

Logistic Regression

Till now: ml \rightarrow supervised \rightarrow classification \rightarrow Logistic Reg.

e.g.: # no. of hours studied \rightarrow Student Pass/Fail

Hours Studied	Pass/Fail
1	F
2	P
3	P
1	P
2	F
4	P
5	P
1	F
2	F
4	F
3	F

What we want

↓ train

new data \rightarrow [Model] \rightarrow P/F or (1/0)

Accuracy ↑

This comes under classification problem

Q1: Predict whether person is fraud/not fraud bank

Q2: Predict whether mail is spam/not

Q3: Predict weather mail is spam/not

Classification

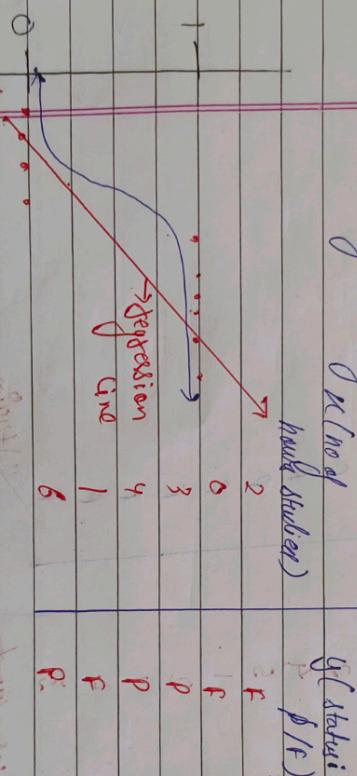
Binaria
→ output → 2 categories

more than 2 categories
→

Water → fail
Risk → high
medium
Pass

There are different algorithm to solve classification problem.

1) Logistic Regression



→ Can we solve it using regression line?

No → because nature of problem statement is diff.
in regression target value should continuous whereas
it is discrete

best fit line change due to presence of outliers

→ due to outliers like extremes

→ line should be this but

best fit line.

$\geq 0.5 \rightarrow 1$

$\geq 0.5 \rightarrow 1$

$< 0.5 \rightarrow 0$

\hookrightarrow after comes corner 0 that due to change
in line model predicted as 1

\hookrightarrow This should be best fit line but to
accommodate outliers, fitted fit line changed
its this.

\hookrightarrow Related to Bad point

\hookrightarrow Range of prediction line $\rightarrow (-\infty, +\infty)$

1. "Classification" $\rightarrow [0, 1]$

• prediction changes rapidly suddenly

\hookrightarrow In above eg. 2.3 \rightarrow Pass \rightarrow Fail

\hookrightarrow This is just maximized in real life, it
should change gradually

How to solve the problem?

① Squaring of error to [0, 1]

George Hinton

① prediction changing suddenly just like we call sigmoid function

Q.

↳ Sigmoid function
↳ Logistic regression

formula for
this func:

$$f(z) = \frac{1}{1 + e^{-z}} \rightarrow \text{Range } [0, 1]$$

Expt. Q. $z = -100$ divide exp/ with exp 0

$$\frac{1}{1 + e^{(-100)}} = \frac{1}{1 + e^{100}} = \frac{1}{1 + 100} = \frac{1}{101} \approx 0$$

↳ anything to power(0) is (0)

Q. $z = 0$

$$\frac{1}{1 + e^{-0}} = \frac{1}{1 + e^0} = \frac{1}{1 + 1} = \frac{1}{2} = 0.5$$

$$\therefore \frac{1}{1 + 0} = 1 \therefore 1$$

for ($z = -\infty$, ≈ 0) predicted

$$f(z = -\infty, z = 0) \approx 0$$

①

$$h_{\theta}(x) = \theta_0 + \theta_1 x \rightarrow \text{best fit line}$$

Now via Logistic regression:

$$h_{\theta}(x) = \frac{1}{1+e^{-(\theta_0 + \theta_1 x)}} \quad z = \theta_0 + \theta_1 x$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \ln(1+e^{-y_i(\theta_0 + \theta_1 x_i)}) \rightarrow \text{Logistic regression}$$

- ② To get best (θ_0, θ_1) we minimize cost function ($J(\theta)$)

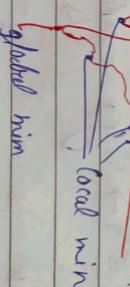
CF of Logistic regression?

