

W04 - p-values, significance. The Student's t-test

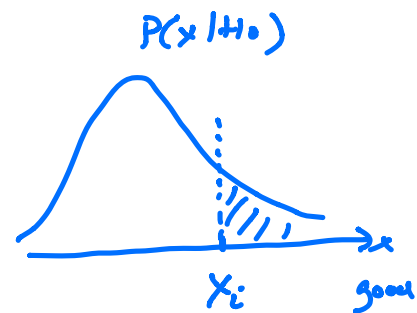
Experiments - $X \rightarrow \text{Data } \{x_i\} = D$

H_1 - interesting hypothesis about D , that you want to test/validate

H_0 - null hypothesis that "could" also explain your results

H_1 is hard to formulate

H_0 is easier to formulate



take one value x_i :

$$\text{p-value}(x_i) = P(x \geq x_i | H_0) \quad \text{(right tail)}$$

$$= 1 - P(x < x_i | H_0)$$

$$= 1 - \text{CDF}_{H_0}(x_i)$$

$$\text{p-value}(x_i) = P(x \leq x_i | H_0) \quad \text{left tail}$$

$$\text{p-value}(x_i) = P(x \geq x_i \wedge -x_i \leq x | H_0) \quad \text{two tailed test}$$

If p-val is small, there is little chance that H_0 describes the data.

p-value does not use any info about the other hypothesis H_1
rejects H_0 does not mean H_1 is true

Example experiment failure rate

$H_0 : f = f_0 = 0.2$ what prob do you

$H_1 : f > f_0$

$D = \{ ssfff \}$

$$\begin{aligned} p\text{-value}(ssfff)_{H_0} &= P(ssfff | f = f_0) \\ &+ P(stfff | f = f_0) \\ &+ P(ftfff | f = f_0) \end{aligned}$$

$$= \frac{5!}{2!3!} f_0^2 (1-f_0)^3 + \frac{5!}{1!4!} f_0^4 (1-f_0) + \frac{5!}{0!5!} f_0^5$$

$$= 0.05792$$

$(f_0 = 0.2)$

What do you do with this result

$$p\text{-val}(ssfff)_{H_0} = 0.05792$$

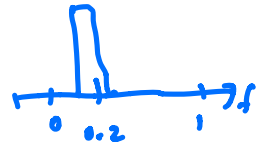
do you reject $f = f_0$ or not?

The Bayesian Approach

$$\frac{P(H_0 | s s f f f)}{P(H_1 | s s f f f)} = \frac{P(s s f f f | H_0) P(H_0)}{P(s s f f f | H_1) P(H_1)} \quad \left. \begin{array}{l} \\ \\ \\ \end{array} \right\} P(H_0) = P(H_1)$$

$$= \frac{\int_0^1 df P(s s f f f | f) P(f | H_0) df}{\int_0^1 df P(s s f f f | f) P(f | H_1) df}$$

$$P(f | H_0) = \begin{cases} \frac{1}{0.21 - 0.19} & 0.19 < f < 0.21 \approx |f - f_0| < 0.1 \\ 0 & \text{otherwise} \end{cases}$$



$$P(f | H_1) = \begin{cases} \frac{1}{1 - 0.2} & f > f_0 = 0.2 \\ 0 & \text{otherwise} \end{cases}$$



$$\frac{P(H_0 | s s f f f)}{P(H_1 | s s f f f)} = \frac{\frac{1}{0.02} \int_{0.19}^{0.21} P(s s f f f | f) df}{\frac{1}{0.90} \int_{0.2}^1 P(s s f f f | f) df}$$

$$= \frac{\frac{1}{0.02} \int_{0.19}^{0.21} f^3 (1-f)^2 df}{\frac{1}{0.90} \int_{0.2}^1 f^3 (1-f)^2 df} \approx \frac{20}{90}$$

} 80:20 chance $f > f_0$

Which method do you prefer?

