} + \epsilon_t \{j(\hat{\theta}^{(t)})\}^{-1} \ell'(\hat{\theta}^{(t)})$

```
optim(par, fn, gr = NULL, ...,
      method = c("Nelder-Mead", "BFGS", "CG", "L-BFGS-B", "SANN", "Brent"),
      lower = -Inf, upper = Inf, control = list(), hessian = FALSE)
```

- Null and alternative hypothesis

- Rejection region: $R \subseteq \mathcal{X}^n$

- Test statistic and critical value

- Type I and Type II error

- Power and Size

$$X_1, \dots, X_n = \underline{x}_n \sim f(\underline{x}; \theta)$$

$$\theta \in \Theta \subseteq \mathbb{R}^k$$

\mathcal{X} sample space

could be infinite

$f(\underline{x})$ "smooth"

non-par

$$\underline{x}_i = m(z_i) + \varepsilon_i$$

"smooth"

$$H_0: \theta \in \Theta_0$$

\equiv

null

$$H_1: \theta \in \Theta_1$$

alternative

↑
come back
later

if $\underline{x} \in R$: "reject H_0 "
 $\underline{x} \notin R$: "don't reject H_0 "
 (retain null)

- Null and alternative hypothesis

• Rejection region

• Test statistic and critical value

• Type I and Type II error

• Power and Size

$$R \subseteq \mathcal{X} \quad \underline{x} \in \mathcal{X}$$

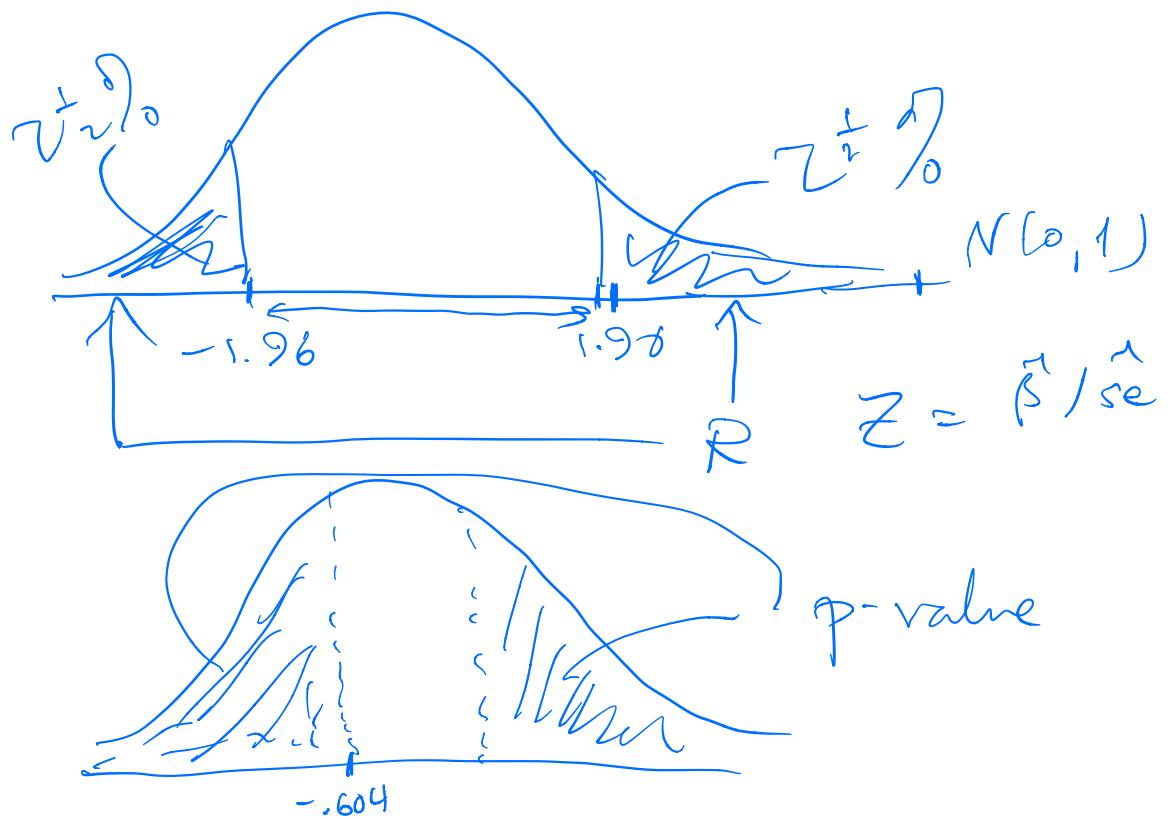
$$R = \{x: t(x_n) > c\} \quad \begin{array}{l} \text{1-dim. summary} \\ \text{c - to be det'd} \end{array}$$