$$\text{In LRM} \rightarrow CF = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (\hat{y}_i = h_{\theta}(x_i))$$

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n \ln(1+e^{-y_i(\theta_0 + \theta_1 x_i)})$$

Convex \Rightarrow curve at 2 point
 $\rightarrow 1 \rightarrow \text{minimum}, 1 \rightarrow \text{maximum}$

- Non convex \Rightarrow more than one minima.



global min

local min

LOG LOSS FUN

$$J(\theta_0, \theta_1) \Rightarrow y_i \log(h_{\theta}(x_i)) - (1-y_i) \log(1-h_{\theta}(x_i))$$

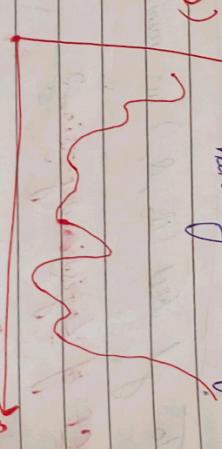
$$\text{If } y = 1 \Rightarrow J(\theta_0, \theta_1) = -y \log(h_{\theta}(x_i))$$

In Logistic Regression

$$J(\theta_0, \theta_1) = \frac{1}{n} \sum_{i=1}^n (y_i - h_{\theta}(x_i))^2$$

$$h_{\theta}(u) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 u)}}$$

notes of CF in Logistic regression



Convex fun. \rightarrow convex cost fun is preferred
that in our case we have non convex

function.
So how is it turned to get a new cost fun.

\hookrightarrow called as C log loss function.)

Log loss fun $\Rightarrow J(\theta_0, \theta_1) \rightarrow -y_i \log(h_{\theta}(x_i)) - (1-y_i) \log(1-h_{\theta}(x_i))$

$$\rightarrow h_{\theta}(u_i) = \frac{1}{1 + e^{-(\theta_0 + \theta_1 u)}}$$

$\rightarrow y_i = \text{actual value}$
 \hookrightarrow possible value $\rightarrow 0, 1$

(1)

$$y \cdot y = 1$$

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x_i)) - (1-y) \log(1-h_{\theta}(x_i))$$

$$J(\theta_0, \theta_1) = -y \log(h_{\theta}(x_i))$$

$$y=0 \quad J(\alpha_0, \alpha_i) = -\log(1 - H_0(u_i))$$

PAGE NO.	
DATE	

$$\begin{cases} \theta \\ g = 0 \end{cases}$$

$$J(\alpha_0, Q) = -g_i \log(H_0(u_i)) - (1-g_i) \log(1-H_0(u_i))$$

$$J(\alpha_0, \alpha_i) = -(1-g_i) \log(1-H_0(u_i))$$

$$J(\alpha_0) = -\log(1 - H_0(u_i))$$

To minimize J we change (α_0, Q) by using
Convergence Algorithm

→ Repeat until convergent.

$$\begin{cases} \theta : g_i - w_i J(\alpha_0, Q) \\ Q \end{cases}$$

To get optimal α_0, Q ,
For one independent variable.

$$\text{for multiple } J' \quad H_0(u) = \frac{1}{1 + e^{-(\alpha_0 + \alpha_1 u_1 + \alpha_2 u_2 + \dots + \alpha_n u_n)}}$$

Note → Logistic function range → [0 to 1]

$$\text{output of regression model} \rightarrow P = \frac{1}{1 + e^{-(\alpha_0 + \alpha_1 u)}}$$

between [0 to 1] Some i.e. P values. Then we take
cut off of 0.5 to put them in class 0 or 1

Logistic Regression with regularization

$$CF = J(\theta_0, \theta_1) - g_i \log(H_0(u_i)) - (1-g_i) \log(1-H_0(u_i))$$

$$(H_0(u_i) - 1)^2 + (1-H_0(u_i))^2 = H_0(u_i) = \frac{1}{1+e^{-f(u_i) + \theta_0 + \theta_1 u_i}} \quad \text{Minimizing}$$

1) Ridge

$$J(\theta_0, \theta_1) = -y_i \log(H_0(u_i)) - (1-y_i) \log(1-H_0(u_i)) + \lambda \sum_{j=1}^d \theta_j^2 \quad \text{to regularization}$$

reduce overfitting

$$\lambda \sum_{j=1}^d \theta_j^2$$

2) Lasso

$$J(\theta_0, \theta_1) = -y_i \log(H_0(u_i)) - (1-y_i) \log(1-H_0(u_i)) + \lambda \sum_{j=1}^d |\theta_j| \quad \text{to feature selection}$$