* $p\text{-value} = 0.0579$

. nothing about H_1

- easier but hard to interpret

$$p\text{-val} \neq P(H_0 \text{ is true})$$

$$1 - p\text{-val} \neq P(H_1 \text{ is true})$$

* Bayes

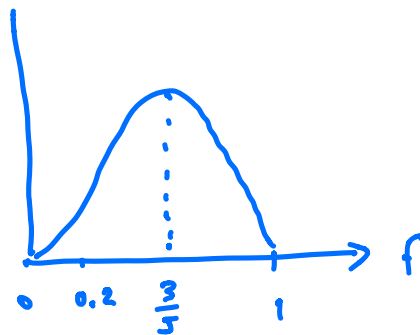
+ model comparison

$$H_1 : H_0 \\ 80 \quad 20$$

+ posteriors

$$P(f | \text{ssfff})$$

$$f^x = \frac{x}{N}$$

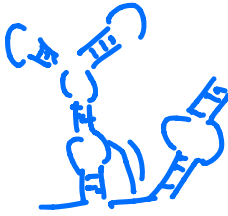
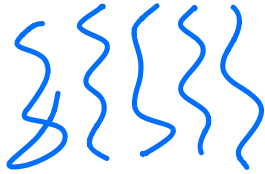


p-value is indirect

If I were to repeat a $N=5$ experiment many times and $f=f_0$, there will be a $\sim 62\%$ chance of obtaining at least 3 failures (ssfff or sffff or fffff)

Student's t-test a widely used p-value

structural RNAs



Given the RNAP sequence, what is the structure?

RNAP RNA
(ribozyme)

base pairs
A:U U:A
C:G G:C
G:U U:G

265 residues

160 pairs (80bp)

105 unpaired

Experimental chemical modification reactivities

$r = \text{reactivity}$ ($0 < r \leq 1$)

r measures "flexibility": is this a synonym of being paired/unpaired?

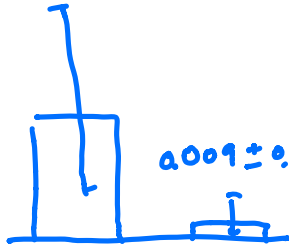
DMS, SHAPE (selective 2'-hydroxyl acylation analyzed by primer extension)

$$D = \{r_1 \dots r_{265}\}$$

$$D_u = \{105 \text{ } r\text{'s}\}$$

$$D_p = \{160 \text{ } r\text{'s}\}$$

0.045 ± 0.075



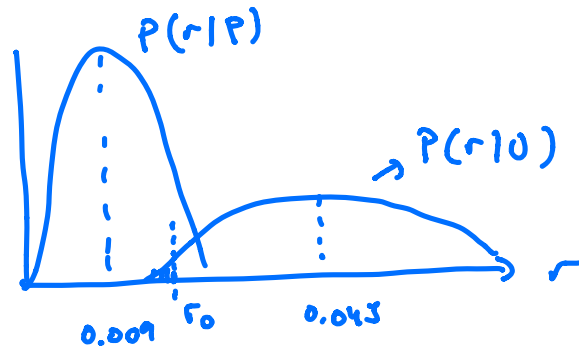
Student-t test
(D_u, D_p)

$p\text{-val} = 4.3 \cdot 10^{-9}$

Great! SHAPE should be enough to distinguish P/U

① Look at the data and assumptions

Ideal scenario



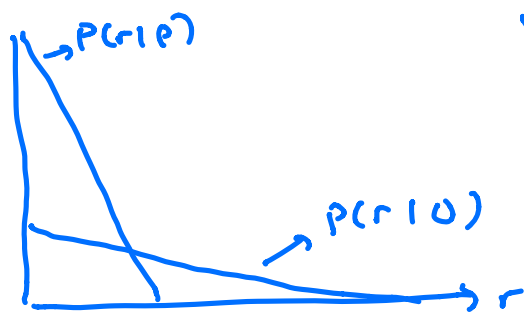
$r < r_0$ call it P

$r > r_0$ call it U

$p\text{-val}(r_0)_U = P(r < r_0 | U) \approx \text{small}$ (False Positives)