$$\Pr_{H_0}\{X^n \in R: \theta \in \Theta_0\} = [\Pr_{H_0}\{X^n \in R\}]^{\text{size}}$$

$$\Pr_{H_A}\{X^n \notin R: \theta \in \Theta_A\} = \Pr_{H_A}\{X^n \notin R\} \quad \text{1-pwr}$$

$$\alpha = \Pr_{\Theta_0}\{X^n \in R\} \quad \text{pr (type 1 error)}$$

$$\beta = 1 - \text{pr (type 2 error)}$$



Example: logistic regression

Coefficients:

	Estimate	Std. Error	z value	Pr(> z)	
(Intercept)	-34.103704	6.530014	-5.223	1.76e-07 ***	
zn	-0.079918	0.033731	-2.369	0.01782 *	
indus	-0.059389	0.043722	-1.358	0.17436	
chas	<u>-0.785327</u>	<u>0.728930</u>	<u>1.077</u>	<u>0.28132</u>	
nox	48.523782	7.396497	6.560	5.37e-11 ***	
rm	-0.425596	0.701104	<u>-0.607</u>	0.54383	$H_0: \beta_5 = 0$ vs. $H_1: \beta_5 \neq 0$
age	0.022172	0.012221	1.814	0.06963 .	
dis	0.691400	0.218308	3.167	0.00154 **	
rad	0.656465	0.152452	4.306	1.66e-05 ***	
tax	-0.006412	0.002689	-2.385	0.01709 *	
ptratio	0.368716	0.122136	3.019	0.00254 **	
black	-0.013524	0.006536	-2.069	0.03853 *	
lstat	0.043862	0.048981	0.895	0.37052	
medv	0.167130	0.066940	2.497	0.01254 *	

Signif. codes: 0 ‘***’ 0.001 ‘**’ 0.01 ‘*’ 0.05 ‘.’ 0.1 ‘ ’ 1

\downarrow p-value $N(0,1)$

	Kreject H_0	don't
H_0	type I error	✓
H_1	✓	type II error

if $p < .05$ ← common
then \textcircled{X}

“statistically
significant”

$$P_{\beta_5=0} (N(0,1) > .607) \quad 0 < p < .001 < p < .01 < p < .05 < p < .1 < p \leq 1$$

$$p < .05 \text{ if } |z| > 1.96 \quad z = \hat{\beta} / \hat{s}_\beta$$

... Example: logistic regression

```
Boston.glmnull <- glm(crim2 ~ 1, family = binomial, data = Boston)
```

```
anova(Boston.glmnull, Boston.glm)
```

Analysis of Deviance Table

$$H_0: \beta_0 = 0$$

$$\boldsymbol{\beta} = (\beta_0, \beta_1, \dots, \beta_p) = \{f_0, \beta_{c,1}, \dots\}$$

Model 1: crim2 ~ 1

Model 2: crim2 ~ (crim + zn + indus + chas + nox + rm + age + dis + rad + tax + ptratio + black + lstat + medv) - crim

$$\text{crim2} \sim \cdot - \text{crim}$$

	Resid. Df	Resid. Dev	Df	Deviance
--	-----------	------------	----	----------

1	505	701.46		
---	-----	--------	--	--

2	492	211.93	13	489.54
---	-----	--------	----	--------

resid. dev. 701
212

```
> pchisq(489.54, 13, lower.tail = F)  
[1] 2.435111e-96
```

$$\text{diff} \sim \chi^2_{13}$$

LRT (need to prove) (coming)

... Example: logistic regression

```
Boston.glmpart <- glm(crim2 ~ . - crim - indus - chas - rm - lstat,  
                      data = Boston, family = binomial)
```

```
anova(Boston.glmpart, Boston.glm)
```

Analysis of Deviance Table

Model 1: crim2 ~ (crim + zn + indus + chas + nox + rm + age + dis + rad +
tax + ptratio + black + lstat + medv) - crim - indus - chas -
rm - lstat

Model 2: crim2 ~ (crim + zn + indus + chas + nox + rm + age + dis + rad +
tax + ptratio + black + lstat + medv) - crim