①

3) Elastic Net

$$\lambda \sum_{j=1}^d |1\theta_j|$$

$$J(\theta_0, \theta_1) = CF + L_1 + L_2$$

L_1 Reg. L_2 Reg.

both reduce overfitting & get better solution

in logistic regression

parameters we have $[C = \lambda]$

1

Classification Evaluation Metrics

- Classification measures

Recall

- $\frac{TP}{TP + FN}$
- True Positive Rate (Sensitivity)
- False Positive Rate
- True Negative Rate (Specificity)
- $\frac{TN}{TN + FP}$

Precision

Precision-Recall

Conclusien Mefric

Based on Some
work

class A & B are separately
dissolved

- The ~~jack~~ ^{gated} ~~o.~~ \rightarrow False/wrong classification

1

1

1

1

100

卷之三

$$\text{Gard} = 1 \quad \text{Spray} = 1$$

you = 0, group = 1

$\rho_{\text{ref}} = 0$, $\eta_{\text{ref}} = 0$ (aprox)

$$f = 1 \quad \text{and} \quad g = 0$$

Dudu cheet

Julie

卷之三

To make this in number of representing

Achmea

10

Value

negative

1

PAGE No. _____
DATE _____

$$\text{Total Prediction} = \frac{TP + TN}{TP + TN + FP + FN}$$

		0	1
y pred	0	2	
	1	1	2
			$\rightarrow TP + TN$

② Accuracy \Rightarrow How many are correctly predicted from all data point.

$$\begin{aligned}\text{Accuracy} &= \frac{TP + TN}{\text{All prediction}} \\ &= \frac{2 + 2}{4} \\ &= 1\end{aligned}$$

$$\begin{aligned}\underline{\text{Accuracy}} &= \frac{TP + TN}{TP + FP + TN + FN} \\ &= \frac{4}{7}\end{aligned}$$

Opposite of Accuracy is misclassification $\Rightarrow \frac{FP + FN}{\text{Total Prediction}}$

$$\Rightarrow \underline{FP + FN}$$

$$\Rightarrow \frac{FP + FN}{TP + FP + TN + FN}$$

$$\Rightarrow 1 - \underline{\text{Accuracy}}$$

③ Precision

why precision?

Before going deep let's understand why precision.

1000 pairs of data

$\rightarrow 900 \rightarrow \text{class } 0$

$\rightarrow 100 \rightarrow \text{class } 1$

Goat

1	0	6
0	100	900
100	900	

confusion matrix \rightarrow $\frac{1}{100} = 1\%$

PAGE No.

DATE		
------	--	--

Here our model predicted all data points as zero, all class 1 are wrongly predicted even after that due to the class imbalance our accuracy comes out to be 90%.

$$\text{Accuracy} = \frac{900}{1000} = 0.90 \quad \left. \begin{array}{l} \text{but it} \\ \text{doesn't make} \\ \text{any sense as} \\ \text{it does not predict as 1 class.} \end{array} \right\}$$

for this we have precision or other classification metrics

(3) Precision $\Rightarrow \frac{TP}{TP+FP}$ (out of all actual values how many were correctly predicted)

\hookrightarrow for class 1, out of all predicted 1, How many are actually one.

(3) Recall $= \frac{TP}{TP+FN}$ (out of all predicted value How many are correctly predicted)

Note \rightarrow Both Precision & Recall we have seen for Positive class (1). Similarly we can see for class (0). but we generally interested in knowing 1st positive class.

Note \rightarrow when FN is imp \rightarrow Recall
when FP is imp \rightarrow Precision

PAGE No.	
DATE	

When to use Precision / Recall.

Clear cases \rightarrow ① spam classification \rightarrow Spam / ham?

\hookrightarrow Here what is imp $\boxed{FP / FN} / TP + TN$

Ans \Rightarrow FP \rightarrow because lets say i worked

a job and selected at one company, they sent me offer letter but it went to spam. and i haven't checked spam folder. so here we give

Precision to FP

$$\text{Precision} = \frac{\underline{TP}}{\underline{TP} + \underline{FP}} \quad | \quad \text{Recall} = \frac{\underline{TP}}{\underline{TP} + \underline{FN}}$$

② \rightarrow where FN is important

\rightarrow Person is diabetic or not Diabetic / Not Diabetic

\bullet in this case when you visit doctor, he says you are not diabetic but you are, this can be harmful for you.

