# Lecture 19

1. Exam Review

2. Linear Algebra!

1.



$H_0$ : negative

$H_1$ : "positive"

$\mu_0$ false negative

$\gamma$

$\mu_1$

$P_{fa}$ false positive (Type I)



$H_1$

$H_0$

p (false positive)

ROC
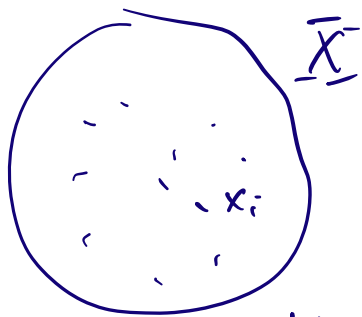
1.0

0

1

$$0.5 \leq AUC\ ROC \leq 1$$

# Hypothesis Testing

## T-test / Z-test

firearm mortality
from rural

$\bar{X}$

$x_i$

vs.

firearm mortality
from urban

$\bar{Y}$

$y_i$

sample mean
for population $\bar{X}$

$$\hat{N}_x = \frac{1}{N_x} \sum_{i=1}^{N_x} x_i$$

sample mean
from population $\bar{Y}$

$$\hat{N}_y = \frac{1}{N_y} \sum_{i=1}^{N_y} y_i$$

Statistic

$$t = \hat{\mu}_x - \hat{\mu}_y$$

$\hookrightarrow$ instantiation of R.V. T

If $t = \hat{\mu}_x - \hat{\mu}_y \neq 0$?

$\hookrightarrow \hat{\mu}_x \neq \hat{\mu}_y$?

## Hypothesis Test:

1) $H_0 : \mu_x = \mu_y$

   $H_1 : \mu_x \neq \mu_y$

   Two sets of data samples $(X, Y)$

2) $H_0 : \mu_x = \mu_{provided}$

   $H_1 : \mu_x \neq \mu_{provided}$

   One set of data samples

$\sigma_x^2 \equiv$ true variance for population $X$

$\sigma_y^2 \equiv$ " " " " " $Y$

1) $T = \mu_x - \mu_y = 0$ (under $H_0$)

Known $\sigma_x^2$ & $\sigma_y^2$

Unknown $\sigma_x^2$ & $\sigma_y^2$

1) $T \sim G\left(0, \dfrac{\sigma_x^2}{N_x} + \dfrac{\sigma_y^2}{N_y}\right)$

$\Rightarrow$ Estimate the variance from data, e.g.

If $\sigma_x^2 = \sigma_y^2 = \sigma^2$,

$T \sim G\left(0, \sigma^2\left(\dfrac{1}{N_x} + \dfrac{1}{N_y}\right)\right)$

1) $S_x^2 = \dfrac{1}{N_x - 1} \sum\limits_{i=1}^{N_x} (x_i - \bar{x})^2$

$S_y^2 = \dfrac{1}{N_y - 1} \sum\limits_{i=1}^{N_y} (y_i - \bar{y})^2$

2) $T \sim G\left(0, \dfrac{\sigma^2}{N}\right)$

$T = \mu_x - \mu$ provided $= 0$
(under $H_0$)

$T \sim$ student's T dist.
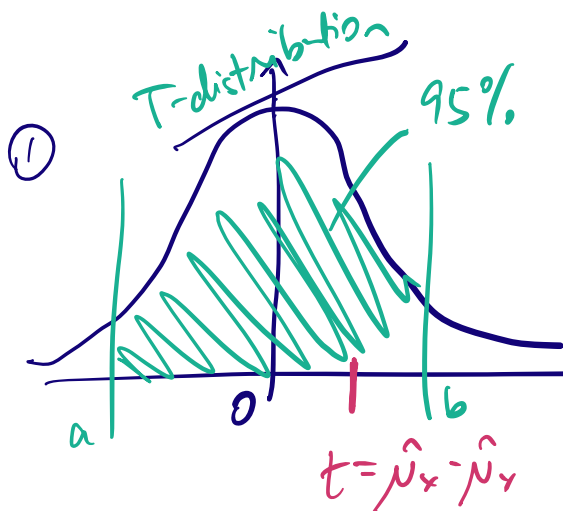$\left(\text{dof} = N_x + N_y - 2\right)$
scale $= \sqrt{\sigma_x^2/N_x + \sigma_y^2/N_y}$
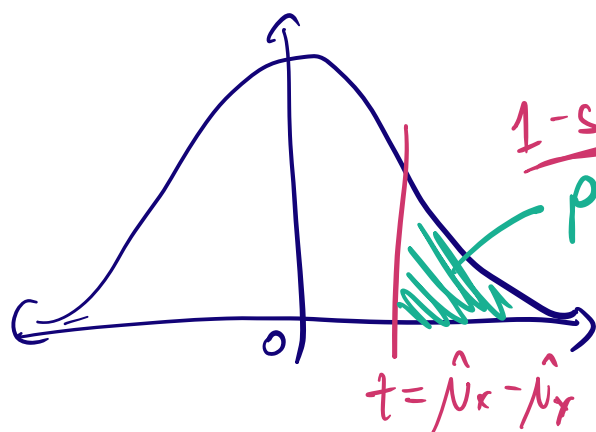
2) $T \sim$ T dist. (dof $= N - 1$)

