

To count transformed features :-

input \rightarrow 4 features (n)

$$X = [x_1 \ x_2 \ x_3 \ x_4]$$

Output \rightarrow features mapped with degree 2 (d)

all monomials (single terms) of input features where the sum of exponents is ≤ 2 .

$\phi(x) =$	$x_1 x_1$	
	$x_1 x_2$	
	$x_1 x_3$	
	$x_2 x_1$	
	$x_2 x_2$	
	$x_2 x_3$	
	$x_3 x_1$	
	$x_3 x_2$	
	$x_3 x_3$	

$$\text{Total count} = \binom{n+d}{d}$$

$$\langle \phi(x), \phi(y) \rangle = 9$$

same order as $\phi(x)$,

$\phi(x) =$	1	2	3	4	5	6	7	8	9

$\phi(y) =$	16	40	72	40	100	180	72	180	324
	$y_1 y_1$								
		$y_1 y_2$							
			$y_1 y_3$						
				$y_2 y_2$					
					$y_2 y_1$				
						$y_2 y_3$			
							$y_3 y_3$		
								$y_3 y_1$	
									$y_3 y_2$

$p = 9$. here,

$$\begin{aligned} \langle \phi(x), \phi(y) \rangle &= 16 + 40 + 72 + 40 + 100 \\ &\quad + 180 + 72 + 180 + 324 \end{aligned}$$

$$\langle \phi(x), \phi(y) \rangle = 1024$$

kernel trick lets us compute the dot product in higher dimension space without explicitly mapping to it.

Date: / /

$$K(x, y) = (\langle x, y \rangle)^2$$

$$K(x, y) = (\langle x, y \rangle)^2 = (x_1 + x_2 + x_3 + \dots)^2 = (4 + 10 + 18)^2 = (32)^2 = 1024$$

able to mimic what happened in higher dimension space on lower dimension space.

Types of kernel $\phi(x) = [\phi_1(x), \phi_2(x), \dots, \phi_D(x)]$

1. linear kernel, $K(x, y) = x \cdot y = \langle x, y \rangle$

2. Polynomial kernel, $K(x, y) = \langle x, y \rangle^d$

3. Radial Basis function (RBF) $K(x, y) = e^{-\gamma \|x - y\|^2}$

4. Laplacian kernel

5. Graph kernel

6. Fischer kernel

$$\theta = \beta_1 \phi(x^{(1)}) + \beta_2 \phi(x^{(2)}) + \dots + \beta_n \phi(x^{(n)})$$

$$\theta = \theta + \alpha \left(y^{(i)} - \theta^\top \phi(x^{(i)}) \right) \phi(x^{(i)})$$

$$\phi(x) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1^2 \\ x_1 x_2 \\ x_2^2 \\ \vdots \\ x_1 x_2 \end{bmatrix} \quad n$$

feature map.

$$\theta = \sum_{j=1}^n \beta_j \phi(x^{(j)})$$

θ as the linear combination of vectors $\phi(x^{(1)}) \phi(x^{(2)})$

11 | 03 | 2025

How is
correlation in kernel?

Example:

→ Dimension of β is $(n \times 1)$

SVM,

$$\begin{array}{c|c|c|c} x_1 & x_2 & y \\ \hline 1 & 2 & 0 \\ 3 & 4 & 0 \\ n=2 & & \end{array} \quad \phi(x) = \begin{bmatrix} x_1^2 \\ x_2^2 \\ x_1 x_2 \end{bmatrix}$$

$$x = [x_1, x_2]^T$$

$$x \in \mathbb{R}^2$$

$$\phi(x) \in \mathbb{R}^3 \quad \text{3 dimension}$$

$$\theta \in \mathbb{R}^3$$

$$\begin{bmatrix} 1 & 3 \\ 2 & 4 \end{bmatrix}_{2 \times 2}$$

$$\rightarrow \begin{bmatrix} \theta_1 & \theta_2 & \theta_3 \end{bmatrix}^T$$

$$\text{Find } \theta = \sum_{i=1}^n \beta_i \phi(x^{(i)}) \quad \beta_i \in \{\beta_1, \beta_2\}_{1 \times 2}^T$$

$$\beta_1 \phi(x^{(1)}) = \begin{bmatrix} 1 \\ 4 \end{bmatrix}, \quad \phi(x^{(2)}) = \begin{bmatrix} 9 \\ 16 \\ 12 \end{bmatrix}_{3 \times 1}$$

$$\theta = \beta_1 \phi(x^{(1)}) + \beta_2 \phi(x^{(2)})$$

$$\begin{aligned} \theta_1 &= \beta_1 + 9\beta_2 \\ \theta_2 &= 4\beta_1 + 16\beta_2 \\ \theta_3 &= 2\beta_1 + 12\beta_2 \end{aligned}$$

for all $i = 1, 2, \dots, n$.

10
grid

eqn (3) says that find the M which
the datapoints should lie on or
above M .

perpendicular.

How to find distances from point to plane?
How to calculate margin?

perpendicular (d) = $\frac{|B_0 + (\beta_1 x_1 + \beta_2 x_2)|}{\sqrt{\beta_1^2 + \beta_2^2}}$

shifting
coefficient;
each plane
weighted by
the coefficient

* three lines (hyperplane)

$$2x_1 + 3x_2 - 5 = 0 \quad - \textcircled{1}$$

[normalization]
factor

$$-x_1 + 4x_2 + 7 = 0 \quad - \textcircled{2}$$

$$5x_1 - 12x_2 + 10 = 0 \quad - \textcircled{3}$$

data	x_1	x_2	y
	3	4	+1
1	2	3	+1
1	-1	-1	-1
-2	-4	+1	-1

find the optimal
hyperplane?