\bullet and if you are non diabetic and he says you are, you might focus more on your health & diet.

Then we can say in this case, False Negative is important than Recall is used.

$$\text{Recall} = \frac{\underline{TP}}{\underline{TP} + \underline{FN}}$$
$$=$$

⑤ F- β Score

$$\frac{(1 + \beta)^2}{\beta} \cdot \frac{\text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$$

when False Positive & False Negative both are imp

Value of β for diff cases

① when both FP & FN are imp

$$\beta = 1, F_{\beta=1} = \frac{\text{P} \times \text{R}}{\text{P} + \text{R}}$$

ϕ

② when FN is imp

$$\beta = 0.5, F_{\beta=0.5} = \frac{1 + 0.25 \times \text{P} \times \text{R}}{\text{P} + \text{R}}$$

③ when ~~FP is imp~~ FP is imp

$$\beta = 2, F_{\beta=2} = \frac{1}{\text{P}}$$

⑥ True Positive Rate

out of all \mathcal{I} , it is actually predicting \mathcal{I}

$$TPR = \frac{TP}{TP + FN} \quad \left. \begin{array}{l} \\ \text{Also called Recall} \end{array} \right.$$

⑦ False Positive Rate

when it is \mathcal{D} how often it is predicting \mathcal{I}

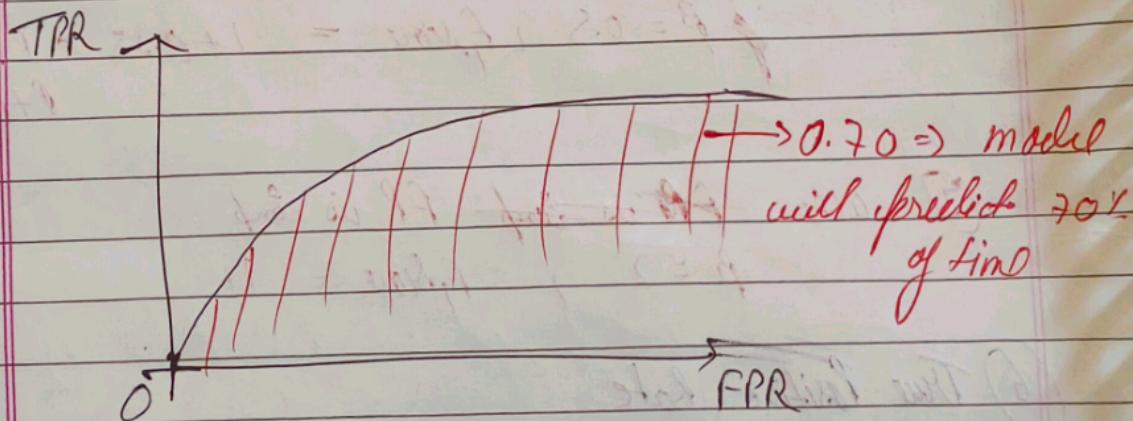
$$FPR = \frac{FP}{TN + FP}$$

8 True Negative Rate \rightarrow out of all actual 0, it is actually predicted 0

$$TNR = \frac{TN}{FP + TN}$$

$$TNR = 1 - FPR$$

9 AUC - ROC \rightarrow Receiver operating Characteristic - Area under curve.



Higher the AUC \rightarrow Better our model is

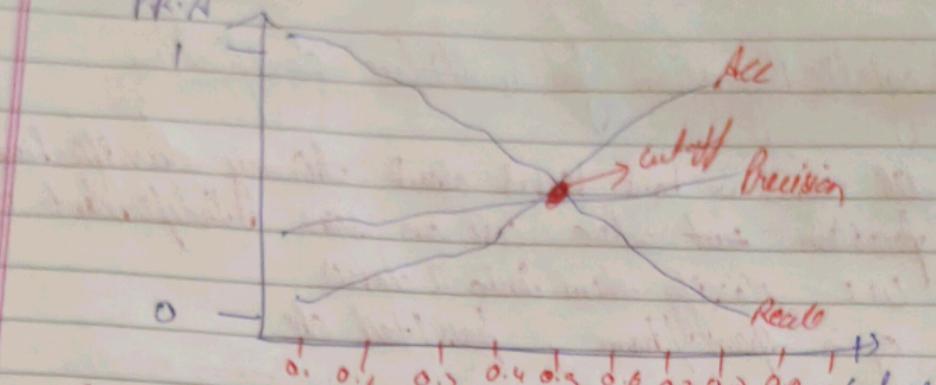
10 \rightarrow Precision - Recall - Accuracy trade off

\hookrightarrow when all three are imp. (P, R, Acc)

\hookrightarrow How to decide ~~cutoff~~ to classify b/w classes?
 \hookrightarrow P-R-C trade off is used for that.

