

Data Science and Artificial Intelligence

Machine Learning



Classification

Lecture No. 4

By- SIDDHARTH SABHARWAL SIR

GATE WALLAH

Recap of Previous Lecture



Topic

logistic Reg.

Topic

Questions

Topic

logit / 1D Question

Topic

Topic

Topics to be Covered



Topic

^{LR}
Likelihood, Logistic Reg.

Topic

Cost fxn, loss function


Topic

Multi class LR

Topic

Logistic Reg: how safe from outliers

Topic

The background of the slide features a person with long dark hair, seen from behind, wearing a brown robe. They are looking out over a body of water towards a sunset with orange and yellow clouds. The text is overlaid on this image.

**Inspiration comes from
within yourself. One
has to be positive.
When you're positive,
good things happen.**

DEEP ROY

Let $x = 1$ if an email subject includes the student's name and $x = 0$ otherwise.

There are 350 emails with $x = 1$ of which 161 were opened ($y = 1$), and 400 emails with $x = 0$ of which 140 were opened.

Fit a logistic regression for the log-odds of opening:

$y=1$ email opened
 $y=0$ email not opened

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 x.$$

Which model is correct?

- A) $-0.6190 + 0.4587x$
- B) $-0.1603 + 0.6190x$
- C) $-0.6190 + 0.1603x$
- D) $0.4587 - 0.6190x$

Answer: A

-H.W

+H.W

@ $x=0$ $\beta_0 = \log_e \frac{P(y=1|x=0)}{P(y=0|x=0)}$

$\Rightarrow 140/400$

$\Rightarrow 260/400$

@ $x=1$

$\beta_1 + \beta_0 = \log_e \frac{P(y=1|x=1)}{P(y=0|x=1)}$

$\Rightarrow 161/350$

$\Rightarrow 189/350$



What is Likelihood.

Example 1: Suppose that X is a discrete random variable with the following probability mass function: where $0 \leq \theta \leq 1$ is a parameter. The following 10 independent observations

p.w.

| X | 0 | 1 | 2 | 3 |
|--------|-------------|------------|-----------------|----------------|
| $P(X)$ | $2\theta/3$ | $\theta/3$ | $2(1-\theta)/3$ | $(1-\theta)/3$ |

were taken from such a distribution: (3,0,2,1,3,2,1,0,2,1). What is the maximum likelihood estimate of θ .

You are walking down Shattuck Ave. when you find a quarter on the ground. You see nothing unusual about this quarter, so you figure it is almost certainly a fair coin, though you realize that manufacturing irregularities in the coin minting process mean that coins are rarely *exactly* fair. You toss the coin 10 times and observe the following outcomes:

P.W

H H H H H H H H T

with H denoting heads and T denoting tails. Assume coin tosses are independent. What is the maximum likelihood estimate of the next toss being heads?

- ☐ $\frac{5}{10}$
- ☐ between $\frac{5}{10}$ and $\frac{9}{10}$
- ☐ $\frac{9}{10}$
- ☐ more than $\frac{9}{10}$

There are 5 balls in a bag. Each ball is either red or blue. Let θ (an integer) be the number of blue balls. We want to estimate θ , so we draw 4 balls **with replacement** out of the bag, replacing each one before drawing the next. We get “blue,” “red,” “blue,” and “blue” (in that order).

- (a) [5 pts] Assuming θ is fixed, what is the likelihood of getting exactly that sequence of colors (expressed as a function of θ)?

h.p.w.



Logistic Regression

What is the meaning of likelihood

We need β values such that

Probab of class '1' for datapoint 1 \Rightarrow close to 1

u u u '1' u u u

2 \Rightarrow Close to 1

u u u '0' u u u

3 \Rightarrow Close to 1

u u u '0' u u u

4 \Rightarrow Close to 1

\rightarrow If this happen it means even on training data i is v. Small

2D

Training data

| x^1 | x^2 | y |
|-------|-------|-----|
| a | b | 1 |
| c | d | 1 |
| e | f | 0 |
| g | h | 0 |
| i | j | 0 |



Logistic Regression

What is the meaning of likelihood

- we know that class '1' Probab
for any data point $P_i = \frac{1}{1+e^{-x_i\beta}}$
Probab of class '0' $\Rightarrow 1 - P_i$
 $\Rightarrow \left(1 - \frac{1}{1+e^{-x_i\beta}}\right)$

2D Training data

| x^1 | x^2 | y |
|-------|-------|-----|
| a | b | 1 |
| c | d | 1 |
| e | f | 0 |
| g | h | 0 |
| i | j | 0 |