hyperplane

$$2x_1 + 3x_2 - 5 = 0$$

$$\beta_1 = 2$$

$$\beta_2 = 3$$

$$\beta_0 = -5$$

Sample: $d_1 = \frac{|B_0 + \beta_1 x_1 + \beta_2 x_2|}{\sqrt{\beta_1^2 + \beta_2^2}}$

dist. of pt. 1 from
hyperplane

$$d_1 = \frac{|-5 + 2 \times 3 + 3 \times 4|}{\sqrt{13}} = \frac{-5 + 6 + 12}{\sqrt{13}} = \frac{13}{\sqrt{13}}$$

$$d_2 = \frac{|-5 + 2 \times 2 + 3 \times 3|}{\sqrt{13}} = \frac{-5 + 4 + 9}{\sqrt{13}} = \frac{8}{\sqrt{13}}$$

$$d_3 = \frac{|-5 + 2 \times 1 + 3 \times (-1)|}{\sqrt{13}} = \frac{|-5 + 2 - 3|}{\sqrt{13}} = \frac{6}{\sqrt{13}}$$

$$d_4 = \frac{|-5 + 2 \times (-2) + 3 \times (+1)|}{\sqrt{13}} = \frac{|-5 - 4 + 3|}{\sqrt{13}} = \frac{6}{\sqrt{13}}$$

minimum distance for plane (1) $\rightarrow \left\{ \frac{6}{\sqrt{13}} \right\} \rightarrow \{d_3, d_4\}$

$$\text{2nd hyperplane, } -x_1 + 4x_2 + 7 = 0$$

$$d_{PL} = d_1 = \frac{|B_0 + B_1 x_1 + B_2 x_2|}{\sqrt{B_1^2 + B_2^2}} = \frac{7 + (-1) \times 3 + 4 \times (4)}{\sqrt{9 + 16}} = \frac{7 - 3 + 16}{\sqrt{25}} = \frac{20}{5} = 4$$

$$d_2 = \frac{7 + (-1) \times 2 + 4 \times (3)}{\sqrt{17}} = \frac{17}{\sqrt{17}}$$

$$d_3 = \frac{7 + (-1) \times 1 + 4 \times (-1)}{\sqrt{17}} = \frac{2}{\sqrt{17}}$$

$$d_4 = \frac{|7 + (-1) \times (-2) + 10 \times (1+1)|}{\sqrt{13}} = \frac{7 + 2 + 4}{\sqrt{13}} = \frac{13}{\sqrt{13}} = \frac{13}{\sqrt{13}}$$

papergrid

Date: / /

$$\text{min dist for plane 2} = d_3 = \frac{2}{\sqrt{17}}$$

$$\text{8th hyperplane: } 5x_1 - 12x_2 + 10 = 0$$

$$d_1 = \frac{|10 + 5(3) - 12(4)|}{\sqrt{25+144}} = \frac{|10 + 15 - 48|}{\sqrt{25+144}} = \frac{23}{\sqrt{25+144}}$$

$$d_2 = \frac{|10 + 5(2) - 12(3)|}{\sqrt{25+144}} = \frac{16}{\sqrt{25+144}}$$

$$d_3 = \frac{|10 + 5(1) - 12(-1)|}{\sqrt{25+144}} = \frac{27}{\sqrt{25+144}}$$

$$d_4 = \frac{|10 + 5(-2) - 12(+1)|}{\sqrt{25+144}} = \frac{-12}{\sqrt{25+144}}$$

$$\text{min dist for plane 3} \rightarrow d_4 = \frac{-12}{\sqrt{25+144}}$$

Step 2. choose the min from computed \rightarrow dist of all hyperplanes.

$$\left[\frac{13}{\sqrt{13}}, \frac{2}{\sqrt{17}}, \frac{12}{\sqrt{25+144}}, \frac{-12}{\sqrt{25+144}} \right] \rightarrow \text{Margin for hyperplanes.}$$

Step 3 - ! choose the maximal margin hyperplane from step 2.

$$\therefore \max \left[1.66, 0.48, 0.923 \right] \Rightarrow 1.66 \text{ i.e. hyperplane (1)}$$

$$2x_1 + 3x_2 - 5 = 0$$

$$p = \text{proj}(a) \cdot b$$

$$= \frac{aa^T}{a^T a} \cdot b$$

$$= \begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1}$$

$$\begin{bmatrix} 1 & 2 \end{bmatrix}_{1 \times 2}^T$$

$$\begin{bmatrix} 3 \\ 4 \end{bmatrix}_{2 \times 1}$$

$$\begin{bmatrix} 1 & 2 \end{bmatrix}_{1 \times 2}$$

$$\begin{bmatrix} 1 \\ 2 \end{bmatrix}_{2 \times 1}$$

$$= \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}_{2 \times 2} \begin{bmatrix} 3 \\ 4 \end{bmatrix}_{2 \times 1}$$

$$\begin{bmatrix} 8 \\ 16 \end{bmatrix}_{2 \times 1} = d$$

$$= \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix}_{2 \times 2} \begin{bmatrix} 3 \\ 4 \end{bmatrix}_{2 \times 1}$$

$$= \begin{bmatrix} 3/5 + 8/5 \\ 6/5 + 16/5 \end{bmatrix}_{2 \times 1}$$

$$= \begin{bmatrix} 11/5 \\ 22/5 \end{bmatrix}_{2 \times 1} = \begin{bmatrix} 2.2 \\ 4.4 \end{bmatrix}$$

Heuristics based approach

Heuristics based approach

Algo random initialization

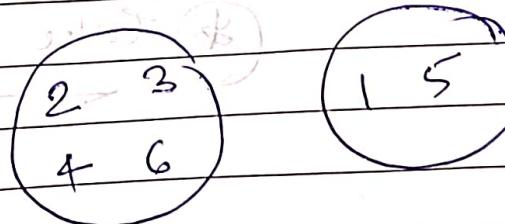
k-means → Greedy optimisation -

- ~~part~~

 - ① Randomly assign a number from 1 to k to each observation. These serve as initial cluster assignments.
 - ② Iterate until the cluster assignments stop changing.
 - a) for each of the k -clusters, compute the cluster centroid
 - b) assign each observation to the cluster whose centroid is closest (Euclidean).

$$k=2 \quad \text{Satz} \quad \boxed{\begin{array}{l} x_1 \\ x_2 \end{array}} \quad \boxed{\begin{array}{l} 13-14N \\ = \end{array}} \quad \boxed{\begin{array}{l} 1) \quad c_1 = \{2, 3, 4, 6\} \\ 2) \quad \dots \end{array}}$$

1	1	1
2	1	2
3	2	2
4	8	8
5	8	9
6	9	8



AFTER 1st iteration; ~~unbold~~

$$\text{centroid}_1 = \left[\frac{1 + 2 + 8 + 8}{4}, \frac{2 + 2 + 8 + 8}{4} \right] \\ = [5, 5]$$

Centroid 2 ; $\left[\begin{array}{c} 1+8 \\ 2 \end{array} \right] = \left[\begin{array}{c} 1+9 \\ 2 \end{array} \right]$

