# Statistics Part II
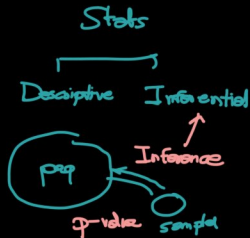
- Hypothesis Test (Inference)
- Binary Classification 0/1
- Clustering → Kmeans , set.seed()
  - kmeans( )
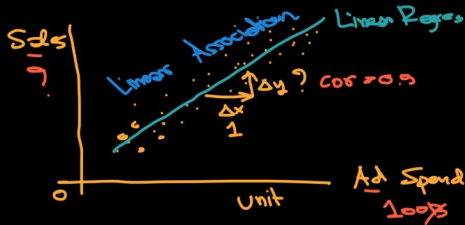
ML
- Supervised (Prediction)
- Unsupervised (Summarise)

Stats
- Descriptive
- Inferential ← Inference

p≥q

p-value    sample

Hypothesis Testing

1. Comparison : AB Test
2. Association ( Pearson Correlation )
3. Prediction : Regression

Correlation does not
Imply causation

& now! 

+
Strength



Sales = f (Ad)

Sales = 500 + 50 Ad

Slope = $\frac{\Delta y}{\Delta x}$

Linear Association     Linear Regres

$\Delta y$ ? cor = 0.9
$\Delta x$
1

Sales

Ad Spend
100$

Unit

0

Ronald Fisher (1925) $\longrightarrow$ Frequentist Approach

SPSS

Identical
Experiment  20 m𝔰𝔰 $\longrightarrow \dfrac{1}{20} = 5\%$ $\longrightarrow$ Reject Hyp.

Arbitrary 1% 5% 10%

Low

Chance

$\downarrow$ p ( H 95% | Fair coin )



$\alpha$  $P(HHH \ldots 95)$

m𝔰𝔰

100 m𝔰𝔰

Fair Coin  50% (0.5)

$\dfrac{1}{2} = 0.5$

p-value ≤ 5%

$$P(H \geq 95 \mid \text{Fair Coin})$$

$0.001 \leq 0.05$

~~Ho: Fair Coin~~

✓ Ha: Not a Fair Coin

Fisher

Reject Ho : p-value ≤ 0.05

Fail to reject Ho : p-value > 0.05

Sales

$y = b_0 + b_1 x$

Ad Spend

Significance Test

Sales $= f(Ad)$

Sales $= b_0 + b_1 Ad$

Sales $= 500 + 50 \cdot Ads$

Two tailed.

$H_0: b_1^{Ad} = 0$

$H_a: b_1 \neq 0$

Null Hyp (H_0)

Fail to reject H_0

Reject H_0

$H_0$ 0.5

50.

Sales

$y = b_0 + b_1 x$

Ad Spend

Significance Test

$Sales = f(Ad)$

$Sales = b_0 + b_1 Ad$

$Sales = 500 + 50 \cdot Ads$

Two tailed.

$H_0: b_1^{Ad} = 0$

$H_a: b_1 \neq 0$

Null Hyp $(H_0)$

Fail to reject $H_0$

Reject $H_0$

$H_0$ 0.5

50.

NHST

Null Hypo $H_0$

0.05

Critical Rejet
0.025

Fail to reject $H_0$

0.025

$\leftarrow$ 0 $\rightarrow$

Reject $H_0$   $\uparrow$ p.value
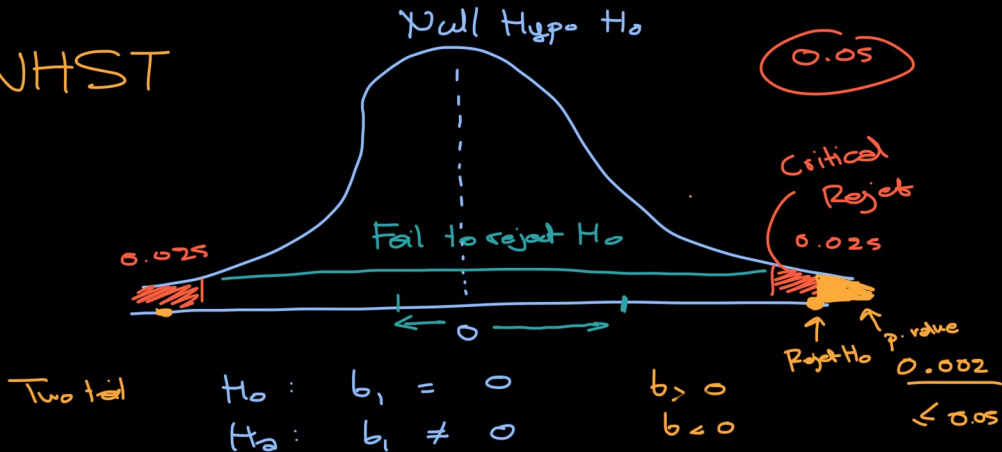0.002
$\leq$ 0.05

Two tail

$H_0$ : $b_1 = 0$        $b > 0$
$H_a$ : $b_1 \neq 0$        $b < 0$

1. p.value $\leq$ alpha (.05) ✓

2. confidence Interval ✓

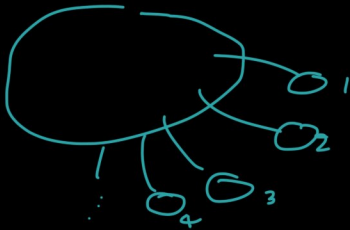$$\text{Sales} = b_0 + b_1 \cdot \text{Ad}$$

$$= 500 + 50 \cdot \text{Ads}$$

estimate

[ 0 ]

CI láuín $H_0$ (0)

↳ Fail to reject Hb

Ad

$H_a: b_1 \neq 0$

CI [ láu o ] Reject Hb

1. Ho, Ha.
2. Collect Data
3. Conclusion

model  Ad
P. $\leq$ 5%
Sig!