Logistic Regression

What is the meaning of likelihood

So we want to maximize the
Probability that Y matrix =

$$\max \text{Probab of } Y \text{ matrix} = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \end{bmatrix} \begin{matrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \end{matrix}$$

2D Training data

| | | |
|-------|-------|-----|
| x^1 | x^2 | y |
| a | b | 1 |
| c | d | 1 |
| e | f | 0 |
| g | h | 0 |
| i | j | 0 |

What is the meaning of likelihood

So we want to maximize the
Probability that Y matrix =

$$\begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix} \begin{matrix} y_1 \\ y_2 \\ y_3 \\ y_4 \\ y_5 \\ \vdots \end{matrix}$$

• $\max P(\text{probab of } Y \text{ matrix}) = \begin{bmatrix} 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ \vdots \end{bmatrix}$

Now $\max P(y_1=1, y_2=1, y_3=0, y_4=0, y_5=0 \dots)$ all y 's are independent

$$\max P(y_1=1) \cdot P(y_2=1) \cdot P(y_3=0) \cdot P(y_4=0) \cdot P(y_5=0) \dots$$

$$\max P_{11} \cdot P_{21} \cdot [1 - P_{31}] [1 - P_{41}] [1 - P_{51}] \dots$$

2D Training data

| x^1 | x^2 | y |
|-------|-------|-----|
| 1 | a | 1 |
| | c | 1 |
| | e | 0 |
| | g | 0 |
| | i | 0 |
| | j | 0 |
| | | 0 |

What is the meaning of likelihood

2D Training data

Now $\max P(y_1=1, y_2=1, y_3=0, y_4=0, y_5=0 \dots)$

$\max P(y_1=1) \cdot P(y_2=1) \cdot P(y_3=0) \cdot P(y_4=0) P(y_5=0) \dots$

$\max P_{11} \cdot P_{21} \cdot [1 - P_{31}] [1 - P_{41}] [1 - P_{51}] \dots$

1

| x^1 | x^2 | y |
|-------|-------|-----|
| a | b | 1 |
| c | d | 1 |
| e | f | 0 |
| g | h | 0 |
| i | j | 0 |

for any datapoint \Rightarrow $\begin{cases} P_{i1} & \text{if } y_i=1 \\ (1 - P_{i1}) & \text{if } y_i=0 \end{cases}$

all y's are independent

\rightarrow for any i th datapoint $(P_{i1})^{y_i} (1 - P_{i1})^{(1-y_i)}$

for any datapoint i^{th} datapoint $\Rightarrow \begin{cases} p_{i1} & \text{if } y_i = 1 \\ (1-p_{i1}) & \text{if } y_i = 0 \end{cases}$

\rightarrow for any i^{th} datapoint $(p_{i1})^{y_i} (1-p_{i1})^{(1-y_i)}$

if $y_i = 1$

$$(p_{i1})^1 (1-p_{i1})^{1-1}$$

p_{i1}

$y_i = 0$

$$(p_{i1})^0 (1-p_{i1})^{1-0}$$

$(1-p_{i1})$

- Likelihood meaning % Think you are an artist

In logistic
regression similar

task, for every data
point we want $\hat{y} = y$

We want Predicted class of
every data = y

Your task is to draw painting of
of the person exactly similar

• We want Keponi ki
Pooni ^{same} matrix generate
Kanne ki Probab max
Karo

What is the meaning of likelihood

2D Training data

Now $\max P_Y$
 $\max P(y_1=1, y_2=1, y_3=0, y_4=0, y_5=0 \dots)$

$\max P(y_1=1) \cdot P(y_2=1) \cdot P(y_3=0) \cdot P(y_4=0) \cdot P(y_5=0) \dots$

$\max P(y_1=1 | x_1; \beta) P(y_2=1 | x_2; \beta) P(y_3=0 | x_3; \beta) \dots$

| x^1 | x^2 | y |
|-------|-------|-----|
| a | b | 1 |
| c | d | 1 |
| e | f | 0 |
| g | h | 0 |
| i | j | 0 |