$\left[\begin{array}{c} 1+8 \\ 2 \end{array} \right] = \left[\begin{array}{c} 4.5 \\ 5 \end{array} \right]$

at each next column \rightarrow iteration continues

and find all three points in next iteration until

iteration ($i=1$)

$s_1 = \left[\begin{array}{c} 1 \\ 1 \end{array} \right] \rightarrow c_1 = \left[\begin{array}{c} 1 \\ 1 \end{array} \right]$

$c_1 = \left[\begin{array}{c} 5 \\ 5 \end{array} \right] \rightarrow c_2 = \left[\begin{array}{c} 4.5 \\ 5 \end{array} \right]$

Distribution method

$$\begin{aligned} d_{11} &= \sqrt{(1-5)^2 + (1-5)^2} = \sqrt{(-3-5)^2 + (4-5)^2} \\ &= \sqrt{16+16} = \sqrt{26.25} \\ &= \sqrt{32} \\ &= 4\sqrt{2} \\ &= 4 \times 1.41 \\ &= 5.64 \end{aligned}$$

$d_{12} = \sqrt{(1-4)^2 + (1-4)^2} = \sqrt{(-3-4)^2 + (4-4)^2} = \sqrt{25} = 5$

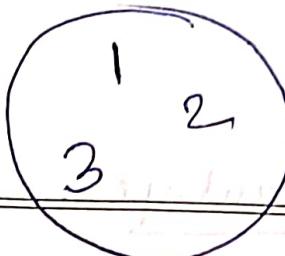
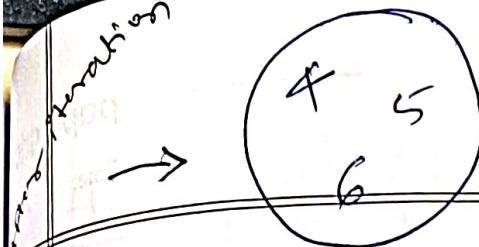
Since $\min(d_{11}, d_{12}) \therefore s_1 [1] \text{ is in cluster 2.}$

(3-1)

Iteration step with

$[8+2+5+5+8+3+5+9] = 45$

$\left[\begin{array}{c} \bar{x} \\ \bar{y} \end{array} \right] =$



$$d_{21} = \sqrt{(1-5)^2 + (2-5)^2}$$

$$= \sqrt{(-4)^2 + (-3)^2}$$

$$= \sqrt{16 + 9} = \sqrt{25} = 5$$

$$d_{22} = \sqrt{(1-4)^2 + (2-5)^2}$$

$$= \sqrt{(-3)^2 + (-3)^2}$$

$$= \sqrt{9 + 9} = \sqrt{18} = 3\sqrt{2} = 4.24 \cdot 10^{-1}$$

$$\min(d_{21}, d_{22})$$

$$d_{31} = \sqrt{(-3)^2 + (-3)^2} = \sqrt{18}$$

$$d_{32} = \sqrt{(-2-7)^2 + (2-5)^2} = \sqrt{62.5 + 9}$$

$$= \sqrt{71.5} = \sqrt{15.25}$$

$$d_{41} = \sqrt{(3)^2 + (3)^2} = \sqrt{18}$$

$$d_{42} = \sqrt{(3-5)^2 + (3)^2} = \sqrt{12.25 + 9}$$

$$= \sqrt{21.25}$$

$$d_{51} = \sqrt{(3)^2 + (4)^2} = \sqrt{25}$$

$$\min(d_{41}, d_{42})$$

$$d_{52} = \sqrt{(3-5)^2 + (4)^2} = \sqrt{12.25 + 16}$$

$$= \sqrt{28.25}$$

d_{51}

2

*. Centroid linkage



~~(1) Linkage based on dissimilarity between dot product.~~

Example: with 5 data points A-E find the dissimilarity matrix

$$\{(A, A) \in, (B, A) \in, (C, A) \in\} \text{ min} = 1.6$$

	x_1	x_2	Single linkage ?
A	1	1	Euclidean distance
B	2	1	as the measure
C	4	3	
D	5	4	nC_2 of pair-wise \rightarrow 10
E	6	5	= 10

Step 1 compute the distance matrix:

	A	B	C	D	E
A	0	1.6	3.6	4.9	6.4
B	1.6	0	2.82	4.24	5.64
C	3.6	2.82	0	1.41	2.82
D	4.9	4.24	1.41	0	1.41
E	6.4	5.64	2.82	1.41	0

→ starts to merge → answer

$$d(A-B) = \sqrt{(2-1)^2 + (1-1)^2} = \sqrt{(1)^2} = 1.$$

$$d(A-C) = \sqrt{(3-1)^2 + (2-1)^2} = \sqrt{13} = 3.6$$

$$d(A-D) = \sqrt{16 + 0.25} = 5$$

$$d(A-E) = \sqrt{25 + 16} = 6.4$$

$$d(B-C) = \sqrt{(-2)^2 + (-2)^2} = \sqrt{8} = 2\sqrt{2}$$

$$= 2 \times 1.41$$

$$= 2.82$$

$$d(B-D) = \sqrt{(-3)^2 + (-3)^2} = \sqrt{18} = 3\sqrt{2}$$

$$= 3 \times 1.41$$

$$= 4.24$$

$$d(B-E) = \sqrt{+16 + 16} = \sqrt{32} = 4\sqrt{2} = 5.64$$

$$d(D-E) = \sqrt{(-1)^2 + (-1)^2} = 2\sqrt{2} = 2 \times 1.41$$

* Merge clusters A & B,

New cluster :- AB C D E

AB C D E

AB	0	2.82	4.24	5.64
----	---	------	------	------

C	2.82	0	1.41	2.82
---	------	---	------	------

D	4.24	1.41	0	1.41
---	------	------	---	------

E	5.65	0	1.41	0
---	------	---	------	---

2.82

$$\min \left[D(c, A), D(c, B) \right] = \min \left[3, 6, 2.82 \right] = 2.82$$

papergrid

$$= 2.82$$

$$S^f_S = S^f_B = f(S-) + f(B-) = (2-3) b$$

~~$S^f_A = \dots$~~ cut line 1 (one cluster)