$t = \hat{\mu}_x - \hat{\mu}_y$

$t$ is statistically significant :

①

T-distribution

95%

$a$     $0$    $b$

$t = \hat{\mu}_x - \hat{\mu}_y$

If $t \in [a, b]$, we cannot reject $H_0$

But if $t \notin [a, b]$, we can reject $H_0$.

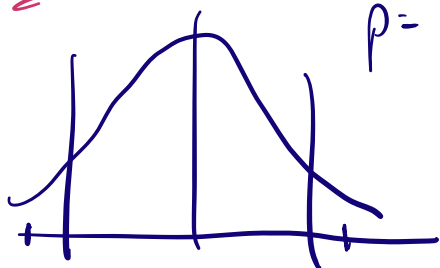1-sided:

$0$

$t = \hat{\mu}_x - \hat{\mu}_y$

$p = P(T \geq t)$ $\equiv$ Survival function @ $t$

If $p < \alpha$ $(\alpha = 0.05)$, then we reject $H_0$.

Otherwise, we cannot reject $H_0$.

2-sided:

$p = P(|T| \geq t)$, assuming $t > 0$

$t$

$p = P(|T| \geq t) = 2P(T \geq t)$

If $p < \alpha$, then we reject $H_0$.

If you're not sure, use a 2-sided hypothesis test.

1-sided:

$$H_1 : \hat{\mu}_x > \hat{\mu}_y$$

$$\hat{\mu}_x > \mu_{provided}$$

CDF:

$$P(X \le x) = \begin{cases} \sum_{k \le x} p_x(k), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{x} f_x(x)\,dx, & \text{if } X \text{ is continuous} \end{cases}$$

# Goodness-of-fit measures

| Discrete R.V. | Continuous |
|---|---|

**Discrete R.V.**
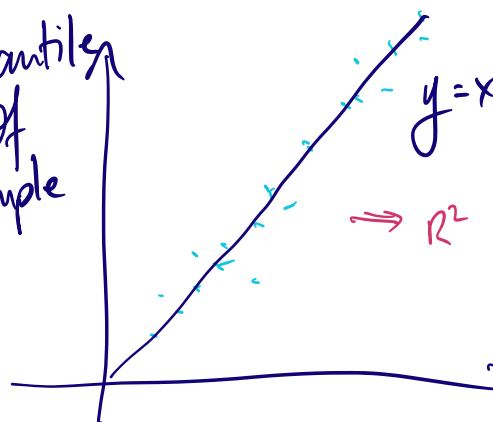
$\chi^2$-test

$T \sim \chi^2(K)$

$\downarrow$ $dof = N_x - 1$

$chi = \sum\limits_{i=1}^{N_x} \dfrac{(x_i - E[x])^2}{E[x]}$

$E[x] = \hat{N}_x$
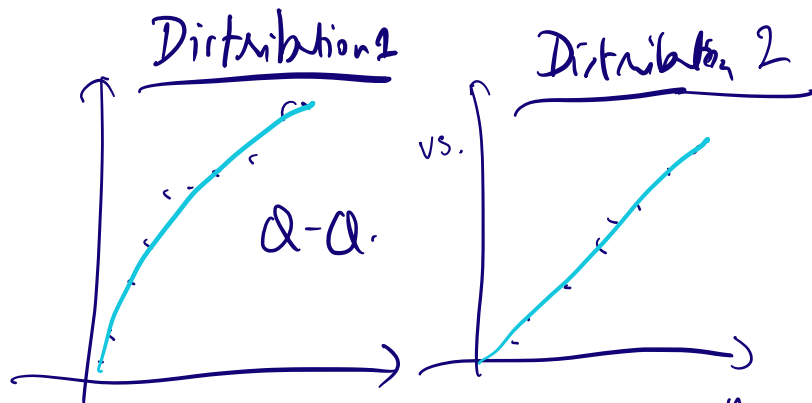
$P(T \geq chi) = p$

**Continuous**

quantile of sample



$y = x$

$\Rightarrow R^2$

quantile of R.V. (true dist.)

$F_x(x)$

1

percentile 0.5

quantile

$x$

## Q-Q plot

$R^2 =$ coefficient of determination

$R^2 > 0.9$, it is a good fit



Dirtribution 1        vs.        Distribution 2

Q-Q.

If $R_1^2 < R_2^2 \Rightarrow$ pick dist$^n$ 2.