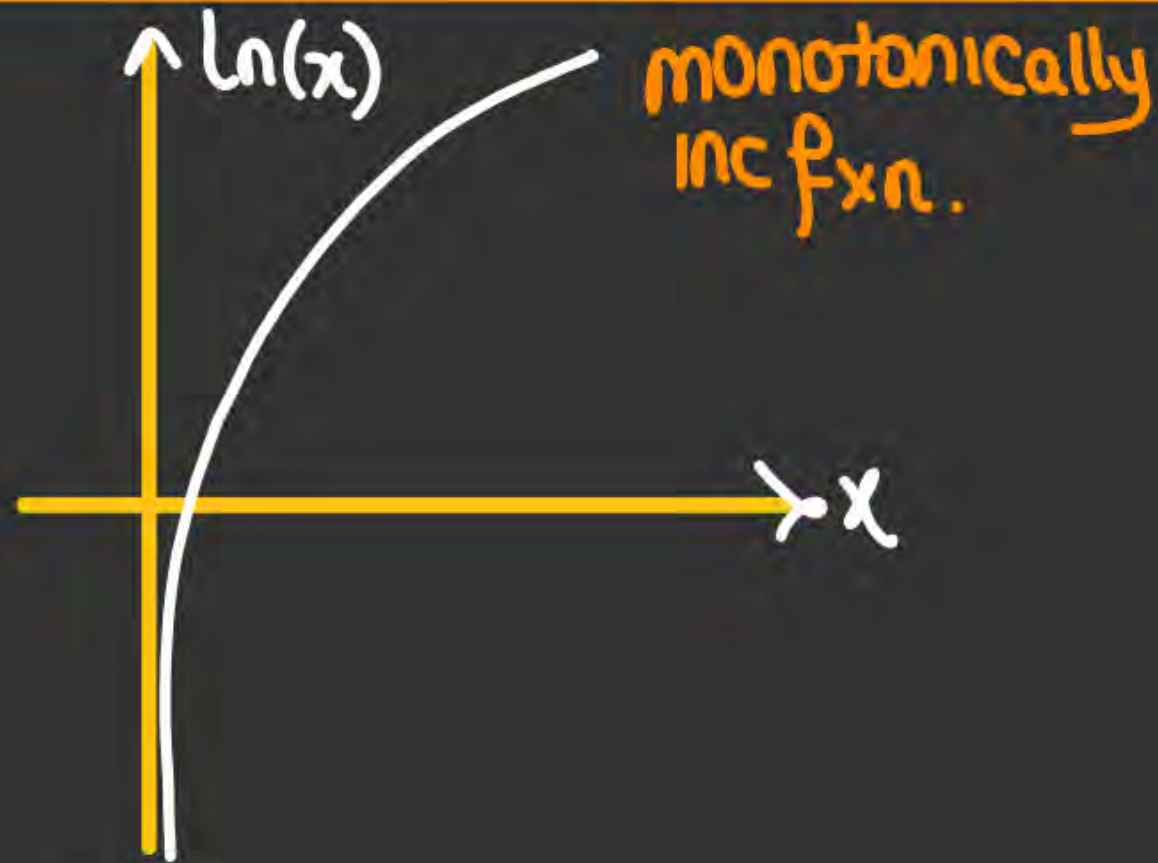
$(P_{11})^{y_1} (1-P_{11})^{1-y_1} (P_{21})^{y_2} (1-P_{21})^{1-y_2} \dots$ all y's are independent