2.83

~~\dots~~ cut line 2 (two clusters)

$$1.41 S = S^f_A + f(A-) = f(S-) + f(A-) = (0-3) b$$

$$A =$$

$$B =$$

$$S = A \cup B \quad S^f_A = f(A-) + f(A+) = (3-2) b$$

#

$$(AB) \quad C \quad (DE) = f(-) - (AB)(CD(E)) \cdot c b$$

$$AB \quad C$$

$b \in \{0, 1\}$ & c far enough

$$AB$$

$$0$$

$$\cancel{0} \quad 0$$

$B \in \{0, 1\}$ & A far enough

$$C$$

$$0$$

$B \in \{0, 1\}$ & C far enough

$$DE$$

$$0$$

$D \in \{0, 1\}$ & E far enough

$S^f_A = 14.4 \quad S^f_B = 18.6 \quad S^f_C = 12.0 \quad S^f_D = 12.0$

$0 \quad 14.4 \quad 18.6 \quad 12.0 \quad 12.0$

x_1	x_2	x_3	class
1	0	1	spam
0	1	0	spam
1	1	0	spam
0	0	1	no spam
0	0	0	no spam
1	1	1	no spam
0	1	0	spam
1	0	0	spam
0	1	0	spam
1	0	0	spam
1	(all 0)	0	not spam

construct prob. table. parameter.
from Training data.

- Each of the row can be assigned a certain probability.
- Total combination possible are $2^3 = 8$; as each feature can take 0 or 1 (Boolean values).

x_1	x_2	x_3	$\text{Prob}(y=1) \rightarrow \text{spam}$	$\text{Prob}(y=0) \rightarrow \text{not spam}$
0	0	0	0	$\frac{1}{4} \rightarrow 1 \text{ in } 4 \text{ not spam.}$
0	0	1	$(1 - \frac{1}{4}) = \frac{3}{4}$	$\frac{1}{4} \rightarrow 1 \text{ in } 4 \text{ not spam.}$
0	1	0	$(2/7)$ → 2 hard same reason in spam	0
0	1	1	2/7	0
1	0	0	2/7	0
1	0	1	1/7	0
1	1	0	1/7	$\frac{1}{4}$
1	1	1	0	$\frac{1}{4}$

can be explained using Binomial (Bernoulli) distribution
as the values are 0 & 1.

1	Free win now (cont)	spam
2	Win a prize	spam
3	Hello How are you	not spam
4	Let's win it	not spam
5	Free lunch today?	not spam

feature

choose 2 words

	word "free"	word "win"	label
1)	Yes	Yes	spam
2)	No	Yes	spam
3)	No	No	not spam
4)	No	Yes	not spam
5)	Yes	No	not spam

Step 1: Calculate prior probability.

$$P(y=1 | x) = p(x|y=1) p(y=1)$$

and

$$\theta_y = \frac{1}{n} \sum_{i=1}^n I(y^{(i)} = 1), \Rightarrow \frac{2}{5} = P(\text{spam}) = 0.4$$

$$\theta_y = 1 - \frac{2}{5} = \frac{3}{5}$$

$$P(\text{not spam}) = \frac{3}{5} = 0.6$$

Step 2:

Compute the likelihood.

For spam class,

Free → appears → 1 of 2 spam class.

$$P(\text{Free} = \text{yes} | \text{spam} = 1) = \frac{1}{2}$$

$$P(\text{win} = \text{yes} | \text{spam} = 1) = \frac{2}{2} = 1$$

For not spam.

Free appears.

$$P(\text{Free} = \text{yes} | \text{not spam}) = \frac{1}{3}$$

$$P(\text{win} = \text{yes} | \text{not spam}) = \frac{1}{3}$$

Inference phase.

Test message = Free Win. ?
 ↓ ↓
 Yes Yes.

$$P(\text{spam} | x) = P(x | \text{spam}) P(\text{spam})$$

$$P(x_1 = \text{yes} | \text{spam} = \text{yes})$$

$$P(x | \text{spam}) = P(x_1, x_2 | \text{spam}) \quad \begin{cases} P(x_2 = \text{yes} | \\ \text{spam} = \text{yes}) \end{cases}$$

$$= \{P(x_1 | \text{spam}), P(x_2 | \text{spam})\}$$

Free = yes

$$= \frac{1}{2} \times \frac{1}{3} = \frac{1}{6}$$

win = yes

$$P(\text{spam} | x) = \frac{1}{2} \times P(\text{spam}) = \frac{1}{2} \times \frac{2}{5}$$

$$= \frac{1}{5} = 0.2$$

$$P(\text{not spam} | x) = \frac{1}{3}, \frac{1}{3}, \frac{3}{5} \stackrel{\text{sum}}{=} \frac{1}{15} = 0.066$$

Date: / / papergrid

$0.2 > 0.066$

~~if compare both prob & assign the max~~

~~↓ message belongs to spam class (0.2)~~

$$\underline{1} = \frac{2}{\Sigma} = (\text{message from } \text{spam} = \text{msg}) \cdot 9$$

~~message from not~~

$$\underline{1} = (\text{message from } \text{not-spam} = \text{msg}) \cdot 9$$

$$\underline{1} = (\text{message from } \text{not-spam} = \text{msg}) \cdot 9$$

~~message from not-spam = message from not~~

$$= (\text{msg}_2) \cdot (\text{msg}_2 \times 1) \cdot 9 = (x | \text{msg}_2) \cdot 9$$

~~message from not~~

$$= (\text{msg}_2) \cdot (\text{msg}_2 \times 1) \cdot 9 = (\text{msg}_2 | x) \cdot 9$$

$$((\text{msg}_2 | x) \cdot 9 \cdot (\text{msg}_2 | x) \cdot 9) =$$

$$= 0.8 \cdot 0.8 = 0.64$$

$$= \sqrt{0.64} = (\text{msg}_2) \cdot 9 \cdot 1 = (x | \text{msg}_2) \cdot 9$$

$$= \frac{1}{2}$$

Date: 21/03/2025

<u>Observation</u>	<u>Actual Labels</u>	<u>Predicted score</u>	<u>Predicted label</u>
1	1	0.85 (TP)	1
2	0	0.60	0.185 * (FP)
3	1	0.70 (TP)	1
4	x	0.40	1.0 0 * (FN)
5	0	0.55	0.1 * (FP)
6	1	0.50	1 0 1
7	0	0.65	0.0 * (FP)
8	0.25	0.35 (TN)	0 0 0
9	x	0.60 (TP)	1 0 1
10	0	0.20 (TN)	0 0 0
(0.5)			(0.5)

- 1) Accuracy
- 2) Precision
- 3) Recall or sensitivity or True positive rate
- 4) Specificity
- 5) False positive rate (1 - specificity)
- 6) F1-score
- 7) ROC plot
- 8) AUC.