Better  Ad = $\boxed{95\%}$ [ +50, +65 ]
Explanation $_{100}\uparrow$

AB Test $\Big<$ $\begin{array}{l} Ad_1 \uparrow \\ 50-65 \quad 95\% \\ Ad_2 \end{array}$

Ho: $Ad_1 = Ad_2$ $\longleftarrow$ T.test /Linear
Ha: $Ad_1 \neq Ad_2$

# p · value definition

$$P\left(\underset{\text{or more extreme}}{\text{observed data}} \geq 95 \;\middle|\; H_0 \text{ is TRUE}\right)$$

$H_0$: Fair coin

$H_a$: Not a fair coin

ความน่าจะเป็น
ที่จะเห็น H ≥ 95 ?

ถ้า เหรียญ เป็น Fair Coin
จริง.

Binary Classification → 0/1 yes, no

Loan

1

$z = b_0 + b_1 x$

Sigmoid function

$\dfrac{5}{6}$

$\dfrac{e^z \, 100}{1 + e^z \, 500} \sim [0, 1]$

$\dfrac{100}{101} \sim .99$

S curve

$\dfrac{0.2}{1.2} \sim 0.$  Sigmoid

0 Saving

Function

$\dfrac{e^z}{1 + e^z} \sim$ Sigmoid

output

$[0, 1] = S(z)$

$$\text{Sigmoid} = \frac{e^z}{1+e^z} \longleftarrow z = b_0 + b_1 x$$

$$= [0, 1] \quad \text{probability}$$



prob
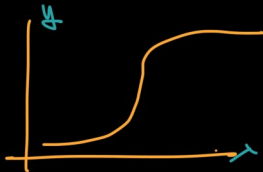
Credit → $0.8$

Credit → $0.2$

Default → $0.1$

$p > 0.5$
↳ No

$0.5$

Green
Bank

X

Sav Hou lu lu , yes, no .

new
users

Binary Classification

Model :

Guide.

probability

0 — 1

Decision

yes , no

Threshold    0.5

glm ( )
↑
family = "binomial"



# Evaluate Binary Model

## Confusion Matrix

|  | predict | |
|---|---|---|
| TRUE Negative | no | yes |
| Actual no | ② | ① → False Positive |
| Actual yes | ① | ① → TRUE POSITIVE |

False Negative

| x | y | ŷ | correct |
|---|---|---|---|
| - | 0 | 0 | 1 |
| - | 0 | 0 | 1 |
| . | 1 | 0 | 0 |
| . | 0 | 0 | 1 |
| . | 1 | 1 | 1 |

N → 0 1 0
P

3/5 = 60%

# CONFUSION MAT

$n = 1000$

2. $\dfrac{\text{Prec} \cdot \text{Recall}}{\text{Prec} + \text{Recall}}$

$F_1$

|  | prediction | |
|---|---|---|
|  | no | yes |
| **no** | TN 300 | FP 120 |
| **yes** | FN 80 | TP 500 |

Actual

no

yes
Default

$\text{Acc} = \dfrac{300 + 500}{1000}$

$= 80\%$

$\text{Recall} = \dfrac{TP}{FN + TP}$

$= \dfrac{500}{80 + 500}$

$= \dfrac{500}{580}$

$= 86\%$

$\text{precision} = \dfrac{TP}{FP + TP} = \dfrac{500}{120 + 500} = \dfrac{500}{620}$

$\approx 80.6\%$

1. Hypothesis · p·value / CI.
2. Binary Classification
3. Clustering — K·means.



Unsupervised Learning

# Regression * ( Parametric )

$$y = b_0 + b_1 \cdot x + b_2 x_2 + \ldots$$



$x_2$

Decision
Boundary

$x_3$  $x_1$

# 3. Clustering

# 3. Clustering K-Means (Iterative) + Random 1



$k = 2$

Centroid

# 3. Clustering



Compact.

$x_2$

Euclidean Distance

$x_1$

1. Random Centroid

2. Label #.
   Distance

Loop

3. move Centroid to centra Cluster

# 3. Clustering

S1 → S2 = kmeans (data)

center = 2

k = 2 ... 15

Optional R

ART

MKT

$x_2$

$x_1$

Business

$(1, 4)$

$$= \sqrt{16+4} = \sqrt{20} = 4.47$$

$$d = \sqrt{(1-5)^2 + (4-2)^2}$$

$$d = \sqrt{(x_1 - x_2)^2 + (y_1 - y_2)^2} \checkmark$$

Euclidean
Distance

$d = 4.47$

$(5, 2)$

$(1, 4, 5) \xrightarrow{\quad\quad d \quad\quad} (2, 10, 9)$

$$d = \sqrt{(1-2)^2 + (4-10)^2 + (5-9)^2} = \dots$$

1. Hypothesis  A/B Test : t.test()    =T.TEST )

2. Binary Classification $^{glm()}$ / Confusion Matrix
                                                    table()

3. Clustering kmeans

Data Skills
Science

Coding    Stats

Business

Model
Business.

# Project (Open).

#stat.ml
Datadose

1. data.world

$\longrightarrow$ churn

Explore data.

2. Build Model

Acc
Recall

Bonus 3. K-means $\Big[\,\longrightarrow$ Logistic Reg.

Precision
F1

$\longrightarrow$ glm( )