Resid.	Df	Resid.	Df	Deviance
--------	----	--------	----	----------

1	496	216.22		
---	-----	--------	--	--

2	492	211.93	4	4.2891
---	-----	--------	---	--------

```
> pchisq(4.2891, 4, lower.tail = F)
```

```
[1] 0.368292
```

$$H_0: \beta_{(\cdot)} = 0$$

$$H_1: \beta_{(\cdot)} \neq 0$$

$$\{ \cdot \} \subseteq \{1, \dots, p\}$$

- Null and alternative hypothesis: $H_0 : \theta \in \Theta_0; H_1 : \theta \in \Theta_1,$

$$\Theta_0 \cup \Theta_1 = \Theta$$

- Rejection region: $R \subset \mathcal{X}$; if $x \in R$ "reject" H_0

$$\alpha = P_{\cap H_0} \{x \in R\}$$

- Test statistic and critical value: $R = \{x \in \mathcal{X} : t(x) > c\}$

c to be chosen

- Type I and Type II error: $\Pr\{t(X) > c | \theta \in \Theta_0\}, \quad \Pr\{t(X) \leq c | \theta \in \Theta_1\}$

Def
for

- Power and Size: $\beta(\theta) = \Pr_{\theta}(X \in R)$

$$\alpha = \sup_{\theta \in \Theta_0} \beta(\theta)$$

function of θ

size

- Optimal tests: among all level- α tests, find that with the highest power under H_1

of size α

level- α means size $\leq \alpha$

Neyman-Pearson '33

Example: Two-sample t -test

EH §1.2

1.2 Hypothesis Testing

Our second example concerns the march of methodology and inference for *hypothesis testing* rather than estimation: 72 leukemia patients, 47 with **ALL** (acute lymphoblastic leukemia) and 25 with **AML** (acute myeloid leukemia, a worse prognosis) have each had genetic activity measured for a panel of 7,128 genes. The histograms in Figure 1.4 compare the genetic activities in the two groups for gene 136.

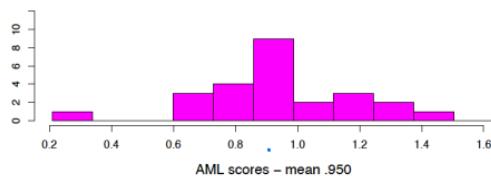
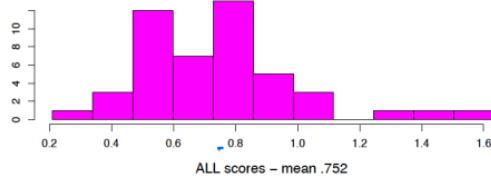


Figure 1.4 Scores for gene 136, leukemia data. Top **ALL** ($n = 47$), bottom **AML** ($n = 25$). A two-sample t -statistic = 3.01 with p -value = .0036.

The **AML** group appears to show greater activity, the mean values being $\text{ALL} = 0.752$ and $\text{AML} = 0.950$. (1.5)

$$x_{i1} \sim N(\mu_1, \sigma^2)$$

$$x_{i2} \sim N(\mu_2, \sigma^2)$$

← gp 1 values

x_{i1} "genetic activity" for patient i on gene \leftrightarrow row (136)

gp 2

$$H_0: \mu_1 = \mu_2$$

$$H_1: \mu_1 \neq \mu_2$$

$$t(\bar{x}) > c$$

then

$$\mu_1 \neq \mu_2$$

$$P_{H_0} \{ t(\bar{x}) > c \}$$

$$= \alpha$$

$$\frac{(\bar{x}_1 - \bar{x}_2)}{\sqrt{s_e(\bar{x}_1 - \bar{x}_2)}} = t$$

... Example 1

AoS Ex.10.8

$$s^2 = \frac{\sum (x_i - \bar{x})^2 + \sum (y_i - \bar{y})^2}{n_1 + n_2 - 2} s^2 \left(\frac{1}{n_1} + \frac{1}{n_2} \right)$$

$$\frac{s_1^2}{n_1} + \frac{s_2^2}{n_2}$$

```
leukemia_big <- read.csv("http://web.stanford.edu/~hastie/CASI_files/DATA/leukemia_big.csv")
```

```
oneline <- leukemia_big[136,]
```

```
one <- c(1:20, 35:61) # I had to extract these manually,
```

```
two <- c(21:34, 62:72) # couldn't figure out the data frame
```

```
n1 <- length(one); n2 <- length(two)
```

```
mean_one <- sum(oneline[1,one])/n1. ## [1] 0.7524794
```

```
mean_two <- sum(oneline[1,two])/n2. ## [1] 0.9499731
```

```
var_one <- sum((oneline[1,one]-mean_one)^2)/(n1-1)
```

```
var_two <- sum((oneline[1,two]-mean_two)^2)/(n2-1)
```

```
pooled <- ((n1-1)*var_one + (n2-1)*var_two)/(n1+n2-1)
```

```
taos <- (mean_one-mean_two)/sqrt((var_one/n1)+(var_two/n2))
```

```
## [1] -3.132304
```

```
tbe <- (mean_one-mean_two)/sqrt(pooled*((1/n1)+(1/n2)))
```

```
## [1] -3.035455
```