Confusion matrix

		original positive	original -ve
		(TP)	(TN)
predicted	original	1, 3, 9 (4)	8, 10 (3)
(TP)	original +ve	1 (label 6)	2, 7, 5
(TN)	original -ve	4 (label 4)	8, 10 (2)

$$\text{Accuracy} = \frac{(TP + TN)}{\text{Total}}$$

$$= \frac{8+2}{10} = 0.8$$

$$\text{Precision} = \frac{TP}{TP + FP} = \frac{4}{4+2} = 0.57$$

$$\text{Sensitivity} = \frac{TP}{TP + FN} = \frac{4}{4+1} = 0.8 = TPR,$$

$$\text{Specificity} = \frac{TN}{TN + FP} = \frac{5}{5+2} = 0.6$$

$$\text{False positive rate} = (1 - 0.6) = 0.4 = 0.6.$$

$$\text{F1-score} = \frac{2 \cdot 0.57 \cdot 0.8}{0.57 + 0.8} = 0.66.$$

Threshold

	TPR	FPR
0.7	0.4	0
0.5	0.8	0.6
0.4	(4)	0.6
0.2	1	1

	OP	ON	TPR	FPR	specificity
TP	2	0	$\frac{2}{2} = 0.4$	$\frac{2}{5} = 0.4$	$\frac{5}{5} = 1$
PN	3	5			$FPR = 0$

	OP	ON
--	----	----

PP	5	3
PN	0	2

precision = $\frac{5}{8} = 0.625$

sensitivity = $\frac{5}{8} = 0.625$

FPR = $1 - 0.625 = 0.375$

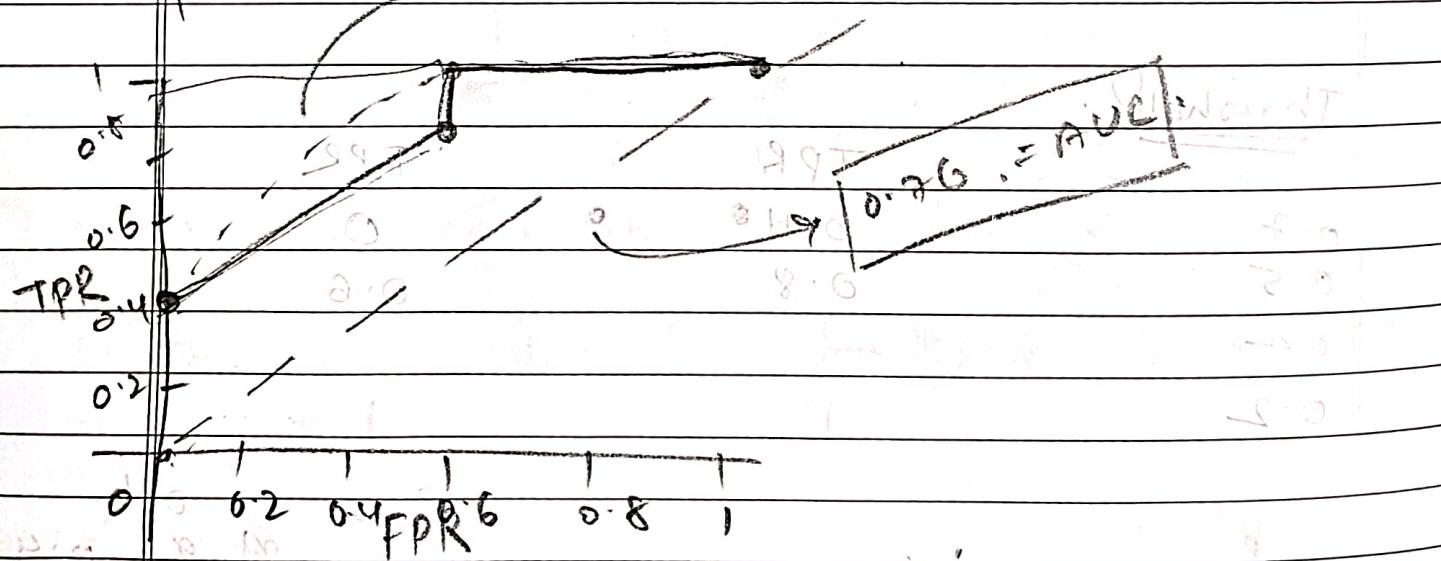
FPR = $1 - 0.6 = 0.4$

AUC = $(0.625 + 1) \times 0.25 = 0.875$

	OP	ON
--	----	----

PP	5	5
PN	0	0.25

FPR = 0.25



$AUC = \frac{1}{4} \sum_{i=1}^4 (FPR_i - FPR_{i-1}) * (TPR_i + TPR_{i-1})$

$$\sum_{i=1}^4 (0.625 - 0) * (0.625 + 0.4) = 0.375 * 1.0625 = 0.400625$$

(1) $(0 - 0) * = 0$

(2) $(0.625) * (0.875 + 0.625) = 0.375 * 1.5 = 0.5625$

x_1	x_2	y
1	2	3
2	1	4
3	3	5

$$h_0(x) = \theta_0 + \theta_1 x_1 + \dots$$

$$\theta_0 = \theta_1 = \theta_2 = 0$$

$$\alpha = 0.1$$

$$(1) h_0((x_1, x_2)) = \theta_0 + \theta_1 x_1 + \theta_2 x_2$$

① compute the prediction. (for each x_i , every sample).

$$\begin{aligned} h_0(x^{(1)}) &= \theta_0 + \theta_1 x_1 + \theta_2 x_2 \\ &= 0.1 + 0.1 + 0.2 \\ &= 0.4 \end{aligned}$$

$$\begin{aligned} h_0(x^{(2)}) &= 0 + 0 + 0 \\ &= 0 \end{aligned}$$

$$\begin{aligned} h_0(x^{(3)}) &= 0 + 0 + 0 \\ &= 0 \end{aligned}$$

$$f(\theta) = \frac{1}{2} [(0-3)^2 + (0-4)^2 + (0-5)^2]$$

② compute the gradients.

i) compute errors / residuals.