P_{11} = datapoint 1
 Class 1 probab

$\left\{ \max \prod_{i=1}^N (P_{i1})^{y_i} (1-P_{i1})^{1-y_i} \right\}$

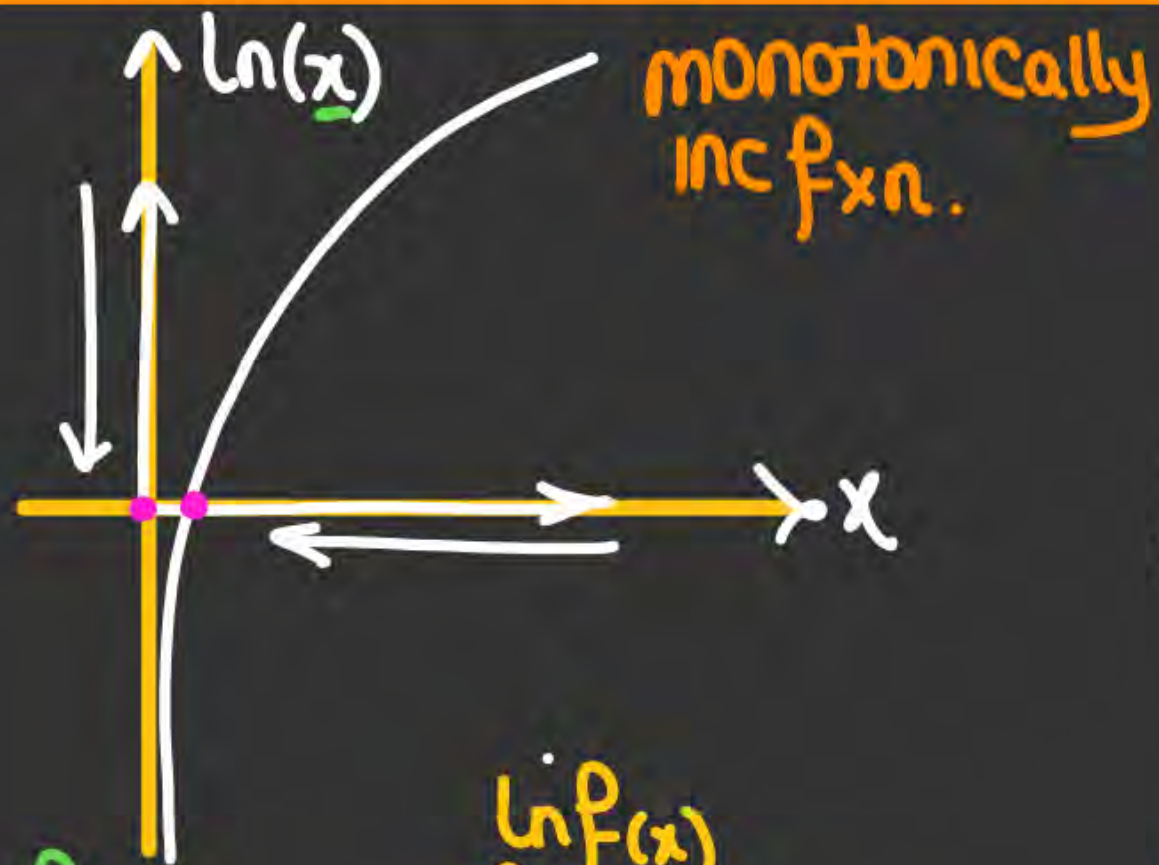
max likelihood \Rightarrow .