"reject H_0 "

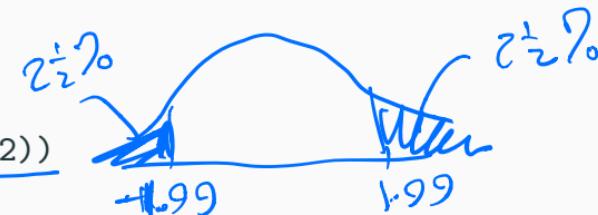
at level .05

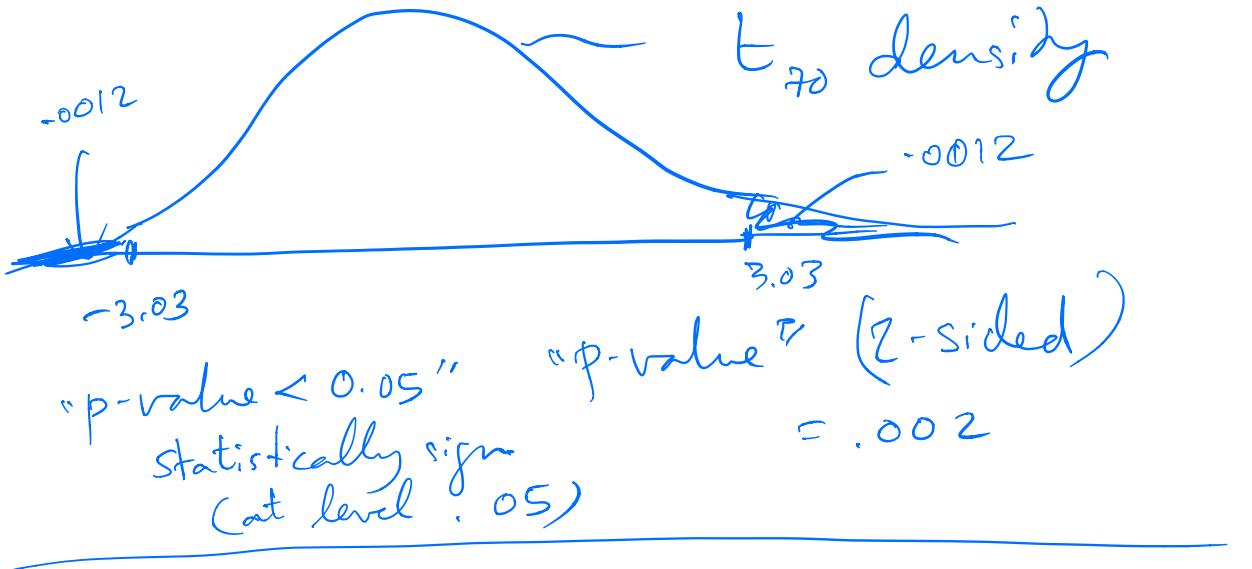
$$n_1, n_2 = 72 \quad d.f. \quad [n_1 + n_2 - 2 = 70]$$

P_{H_0} (reject H_0)

$$= p_{H_0}(|T_{70}| \geq c_\alpha)$$

$$c_\alpha = 1.99$$





$X_1, \dots, X_n \text{ iid } N(\mu, \sigma^2)$ σ^2 known

Let's use $R = \{\bar{X} > c\}$ sample mean

$$P_{H_0}(\bar{X} > c) = \alpha \quad \text{under } H_0: \mu \leq 0$$

H_0 unif. most powerful

$$A_c: \mu \geq 0$$

$$\beta(\mu) = P_{H_0}(\bar{X} > c)$$

$$\begin{aligned} \beta(\mu) &= P_{H_0}\left(\frac{\bar{X} - \mu}{\sigma/\sqrt{n}} > \frac{c - \mu}{\sigma/\sqrt{n}}\right) \\ &= P_{H_0}\left(Z > \frac{c - \mu}{\sigma/\sqrt{n}}\right) = \alpha \quad \frac{c - \mu}{\sigma/\sqrt{n}} = -1.96 \end{aligned}$$

nope? μ \bar{X} c

$$\tilde{\mu} = \text{median } (X_1, \dots, X_n) \quad P_{H_0}(\tilde{\mu} > c) = \alpha$$