PR.A

PAGE NO.
DATE



cutoff is decided as that point cuts where all (PR,A) three intersect each other.

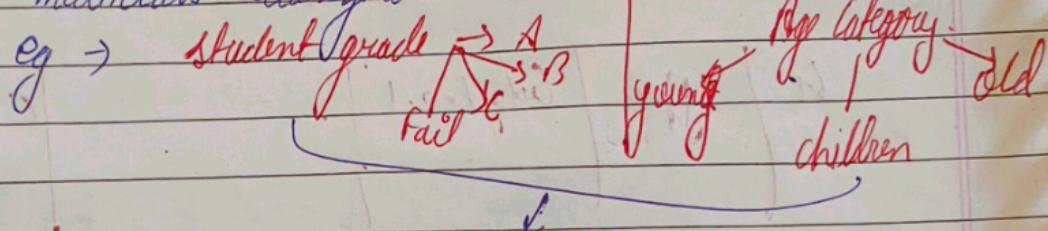
$$C_i / Y_{true} \left[Cutoff - 0.1 \right]_+ - [0.9]$$

for all cut off we will calculate (P, R, F) , then we plot them all. wherever they intersect this will be our point.

Logistic Regression for multiclass classification

till now we have learnt about binary classification problem but where base on some independent variables we classify them into 2 classes.
eg → spam/ham, pass/fail etc.

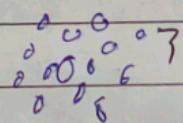
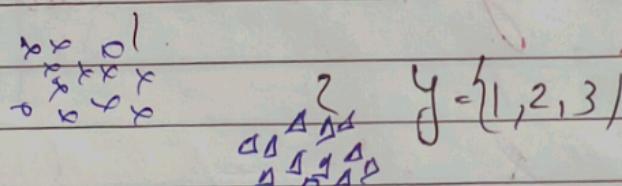
but when we have more than 2 classes it is called as multiclass classification.



for them we have different algo.
 $Sol^n \Rightarrow OVR$ (One vs Rest)

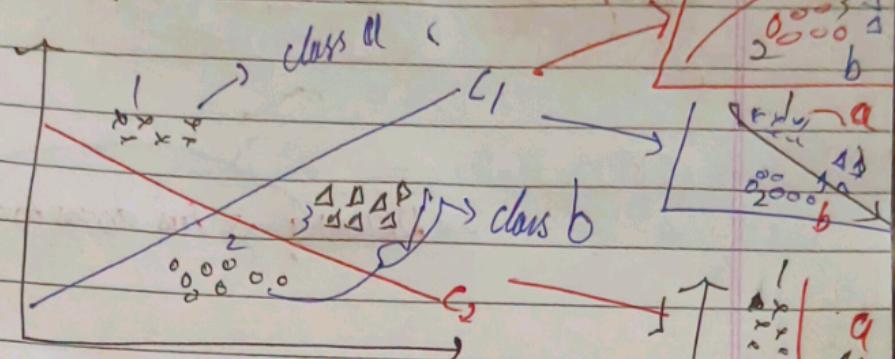
Current Scenario

Multinomial (optional)



- ① One vs Rest → we will try to modify the logistic in such a way that we were able to solve multi class problem.

How to solve it?