$$e^{(1)} = h_0(x^{(1)}) - y^{(1)} = 0.4 - 3 = -2.6$$

$$e^{(2)} = 0 - 4 = -4$$

$$e^{(3)} = 0 - 5 = -5$$

ii) compute gradients.

$$\frac{\partial f}{\partial \theta_0} = \sum_{i=1}^n (h_0(x^{(i)}) - y^{(i)}) x_0^{(i)}$$

$$= e^1 \cdot x_0^1 + e^2 \cdot x_0^2 + e^3 \cdot x_0^3$$

$$= -3 \times 1 + (-4) \times 1 + (-5) \times 1 \\ = -12$$

$$\frac{\partial J}{\partial \theta_1} = \sum_{i=1}^3 (h_\theta(x^{(i)}) - y^{(i)}) x_1^{(i)}$$

$$= -3 \times 1 + -4 \times 2 + -5 \times 3 \\ = -26$$

$$\frac{\partial J}{\partial \theta_2} = \sum_{i=1}^3 (h_\theta(x^{(i)}) - y^{(i)}) x_2^{(i)}$$

$$= -3 \times 2 + -4 \times 1 + -5 \times 3 \\ = -6 - 4 - 15 \\ = -25$$

③

Update parameters

$$\theta_0 := \theta_0 - \alpha \frac{\partial J}{\partial \theta_0}$$

$$= 0 - 0.1 \times (-12) \\ = 0 + 1.2$$

$$\boxed{\theta_0 = 1.2}$$

$$\theta_1 := \theta_1 - \alpha \frac{\partial J}{\partial \theta_1}$$

$$\boxed{\theta_1 = 2.6}$$

$$\theta_2 := \theta_2 - \alpha \frac{\delta J}{\delta \theta_2} = 0 - 0.1 \times (-2r)$$

$$\boxed{\theta_2 = +2.5}$$

New parameters $\rightarrow \theta_0 = 1.2$

$$\theta_1 = 2.6$$

$$\theta_2 = 2.5$$

After 1st iteration.

(P)

compute J value? $\rightarrow J = 94.95$ → missed the global minimum.

due to the step size i.e. 0.06.

After 2nd iteration

$$\theta_0 = -1.02, \theta_1 = -2.41, \theta_2 = -2.6.$$

$$J = 25678.56.$$

Reduce the learning rate α from 0.1 to 0.01
iterate again.

$$\theta_0 = 0, \theta_1 = 0, \theta_2 = 0$$

$$J = 10.81$$

After 1st iteration: $\theta_0 = 0.19, \theta_1 = 0.45, \theta_2 = 0.41$

$$\boxed{J = 5.4}$$

After 2nd iteration: $(\theta_0 = 0.25, \theta_1 = 0.55, \theta_2 = 0.53)$

$$\boxed{J = 3}$$

↳ decreases!

k-fold cross-validation

	x_1	x_2	y	3-fold CV
1	5	-1	1	Acc.
2	0.5	1.2	0	std.
	1	2	1	
K	-3	-2	1	
	4	0.1	0	

$$\text{Fold 1 } [0_1, 0_2] = [-1.8, 2.8]$$

$$\text{Fold 2 } [0_1, 0_2] = [2.1, 3.1]$$

$$\text{Fold 3 } [0_1, 0_2] = [1.9, 4]$$

Compute Avg. Acc & Avg. std deviation?

Fold I $\xrightarrow{k=3}$ train set: {I, II, III}

$\overline{x_1}$	$\overline{x_2}$	\overline{y}	test set: {III}
2	-1	1	
0.5	1.2	0	Acc = ? 100%

I	2	1	3	Fold 2
-3	-2	1		50%

train set: {I, III}

Fold 4 II 0.1 Ratio $\xrightarrow{\text{train set: } \{I\}, \text{ test set: } \{II\}}$

Test set: {II}

Fold 2 0.50% 50%

Fold 3 Train set: {I} Train set: {II, III}

Test set: {III} Test set: {I}

$$\text{Avg} = 66.67\%$$

compute $\underline{\theta^T x}$

$$\underline{\theta^T x} =$$

$$= \begin{bmatrix} -1.8 \\ 2.8 \\ 1.2 \end{bmatrix} \times \begin{bmatrix} 2.3 \\ -1 \\ 1.2 \end{bmatrix} = 0.5 \times 1.2 + 1.2 = 0$$

$$\theta^T x = \begin{bmatrix} 2.4 - 1.8 & + (-1) \times 2.8 \\ 0.5 \times -1.8 & + 2(1.2) \times 2.8 \\ 1 \times -1.8 & + 1.02 \times -2.8 \\ -3 \times 1.2 & + (-2) \times 2.8 \end{bmatrix} = \begin{bmatrix} -6.4 \\ 3.46 \\ 3.18 \\ -0.2 \end{bmatrix}$$

$$g(z) = \frac{1}{1 + e^{-6.4}} = 0.001658$$

$$g(z) = \begin{bmatrix} 0.001658 \\ 0.998 \\ 0.998 \\ 0.998 \end{bmatrix}$$

threshold = 0.15

$$\hat{y} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$y = \begin{bmatrix} 1 \\ 0 \\ 1 \\ 1 \end{bmatrix}$$

Test fold I

$$x = \begin{bmatrix} 4 & 0.1 \end{bmatrix} \quad \delta^T = \begin{bmatrix} -1 & \text{papergrid} \\ 0 & \text{Date: } \end{bmatrix}$$

$$\begin{aligned} x &= \\ &= -1.8 \times 4 + 2.8 \times 0.1 \\ &= -7.2 + 0.28 \\ &= -6.92 \end{aligned}$$

$$g(z) = 0.0009$$

$$\hat{y} = 0 \quad | \quad y = 0$$

Acc = 100%

$$\textcircled{II} \quad \alpha = \begin{bmatrix} 2 & -1 \\ 0.5 & 1.2 \end{bmatrix} \quad \delta^T = \begin{bmatrix} 1 \\ -3 \end{bmatrix} \quad \begin{bmatrix} -2.1 \\ 3.18 \end{bmatrix} \quad 2 \times 2$$

$$\delta^T x = \begin{bmatrix} -2.1 \\ 3.18 \end{bmatrix}_{1 \times 2} \quad \begin{bmatrix} 2 & -1 \\ 0.5 & 1.2 \end{bmatrix}_{2 \times 2}$$

$$\begin{bmatrix} 1.1 \\ 4.72 \end{bmatrix}$$

$$g(z) = \begin{bmatrix} 0.249 \\ 0.0431 \end{bmatrix}$$

$$\hat{y} = \begin{bmatrix} 0 \\ 0 \end{bmatrix} \quad y = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