R

Example: Likelihood inference

X_1, \dots, X_n i.i.d. $f(x; \theta)$; $\hat{\theta}(X_n)$ is maximum likelihood estimate. From last week:

AoS

$$(\hat{\theta} - \theta) / \widehat{se} \stackrel{\text{AoS}}{\approx} "Wald"$$

To test $H_0 : \theta = \theta_0$ vs. $H_1 : \theta \neq \theta_0$ we could use

$$T = W(X_n) = (\hat{\theta} - \theta_0) / \widehat{se},$$

approx.

The critical region will be $\{x : |W(x)| > z_{\alpha/2}\}$, i.e. "reject" H_0 when $|W| \geq z_{\alpha/2}$

This test has approximate size α :

$$\alpha = .05 \quad z_{\alpha/2} = 1.96$$

$$\Pr_{H_0}(|W| > z_{\alpha/2}) = \alpha.$$

Power? See Figure 10.1 and Theorem 10.6

logistic regression

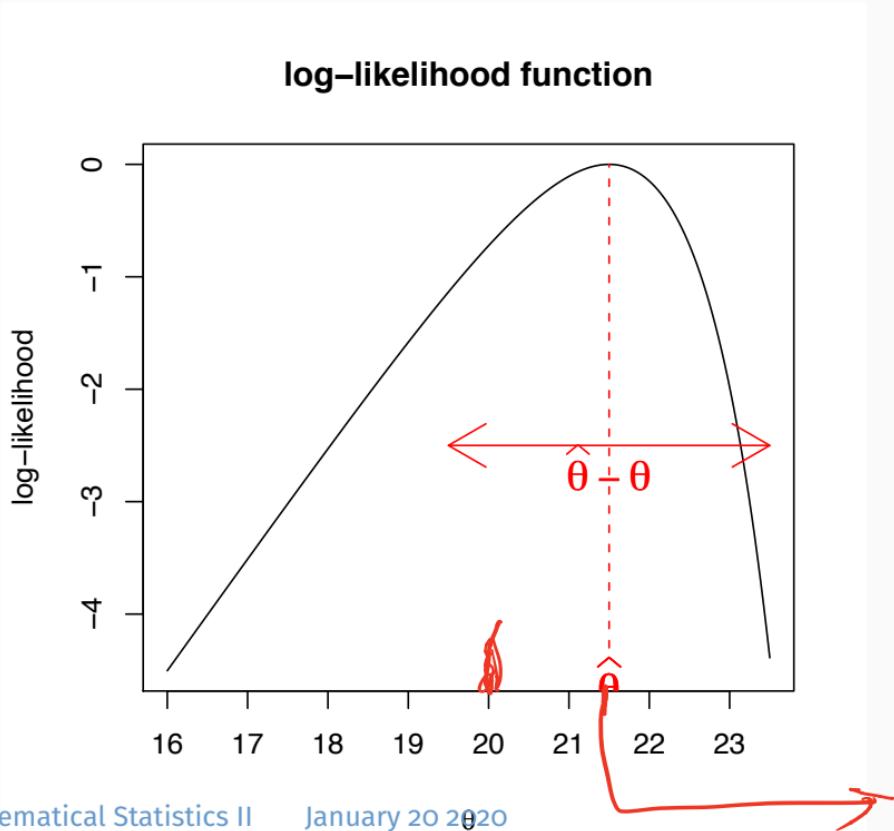
$$\widehat{se} = \sqrt{f^{-1}(\hat{\theta})}$$

$$\text{or } = \sqrt{I_n^{-1}(\hat{\theta})}$$

Wald means mle $\hat{\theta}$

|W| "large" when H_0 is false.