- Drew a line 1 and said $1 \rightarrow \text{class } a$
 $(2, 3) \rightarrow \text{class } b$
- Line 2 $\rightarrow 2 \rightarrow \text{class } a$
 $(1, 3) \rightarrow \text{class } b$

$$\text{Line 3} \rightarrow (1, 2) \rightarrow b \\ (3) \rightarrow a$$

Here No. of logistic regression model = no of classes
 built present

m_1	m_2	m_3
-	-	y_1
-	3	0
-	2	0
-	1	1
-	2	0
-	3	0

$m_1 = \text{class 1 vs rest}$

\Rightarrow Now we have made the model $m_2 = \text{class 2 vs rest}$
 but How will the prediction
 made? $m_3 = \text{class 3 vs rest}$

example →

Viral Kohli → has access to

Batting Coach

Bowling Coach

Fielding Coach

Batting Coach

Bowling Coach

Fielding Coach

Bhuvneshwar Kumar → has access to

Batting

Bowling

Fielding

Jadeja → has access to

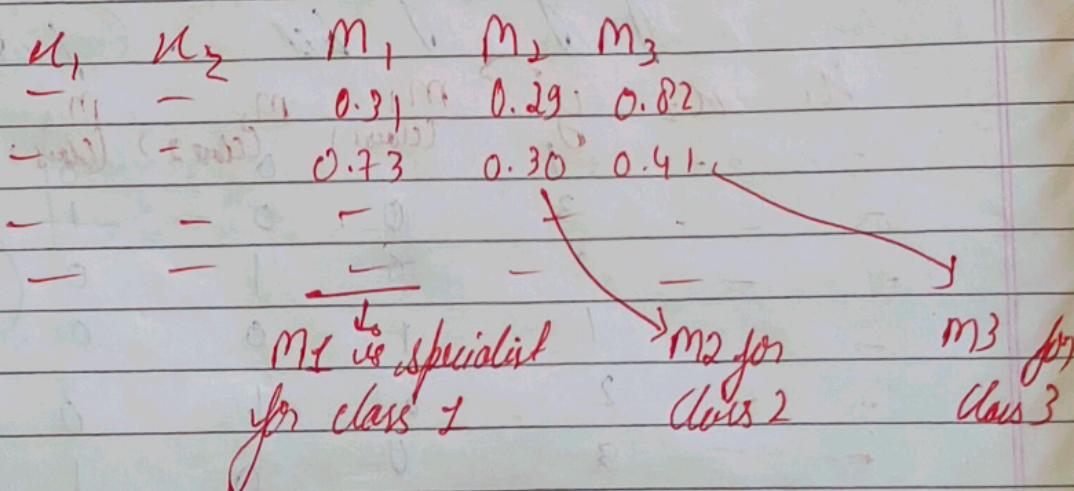
Batting

Bowling

Fielding

Here all players goes to that coach who has highest knowledge about their skills.

Similarly in One vs Rest algo we have n no of models where n is the no. of classes



Now lets say for $u_1 = 3$, $u_2 = 2$ we have 3 Model with diff. prob.

u_1	u_2	M_1	M_2	M_3
-------	-------	-------	-------	-------

3	7	0.31	0.27	0.87
---	---	------	------	------

M_3 is 27%

M_3 is 87%

M_1 is 37% confident

u_1, u_2 belongs to class 1

Confident

Confident

now for x_i, y_i we have m_j is highest confident that
it belongs to class j . So we will
choose it.

Note → Final prediction is that class / model which
is higher confident about any data point.

Draw back → we were creating n models for
1 problem making it computationally
expensive as binary classification is
trained for each class

Step 1 → no of model = no of class

Step 2 → make model for each and get the
probability for each model

Step 3 → Attach to that class which has
highest probability.

→ This problem is solved using a
multinomial method / softmax regression

- Multinomial method / softmax regression
→ we don't decompose the problem into binary
classification instead we modify cost function

→ Single model to reduce cost

$$\text{Sigmoid } (8) = \frac{1}{1 + e^{-z}} \quad \left| \begin{array}{l} \text{Softmax } \sigma(z_j) = \frac{e^{z_j}}{\sum_{j=1}^n e^{z_j}} \\ \text{fun} \end{array} \right. \quad j - \text{no of class}$$

$$Z = (1, 2, 3)$$

$$\sigma(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}, \quad g(z_1) = \frac{e^{z_1}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$g(z_3) = \frac{e^{z_3}}{e^{z_1} + e^{z_2} + e^{z_3}}$$

$$\text{Cost fun} = -\frac{1}{n} \sum_{i=1}^n y_i \log(\bar{y}_i) + (1-y_i) \log(1-\bar{y}_i)$$

→ after modification for
the softmax

$$\text{Cost fun} = -\frac{1}{n} \sum_{i=1}^n \sum_{k=1}^K y_k^{(i)} \log \bar{y}_k^{(i)}$$

n = nof Data point
 K = no of class.