- $$\max \prod_{i=1}^N (p_{i1})^{y_i} (1-p_{i1})^{(1-y_i)}$$



if we want to solve $\max f(x)$

- we want to find value of x where $f(x)$ is max

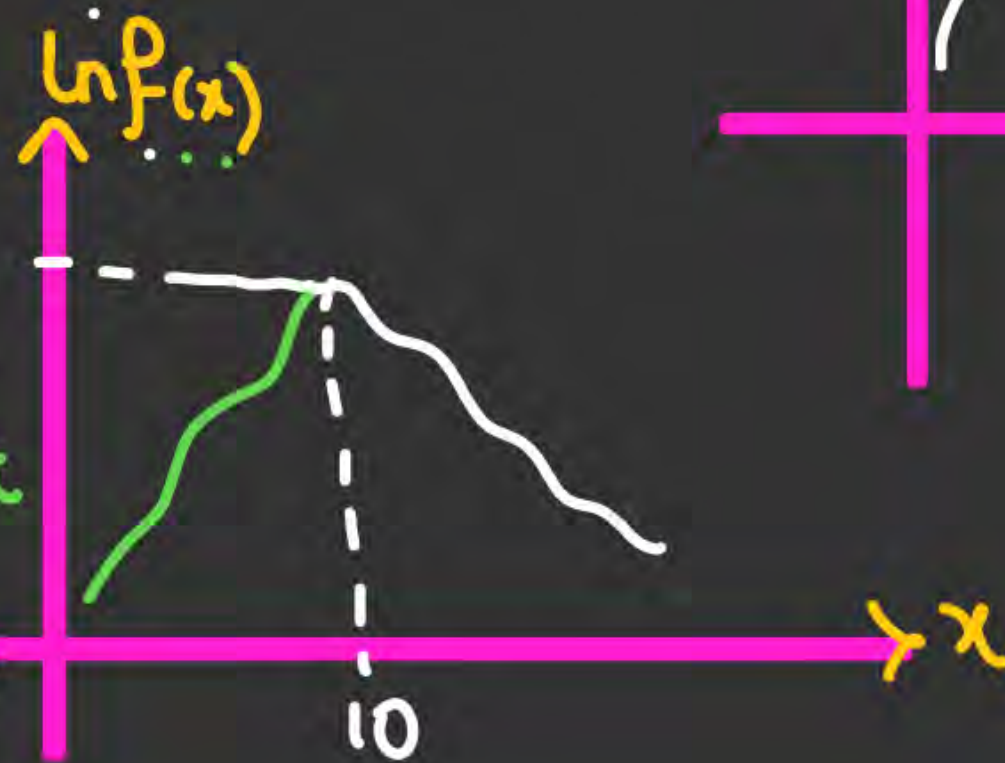


if we want to solve $\max f(x)$

- we want to find value of x where $f(x)$ is max



- log is inc fcn
- Jab $f(x)$ inc karega, $\log f(x)$ will also inc
- when $f(x)$ dec then $\log f(x)$ will dec



max likelihood \Rightarrow

*
$$\max \prod_{i=1}^N (p_{i1})^{y_i} (1-p_{i1})^{(1-y_i)}$$

✓
$$\ln(abcd) \Downarrow \ln a + \ln b + \ln c + \ln d$$

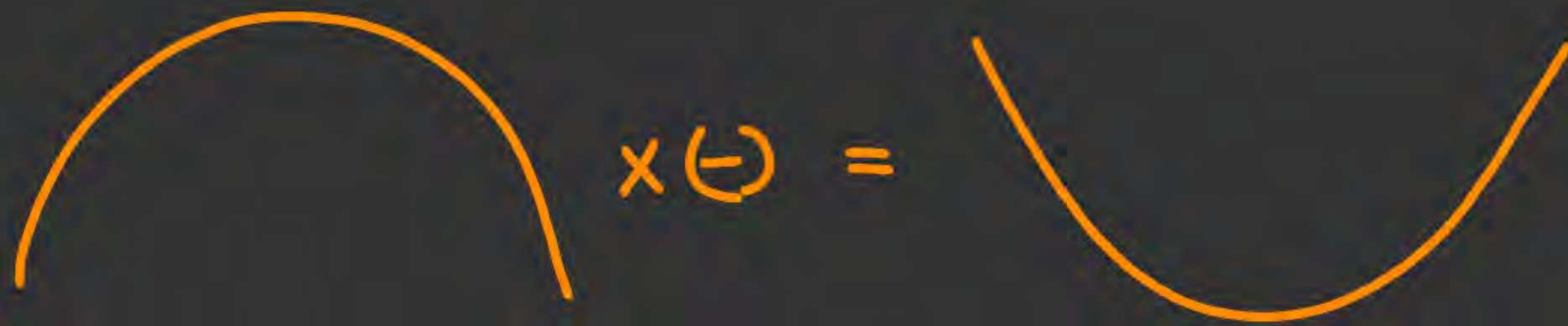
Same as $\max \ln \prod_{i=1}^N (p_{i1})^{y_i} (1-p_{i1})^{1-y_i}$

* \max
Log Likelihood \Rightarrow

$$\max \sum_{i=1}^N y_i \ln(p_{i1}) + (1-y_i) \ln(1-p_{i1})$$

Concave Curve

aise β 's
Chahiye
that max
this equation.



So loss function \Rightarrow - log likelihood.

Convex
Curve

Cross Entropy = $-\left[\sum_{i=1}^N y_i \ln p_{i1} + (1-y_i) \ln (1-p_{i1}) \right]$

• Average Cross Entropy $\Rightarrow -\frac{1}{N} \sum_{i=1}^N y_i \ln p_{i1} + (1-y_i) \ln (1-p_{i1})$



Linear Classification



Logistic Regression

- The Cost function

done

→ log likelihood

How can we
use log into
this function



Linear Classification



Logistic Regression

- **The Cost function**



Linear Classification



Logistic Regression

- Log likelihood and cross entropy loss function.

done

- Why linear classification is more effected by outlier
- $\max \sum y_i x_i \beta$
- * Outlier has large $x_i \beta$
- * $y_i x_i \beta$ will be large -ve
hugely effect model

- logistic negnention is less effected by the outlier.

CE

$y_i = 1, 0$

$$\left[\sum_{i=1}^N y_i \log p_{i1} + (1-y_i) \log (1-p_{i1}) \right]$$

- logistic negnemion is len effected by the outlier.