... likelihood inference



Example: comparing two binomials

AoS Ex.107

X ind't of Y

$$X \sim Bin(n_1, p_1), \quad Y \sim Bin(n_2, p_2), \quad \delta = p_1 - p_2, \quad H_0 : \delta = 0$$

mle $\hat{p}_1 = \frac{X}{n_1}$ $\hat{p}_2 = \frac{Y}{n_2}$ $\hat{\delta}_{mle} = \frac{X}{n_1} - \frac{Y}{n_2} = \hat{p}_1 - \hat{p}_2$

$$\begin{aligned}\hat{se}(\hat{\delta}) &= \sqrt{\text{var}(\hat{p}_1 + \hat{p}_2)} \leftarrow \text{known} \\ &= \sqrt{\frac{\hat{p}_1(1-\hat{p}_1)}{n_1} + \frac{\hat{p}_2(1-\hat{p}_2)}{n_2}}\end{aligned}$$

~~(approx)~~ Normal approx \Rightarrow

"Wald" test at level α reject $H_0: \delta = 0$
 $| \hat{\delta} / \hat{se}(\hat{\delta}) | > 1.96$

power $P_{\alpha} \left\{ \frac{\hat{\delta}}{\text{se}(\hat{\delta})} > 1.96 \right\}$ equality of means; equality of medians; Wald test

$\downarrow \hat{\delta} = \hat{\delta} - \delta$

 $= P_{\alpha} \left\{ \hat{\delta} > 1.96 \cdot \text{se}(\hat{\delta}) \right\} \xrightarrow{\text{N approx.} \sim \text{binom.}}$

$\downarrow \beta(\delta) = P_{\alpha} \left\{ \frac{\hat{\delta} - \delta}{\text{se}(\hat{\delta})} > 1.96 \cdot \text{se}(\hat{\delta}) - \frac{\delta}{\text{se}(\hat{\delta})} \right\}$

~~Test~~

10.9 X_1, \dots, X_n iid $f_c(\cdot)$ Y_1, \dots, Y_m iid $f_r(\cdot)$

$H_0: \text{med}_1 = \text{med}_2 \quad H_A: \text{med}_1 \neq \text{med}_2$
 $m_1 - m_2 = \Delta$

find a f: $T = t(x^n)$

know at least $T \underset{H_0}{\sim} (\text{some density})$

"Reject" H_0 iff $\underset{\text{size}}{\sup}_{H_0} P_{n, H_0}(T \in R) \leq \alpha$

need T to be "sensitive" to H_1 $\underset{H_1}{\inf} P_{n, H_1}(T \notin R) \geq 0.95$

The formal theory of testing imagines a decision to "reject H_0 " or not, according as $X \in R$ or $X \notin R$, for some defined region R (e.g. $|Z| > 1.96$) "at level .05" $p_{A_0}(T \in R)$
 This is useful for deriving the form of optimal tests, but not useful in practice.

Doesn't distinguish between $Z = 1.97$ and $Z = 19.7$, for example.

P-values give more precise information about the null hypothesis

AoS definition: p-value = $\inf\{\alpha : T(X_n) \in R_\alpha\}$ $p \approx .048$

Def 10.11

SM definition $p_{obs} = \Pr_{H_0}\{T(X_n) \geq t_{obs}\}$

← T observed value if H_0 true

" \Pr {getting a result as or more extreme than the data, if H_0 is true}"

$$\Pr(\bar{X} > 7; \mu=0) = \Pr(N(0,1) > \frac{7}{\sigma}) \leq 10^{-7}$$

(but $\sigma=1$)

$$H_0: \mu=0$$

$$H_1: \mu \neq 0$$

$$t^{obs} = t(\bar{X}_n) = \frac{1}{n} \sum (X_i - \bar{X})$$

Example: exponential

SM Ex.7.22

X_1, \dots, X_n i.i.d.

$$f(x; \lambda) = \lambda e^{-\lambda x}$$

$$L(\lambda; \underline{x}) = \lambda^n e^{-\lambda \sum x_i}$$

$$H_0 : \lambda = \lambda_0$$

$$\checkmark t(\underline{x}_n) = \sum x_i \sim P(n, \lambda)$$

$$P_{\mathcal{H}_0}(t(\underline{x}_n) > t_{\text{obs}}) = \dots$$

... Example: logistic regression

```
Boston.glmnull <- glm(crim2 ~ 1, family = binomial, data = Boston)
anova(Boston.glmnull, Boston.glm)
```

Analysis of Deviance Table

Model 1: crim2 ~ 1

Model 2: crim2 ~ (crim + zn + indus + chas + nox + rm + age + dis + rad +
tax + ptratio + black + lstat + medv) - crim

	Resid. Df	Resid. Dev	Df	Deviance
1	505	701.46		
2	492	211.93	13	489.54

```
> pchisq(489.54, 13, lower.tail = F)
[1] 2.435111e-96
```