features (x_1, x_2, y)			$y_{k=1}$	$y_{k=2}$	$y_{k=3}$
x_1	x_2	y	1	1	0
1	1	1	0	0	1
2	2	2	0	1	0

$$\begin{aligned}
 & \cancel{y_{k=1}} = \\
 & y_{1 \rightarrow \text{now}} \cdot \log \bar{y}_1 + y_{2 \rightarrow \text{now}} \cdot \log \bar{y}_2 + y_{3 \rightarrow \text{now}} \cdot \log \bar{y}_3 \\
 & y_{1 \rightarrow \text{class}} \cdot \cancel{\log \bar{y}_1} + y_{2 \rightarrow \text{class}} \cdot \cancel{\log \bar{y}_2} + y_{3 \rightarrow \text{class}} \cdot \cancel{\log \bar{y}_3} \\
 & y_{1 \rightarrow \text{now}} \cdot \cancel{\log \bar{y}_1} + y_{2 \rightarrow \text{now}} \cdot \cancel{\log \bar{y}_2} + y_{3 \rightarrow \text{now}} \cdot \cancel{\log \bar{y}_3} \Rightarrow 1
 \end{aligned}$$

now cost fun becomes $C_F = y_1^{(1)} \log \hat{y}_1^{(1)} + y_2^{(2)} \log \hat{y}_2^{(2)} + y_3^{(3)} \log \hat{y}_3^{(3)}$

$$\text{now we have } \hat{y}_1^{(1)} + \hat{y}_2^{(2)} + \hat{y}_3^{(3)}$$

$$\begin{aligned} \hat{y}_1^{(1)} &= \sigma(\theta_0^{(1)} + \theta_1^{(1)} u_{11} + \theta_2^{(1)} u_{12}) \\ \hat{y}_2^{(2)} &= \sigma(\theta_0^{(2)} + \theta_1^{(2)} (u_{21}) + \theta_2^{(2)} (u_{22})) \\ \hat{y}_3^{(3)} &= \sigma(\theta_0^{(3)} + \theta_1^{(3)} (u_{31}) + \theta_2^{(3)} (u_{32})) \end{aligned}$$

gradient descent

$$\frac{\partial L}{\partial \theta_0^{(1)}}, \frac{\partial L}{\partial \theta_1^{(2)}} \dots \dots \rightarrow \text{diff.}$$

$$\text{convergence } \theta_i^{(1)} = \theta_i^{(1)} - \eta \cdot \frac{\partial L}{\partial \theta_i^{(1)}}$$

Simultaneously change θ

now cost fun becomes $CF = \hat{y}_1^{(1)} \log \hat{y}_1^{(1)} + \hat{y}_2^{(2)} \log \hat{y}_2^{(2)} + \hat{y}_3^{(3)} \log \hat{y}_3^{(3)}$

$$\text{now we have } \hat{y}_1^{(1)} + \hat{y}_2^{(2)} + \hat{y}_3^{(3)}$$

$$\begin{array}{c} \alpha_0^{(1)} \quad \alpha_1^{(1)} \quad \alpha_2^{(1)} \\ \alpha_0^{(2)} \quad \alpha_1^{(2)} \quad \alpha_2^{(2)} \\ \alpha_0^{(3)} \quad \alpha_1^{(3)} \quad \alpha_2^{(3)} \end{array} \quad \begin{array}{l} \hat{y}_1^{(1)} = \sigma(\alpha_0^{(1)} + \alpha_1^{(1)} u_{11} + \alpha_2^{(1)} u_{12}) \\ \hat{y}_2^{(2)} = \sigma(\alpha_0^{(2)} + \alpha_1^{(2)} (u_{21}) + \alpha_2^{(2)} (u_{22})) \\ \hat{y}_3^{(3)} = \sigma(\alpha_0^{(3)} + \alpha_1^{(3)} (u_{31}) + \alpha_2^{(3)} (u_{32})) \end{array}$$

gradient descent

$$\frac{\partial L}{\partial \alpha_0^{(1)}} , \frac{\partial L}{\partial \alpha_0^{(2)}} \dots \dots \dots \text{9 diff.}$$

$$\text{Convergence } \alpha_1^{(1)} = \alpha_1^{(1)} - \eta \frac{\partial L}{\partial \alpha_1^{(1)}}$$

Simultaneously change α