Acc = 50%

IV

$$D^T = \begin{bmatrix} 1.9 \\ 4 \end{bmatrix} \quad X = \begin{bmatrix} 1 & 2 \\ -3 & -2 \end{bmatrix}$$

$$D^T X = \begin{bmatrix} 1 & 2 \\ -3 & -2 \end{bmatrix} \begin{bmatrix} 1.9 \\ 4 \end{bmatrix}$$

2x2

1x2 2x1

$$D^T X = \begin{bmatrix} 9.9 \\ -13.7 \end{bmatrix}$$

$$g(2) = \begin{bmatrix} 0.999 \\ 0.0001 \end{bmatrix} \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$

$$ACC = 50\%$$

$$\Rightarrow Avg. ACC = \frac{50\% + 100\% + 70\%}{3} = \frac{220}{3}\% \quad 66.67\%$$

Std. deviation

why did we not use SSE in logistic?

Entropy

$$H(A) = - \sum_{i=1}^K p_i \log_2 (p_i)$$

A good split reduces entropy of the label.
∴ ↑ IG.

* In restraint data set, how to choose between Type & pattern to be the starting node?

for the tree

→ Measure the homogeneity[↑], by measuring the entropy.

IG = Information Gain \rightarrow always true, converges to zero at extreme thresholds (very low or very high)

$p = 6, n = 6.$ $p \rightarrow +ve$ samples
 $n \rightarrow -ve$ samples
0, if both labels are of same class

$$H\left(\frac{p}{p+n}, \frac{n}{p+n}\right) = -\frac{p}{p+n} \log_2 \frac{p}{p+n} - \frac{n}{p+n} \log_2 \frac{n}{p+n}$$

$$\text{entropy } H(A) = -\sum_{i=1}^K p_i \log_2 p_i \geq \text{bits encoding}$$

$$H(A) = -\sum_{i=1}^K p_i \log_2 p_i$$

$$H\left(\frac{6}{12}, \frac{6}{12}\right) = -\frac{6}{12} \log_2 \frac{6}{12} - \frac{6}{12} \log_2 \frac{6}{12}$$

$$= -0.5 \log_2 (0.5) - 0.5 \log_2 (0.5)$$

$$= -0.5 \times (-1) - (0.5 \times -1)$$

$$= 0.5 + 0.5$$

$$= 1$$

12 data points

A ← patrons?

none

some

full

$$E(A) =$$

$$E(A)$$

00

00

0000

000000

000000

Type?

French

Burger

00

00

Expected Entropy

$$EH(A) = - \sum_{i=1}^{p+n} \left(\frac{P_i + n_i}{p+n} \right) H \left(\frac{P_i}{P_i + n_i}, \frac{n_i}{P_i + n_i} \right)$$

$I(A)$

$$= H \left(\frac{P}{P+n}, \frac{n}{P+n} \right) - EH(A)$$

information

gain

how much impurity I can shed?

(we want maximum information)

criteria to choose the feature

Date: 11/02/2025

$$H(\text{patrons}) = H\left(\frac{P}{P+n}, \frac{n}{P+n}\right)$$

$$= H\left(\frac{6}{12}, \frac{6}{12}\right) \quad \checkmark$$

$$= 1$$

$$EH(\text{patrons}) = \sum_{i=1}^3 \left(\frac{p_i + n_i}{P+n} \right) H\left(\frac{p_i}{p_i + n_i}, \frac{n_i}{p_i + n_i}\right)$$

$$\begin{aligned} EH(\text{patrons}) &= - \left[\frac{p_1 + n_1}{12} H\left(\frac{0}{2}, \frac{2}{2}\right) + \frac{p_2 + n_2}{12} H\left(\frac{4}{4}, \frac{4}{4}\right) \right. \\ &\quad \left. + \frac{p_3 + n_3}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right] \\ &= - \left[\frac{2+0}{12} H\left(\frac{0}{2}, \frac{2}{2}\right) + \frac{4+0}{12} H\left(\frac{4}{4}, \frac{0}{4}\right) + \frac{6}{12} H\left(\frac{2}{6}, \frac{4}{6}\right) \right] \end{aligned}$$

~~0.67~~

~~so 4 different cases and only 2 cases are valid right?~~

$$2 \times \left(-\frac{0}{2} \log_2 \frac{0}{2} - \frac{0}{2} \log_2 \frac{2}{2} \right)$$

~~so 4 different cases and only 2 cases are valid right?~~

$$= -\frac{1}{6} (-0 \log_2 0 - 1 \log_2 1) + \frac{1}{3} (-1 \log_2 1 - 0 \log_2 0)$$

$$+ \frac{1}{2} \left(-\frac{1}{3} \log_2 \frac{1}{2} - \frac{2}{3} \log_2 \frac{2}{3} \right)$$

$$= 0.1 \text{ (with 2nd case)} + 1 \log_2 (3) - 2 \log_2 (0.67)$$

$$= \frac{1}{6} \left[\frac{1}{3} \times 1.5849 - \frac{2}{3} \times (-0.577) \right]$$

$$EH(\text{Patrons}) = 0.546, 0.456$$

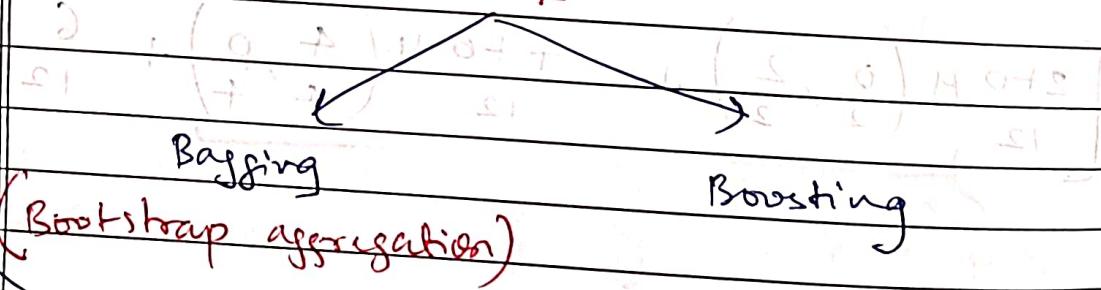
$$\begin{aligned} IG(\text{Patrons}) &= H(\text{Patrons}) - EH(\text{Patrons}) \\ &= 1 - 0.456 \\ &= 0.541 \end{aligned}$$

