

p-value

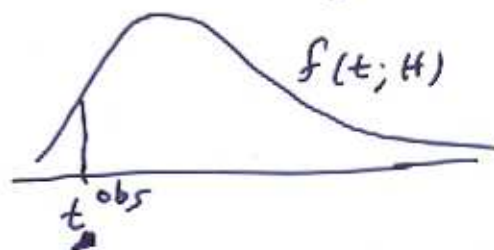
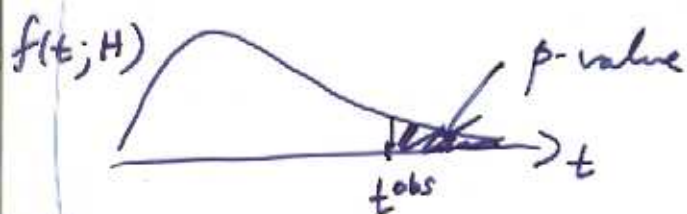
If we have a function of data $t(Y) = T$, say, and a model for the data, hypothesis, then p-value for hypothesis, based on t ,

$$\text{probability } (T \geq t^{\text{obs}}; \text{ model})$$

$$t^{\text{obs}} = t(y)$$

(1-sided p-value)

probability (results as or more extreme than observed) \rightarrow both or one?



$$1 - F_T(t^{\text{obs}})$$

$$\text{dist} = f = \int T$$

we would probably use

$$p = \Pr(T \leq t^{\text{obs}}) \text{ bec. evidence in other direction}$$

- we need a 1-dim $t(\cdot)$ for this to work

$$\checkmark \text{ t-stat } \frac{\bar{y} - \bar{x}}{s\sqrt{\dots}}$$

- we need to know its dist- under model $\checkmark t_{n+m-2}$

Density Estimation § 5.6

Sample y_1, \dots, y_n assume i.i.d. $f(y)$

$f(y)$ = density f for y i.e.

$f(y)dy = \Pr(Y \in (y, y+dy))$ where Y is a r.v.
 ~~with same~~

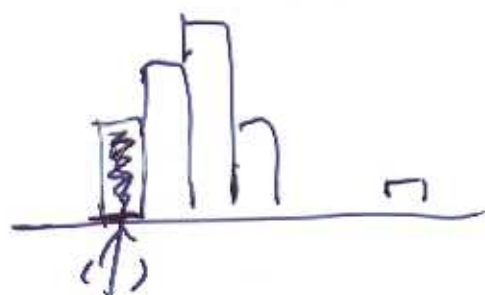
$$f(y) = \frac{d}{dy} \Pr(Y \leq y)$$

How to estimate $f(y)$?

- histogram is simplest

> $\text{hist}(y)$ counts

> $\text{hist}(y, \text{prob} = T)$ freq.



- not smooth, even if $f(\cdot)$ is

$$\frac{\# \text{ } y_i \text{'s in } (,)}{n}$$

- depends a lot on how you choose bins; # + start pt.

> `truehist` in MASS library is better; easier to specify bins

Fig 5.8

- kernel density estimator (better)

$$\hat{f}(y) = \frac{1}{nb} \sum_{i=1}^n K\left(\frac{y-y_i}{b}\right)$$

$K(\cdot)$ is some function e.g. $K(x) = \frac{1}{\sqrt{2\pi}} e^{-\frac{1}{2}x^2}$

b is some parameter "bandwidth"

- default for $K(\cdot)$ is Gaussian (normal), but there are others (but not usually so crucial)

- choice of b - called bandwidth (bw in R) is very crucial

$$E \int \{ \hat{f}(y) - f(y) \}^2 dy$$

$$\hat{f}(\cdot) = \hat{f}(\cdot; y_1, \dots, y_n)$$

fn. of data

Integrated mean squared error

$$= \frac{1}{nb} \int K^2(y) dy + \frac{1}{4} b^4 \int f''(y)^2 dy \left[\int y^2 K(y) dy \right]^2 + O\left(\frac{1}{nb} + b^4\right)$$

$$n \rightarrow \infty, b \rightarrow 0, nb \rightarrow \infty$$

"If we neglect $O(\cdot)$, then optimal b is

$$b = \left[\frac{\int K^2(y) dy}{n \int f''(y)^2 dy \left(\int y^2 K(y) dy \right)^2} \right]^{1/5}$$

- now 2 possibilities
- ① use an est. of $f''(\cdot)$
to find the b e.g. $\hat{f}''(\cdot)$
 - ② solve this eqⁿ with b on both sides

- these are called

bw = "SJ. dpi" \leftarrow "direct plug-in" ①

or bw = "SJ. ste" \leftarrow "solve the eq'n" ②

default is neither, but

bw = ~~$\frac{1}{n} \int f''(y)^2 dy$~~

which ~~means~~ is ~~really~~ a slight modⁿ of

$$b = 0.9 \min\left(\hat{\sigma}, \frac{IQR}{1.34}\right) n^{-1/5}$$

\nearrow "ad-hoc" suggestion that's
after ok

Default choice may need improving;

next best is $bw = SJ \cdot dpi$
or ste //

Chapter 7 Generalized Linear Models

Linear model $y_i = x_i^T \beta + \sigma e_i$ $e_i \sim f_0(\cdot)$

$$f(y_1, \dots, y_n) = \prod_{i=1}^n f_0\left(\frac{y_i - x_i^T \beta}{\sigma}\right) \cdot \frac{1}{\sigma^N}$$

~~data~~ regression-scale

Gen. lin model $f(y_1, \dots, y_n) = \prod_{i=1}^n f_1(y_i; x_i^T \beta, \sigma)$

$f_1(\cdot)$ is a density f_0 , but not location-scale
as above

in exponential family \leftarrow to be defined

- extend linear inference, or linear regression
to more choice of models