

### Week 9.3: Evaluation in ML

**Objective:** 1) Is model A better than B ? 2) Is the difference between the results statistically significant?

**A/B Testing:** Simple Control Experiments: 1) Randomly split traffic between two or more version eg.(A) Control, (B) Treatment (statistical test) เพื่อ confirm ความต่างตามเมื่อ มี feature ใหม่  $var(X) = \frac{\sum_{i=1}^n (x_i - \bar{X})^2}{n-1}$  Hypothesis Testing:  $H_0$ ,  $H_A$  เพื่อหาความแตกต่าง

Type I - reject  $H_0$  when in fact it True.(FP)  
"False Alarm"

Type II - FTR  $H_0$  when in fact False.(FN)

P-value - ถ้า p-value  $\leq \alpha$ : then reject  $H_0$  at level alpha

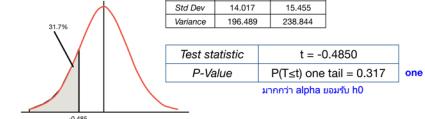
ทำ p-value testing 1. calculate test statistic  
2. Convert the result to a p-value by comparing its value to the distribution of test statistics under the null hypothesis.

3. Decide ว่า reject หรือไม่ reject  $H_0$ .

$$t = \frac{\bar{X}_A - \bar{X}_B}{\sqrt{\frac{var(A)}{n_A} + \frac{var(B)}{n_B}}} = \frac{\bar{X}_D \times \sqrt{n}}{\sigma_D}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800
Std Dev	14.017	15.455



Test statistic	t = -0.4850
P-Value	P(T<t) = one tail = 0.317

หมายเหตุ alpha ของห 0

### Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

### McNemar's Test

Measure for comparing paired proportions.

e.g. Which is better, classifier C2 or C3 ?

Applied to 2x2 contingency table.

Test captures two key differences:

$n_{01}$ : number misclassified by 1<sup>st</sup> but not 2<sup>nd</sup> classifier.

$n_{10}$ : number misclassified by 2<sup>nd</sup> but not 1<sup>st</sup> classifier.

Contingency for C2 v C1

3	2	$n_{00}$	$n_{01}$
0	5	$n_{10}$	$n_{11}$

$\chi^2 = \frac{(|n_{01} - n_{10}| - 1)^2}{n_{01} + n_{10}}$

Note: For test to be applicable require  $(n_{01} + n_{10}) > 10$

Contingency for C3 v C1

1	2	$n_{00}$	$n_{01}$
4	3	$n_{10}$	$n_{11}$

$\chi^2 = \frac{1}{6} = 0.1666$

$\chi^2 > 3.84$  required for statistical significance at 95%. So neither classifier significantly better!

### Week 10.1: Decision Analytic Thinking

Acc = 1- error rate หรือ Correct/Total Predict

Misclassification rate = 1 - Acc

**True Positive Rate:**

Focus on TPs  $TPRate = \frac{TP}{TP + FN}$

Also called **Sensitivity** ความแม่นยำในการพยากรณ์ว่าคนมี疾患จริง

or called **Recall**, aka. Probability of Detection

**False Positive Rate:**

Focus on FPs  $FPRate = \frac{FP}{FP + TN}$

**True Negative Rate:**

Focus on TNs  $TNRate = \frac{TN}{FP + TN}$

Also called **Specificity**

精度 =  $\frac{TP}{TP + FN}$  ความแม่นยำ 2 สมมติฐาน  $H_0$ ,  $H_A$  เพื่อหาความแตกต่าง

Type I - reject  $H_0$  when in fact it True.(FP)  
"False Alarm"

Type II - FTR  $H_0$  when in fact False.(FN)

P-value - ถ้า p-value  $\leq \alpha$ : then reject  $H_0$  at level alpha

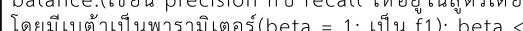
ทำ p-value testing 1. calculate test statistic  
2. Convert the result to a p-value by comparing its value to the distribution of test statistics under the null hypothesis.

3. Decide ว่า reject หรือไม่ reject  $H_0$ .

$$t = \frac{\bar{X}_A - \bar{X}_B}{\sqrt{\frac{var(A)}{n_A} + \frac{var(B)}{n_B}}} = \frac{\bar{X}_D \times \sqrt{n}}{\sigma_D}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800
Std Dev	14.017	15.455



Test statistic	t = -0.4850
P-Value	P(T<t) = one tail = 0.317

หมายเหตุ alpha ของห 0

paired(มี)  $t = \frac{\bar{X}_D \times \sqrt{n}}{\sigma_D}$  Look at the mean and standard deviation of the differences (deltas)

D(delta))  $t = \frac{-3.2 \times \sqrt{10}}{5.2} = -1.946$

Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

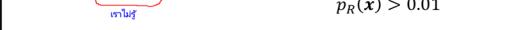
Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800



Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

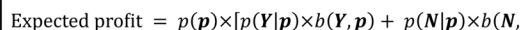
Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800



Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800



Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800



Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$

$$t \text{ statistic} = \frac{\text{Difference in proportions}}{\text{Standard error}} = \frac{c_1 + c_2}{\sqrt{(1-p) \times (\frac{1}{n_1} + \frac{1}{n_2})}}$$

ดูตารางตามค่า significant ที่ได้ ถ้าเป็นเล็กๆ ก็ เอาไป -1 ด้วย

	Team A	Team B
N	10	10
Mean	20.600	23.800



Test statistic	t = -1.946
P-Value	P(T<t) = 0.084

two tails

Difference in proportions

A t-test is sometimes used to analyse differences in proportions e.g. comparison of conversion rates in A/B testing.

Requires a number of assumptions about the population which are usually not true.

	Control	Treatment
Samples	$n_1$	$n_2$
Conversions	$c_1$	$c_2$