* $IG(\text{Type}) = 0.$

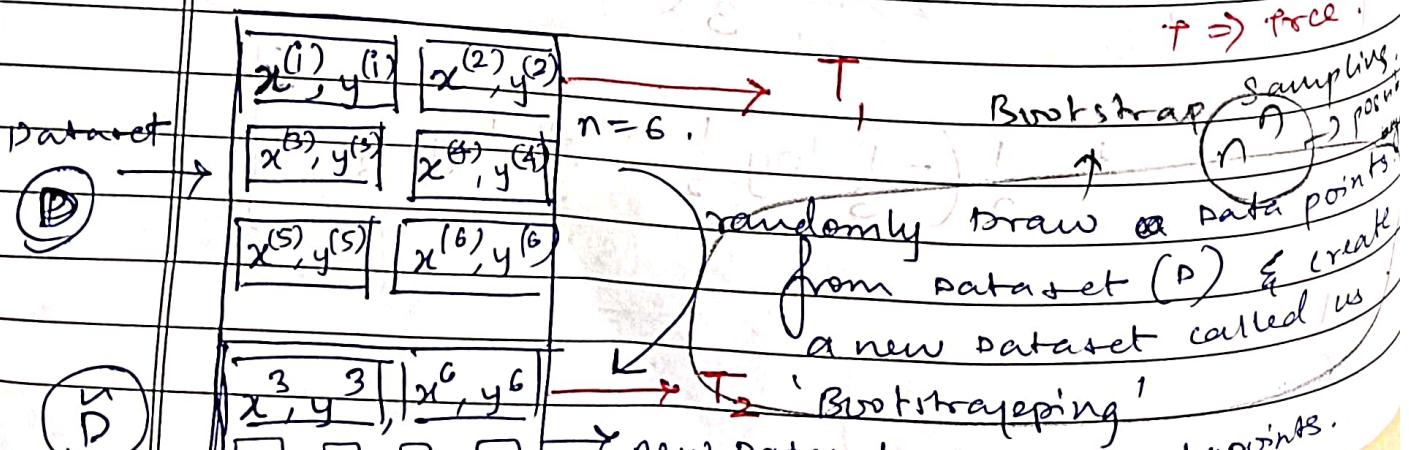
Gini index: → measure like Information gain.

IF the Information gain is equal for two features then we choose arbitrarily.

Ensemble methods.



multiple trees or models can be constructed for given dataset. → individual tree or model are referred to as weak learners → we try to aggregate the learning using Ensemble method.



Dataset

x_1	x_2	y
1	2	4
2	3	5
3	4	6
4	5	7

i) Initialize $F(x) = 0$, $\alpha = (x)$

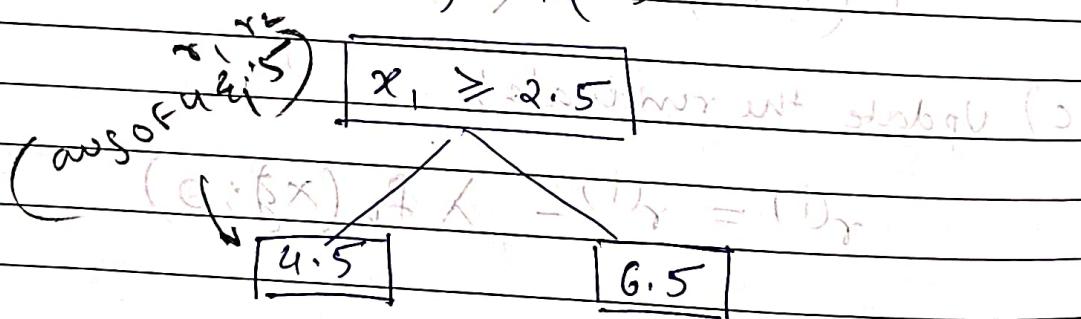
$$f_0(x) = 0$$

(a) $\gamma = r - \alpha f_0(x)$

Iteration 1

	x_1	x_2	r
$\gamma \leftarrow (x)$	1	2	4
	2	3	5
D	3	4	6
	4	5	7

(a) Fit a model, $x \rightarrow f_1(x) \rightarrow \gamma$



b) update the model, $\lambda = 0.1$

$$F(x) = f_0(x) + \lambda f_1(x) \quad \text{assumed}$$

c) update the residuals.

$$\gamma^{(1)} = r^{(1)} - \lambda f_1(x^{(1)}) = 4 - (0.1 \times 4.5) = 3.55$$

$$\gamma^{(2)} = r^{(2)} - \lambda f_1(x^{(2)}) = 5 - (0.1 \times 4.5) = 4.55$$

$$r^{(3)} = r^{(3)} - \lambda f_1(x^{(3)}) = 6 - (0.1 \times 6.5) = 5.35$$

$$r^{(4)} = r^{(4)} - \lambda f_1(x^{(4)}) = 7 - (0.1 \times 6.5) = 6.35$$

Iteration 2

$$r = \begin{bmatrix} 3.55 \\ 4.55 \\ 5.35 \\ 6.35 \end{bmatrix}$$

x_1	x_2	r
1	2	3.55
2	3	4.55
3	4	5.35
4	5	6.35

a) fit a model.

$$x_1 \geq 2.5$$

Avg
of
3.55

$$4.05$$

Avg
of
5.35 &
6.35

b) update the model.

$$F(x) = f_0(x) + \lambda f_1(x) + \lambda f_2(x)$$

c) update the residuals

$$r^{(1)} = r^{(1)} - \lambda f_2(x^{(1)}) = 3.55 - [0.1 \times 4.05]$$

$$r^{(2)} = r^{(2)} - \lambda f_2(x^{(2)}) = 4.145$$

$$r^{(3)} = r^{(3)} - \lambda f_2(x^{(3)}) = 4.765$$

$$r^{(4)} = r^{(4)} - \lambda f_2(x^{(4)}) = 5.765$$

Residuals decreases !! ✓

Rinal. is $F(x, \theta) = f_0(x) + \lambda f_1(x) + \lambda f_2(x)$
 $+ \lambda f_3(x) + \lambda f_4(x)$.