$P(r < r_0 | P) \approx 0.9$ (True Positives)

Real data



very non Gaussian

T-test assumes data is Gaussian!
and it uses the means only

What does the Student's distribution has to do

with anything? since 2.3, 3.2

assume $\{x_i\}$ follows a Gaussian dist $\mathcal{N}(\mu, \sigma)$

$$P(\mu, \sigma | \{x_i\}) \propto P(\{x_i\} | \mu, \sigma) P(\mu, \sigma)$$

$$= \left(\frac{1}{\sqrt{2\pi\sigma^2}} \right)^N e^{-\frac{\sum_i (x_i - \mu)^2}{2\sigma^2}} P(\mu) P(\sigma)$$

$$P(\mu) = \begin{cases} \frac{1}{\mu_+ - \mu_-} & \mu_- < \mu < \mu_+ \\ 0 & \text{otherwise} \end{cases}$$

$$P(\sigma) = \begin{cases} c & \sigma > 0 \\ 0 & \sigma < 0 \end{cases}$$

$$P(\mu, \sigma | \{x_i\}) \propto \sigma^{-N} e^{-\frac{\sum_i (x_i - \mu)^2}{2\sigma^2}}$$

$$P(\mu | \{x_i\}) \propto \int_0^\infty d\sigma \sigma^{-N} e^{-\frac{\sum_i (x_i - \mu)^2}{2\sigma^2}}$$

$$t = \frac{\sqrt{\sum_i (x_i - \mu)^2}}{\sigma} \quad d\sigma = -\frac{\sqrt{\sum_i (x_i - \mu)^2}}{t^2} dt$$

$$P(\mu | \{x_i\}) \propto \left[\sum_i (x_i - \mu)^2 \right]^{\frac{-N+1}{2}} \int_0^\infty \frac{1}{t} e^{-t^2/2} dt$$

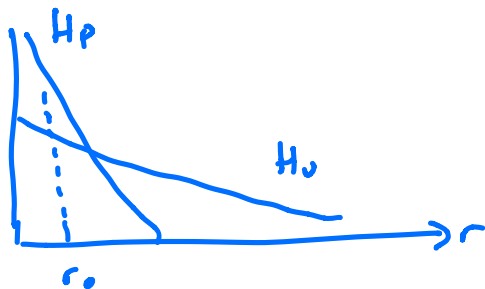
Student's distribution.

IS this p-value that compare the means

really what we want to know?

$\mu_1 = 0.045$ test of $(\mu_1 - \mu_2)^2$
 $\mu_2 = 0.009$ \rightarrow follows Student's dist

What I want to know is if I make the assignment



$r < r_0$ is P

$r > r_0$ is U

What errors do I make?

$p\text{-val}(r_0) = P(r \leq r_0 | H_0) = CDF_{H_0}(r_0)$ prob that a upper
reside has $r \leq r_0$

It test

$n \sim n \cdot p\text{-val}(r_0) = \text{Expected \# FP if all res follow null}$

Then pick p-value based on how many FP you are willing to tolerate

r_{react}	p-value	275 FP
0.0029	0.02	50% FP
0.0034	0.05	100 FP
0.0042	0.10	

N = # of test

p^* = p-value at r^*

F^* = # of N w/ $r \leq r^*$

T^* = $T \cap F^*$

T = 165 (true)

$$FDR = \frac{p^* \cdot N}{F^*}$$

fraction of positive calls
that should be expected
to be U

p^* fraction of tested (N) expected to be false calls
fdr fraction of posits (F^*) expected to be false calls

r	p^*	fdr	sen
0.0029	0.2	17.4	16.9
0.0034	0.05	30.0	23.1
0.0042	0.10	44.8	30.6
0.0063	0.28	66.5	50.0

to get 50% of real rated bases, more than half
of the called paired are going to be wrong