CE

$y_i = 1, 0$

$$\left[\sum_{i=1}^N y_i \log p_{i1} + (1-y_i) \log (1-p_{i1}) \right]$$



outlier
Class 0

$x_i \beta$ for outlier is +ve & large

$$p_{i1} = \frac{1}{1 + e^{-x_i \beta}}$$

$$p_{i1} = 0.899$$

Class 0

Class 1



Linear Classification



Logistic Regression

- **Multiclass Logistic Regression – Softmax Regression**

Prediction Rule

$$\text{Predicted class} = \underset{k}{\operatorname{argmax}} (\mathbf{w}_k \cdot \mathbf{x})$$



Linear Classification



Logistic Regression

- Extending the case for more than 2 classes... (not imp)

We trained a three-way logistic regression and obtained weights

$$w_a = (1, 1, 0), w_b = (-1, 1, 1), w_c = (2, 1, 2).$$

What label would be given to the point $x = (0, 1, 1)$?

- (A) A.
- (B) B.
- (C) C.



Linear Classification



Logistic Regression

- How to turn the value into probability

(2 points) What is the definition of $\text{softmax}(x_1, \dots, x_n)$? Recall that this function takes in a list of n real numbers x_1, \dots, x_n and outputs a list of n real numbers.

- A. $\text{softmax}(x_1, \dots, x_n) = \left[\frac{x_1}{\sum_{i=1}^n x_i}, \dots, \frac{x_n}{\sum_{i=1}^n x_i} \right]$
- B. $\text{softmax}(x_1, \dots, x_n) = \left[\frac{e^{x_1}}{\sum_{i=1}^n e^{x_i}}, \dots, \frac{e^{x_n}}{\sum_{i=1}^n e^{x_i}} \right]$
- C. $\text{softmax}(x_1, \dots, x_n) = \left[\frac{e^{x_1}}{e^{\sum_{i=1}^n x_i}}, \dots, \frac{e^{x_n}}{e^{\sum_{i=1}^n x_i}} \right]$
- D. $\text{softmax}(x_1, \dots, x_n) = [e^{x_1}, \dots, e^{x_n}]$

Q.



- i. [2 Pts] Suppose our true labels are $\vec{y} = [0, 0, 1]$, our predicted probabilities of being in class 1 are $[0.1, 0.6, 0.9]$, and our threshold is $T = 0.5$. Give the total (not average) cross-entropy loss. Do not simplify your answer.

Total CE Loss = $\begin{matrix} 0, 0, 1 \\ y_1=0, y_2=0, y_3=1 \end{matrix} \left\{ \begin{matrix} 0.1, 0.6, 0.9 \\ \text{---} \end{matrix} \right\}$ $CE = - \sum_{i=1}^3 y_i \log p_i + (1-y_i) \log(1-p_i)$

$$= - \left[\left(\cancel{y_1 \log p_1} + (1-y_1) \log(1-p_1) \right) + \left(\cancel{y_2 \log p_2} + (1-y_2) \log(1-p_2) \right) + \left(y_3 \log p_3 + \cancel{(1-y_3) \log(1-p_3)} \right) \right]$$

$$= - \left[1 \log(1-0.1) + 1 \log(1-0.6) + 1 \log 0.9 \right] \Rightarrow 1.127$$

Q.

- i. [2 Pts] Suppose our true labels are $\vec{y} = [0, 0, 1]$, our predicted probabilities of being in class 1 are $[0.1, 0.6, 0.9]$, and our threshold is $T = 0.5$. Give the total (not average) cross-entropy loss. Do not simplify your answer.

Total CE Loss = $\begin{matrix} 0, 0, 1 \\ y_1=0, y_2=0, y_3=1 \end{matrix} \left\{ \begin{matrix} 0.1, 0.6, 0.9 \end{matrix} \right\}$ $CE = - \sum_{i=1}^3 y_i \log p_i + (1-y_i) \log(1-p_i)$

- ii. [2 Pts] For the same values as above, give the total squared loss. Do not simplify your answer.

Squared Loss = $RSS = \sum_{i=1}^3 (y_i - \hat{y}_i)^2$
 $= (0-0)^2 + (0-1)^2 + (1-1)^2 = 1$

| y | p_1 | \hat{y} |
|-----|-------|-----------|
| 0 | 0.1 | 0 |
| 0 | 0.6 | 1 |
| 1 | 0.9 | 1 |

$\hat{y}=1 \quad p_1 > 0.5$
 $\hat{y}=0 \quad p_1 < 0.5$

(b) (2.0 pt) Consider the following three rows from our training data, along with their predicted probabilities \hat{y} for some choice of θ :

| x^1 | x^2 | y | P_1 |
|-------|-------|-----|--------------------------------|
| hue | abv | | $\hat{y} = \sigma(x^T \theta)$ |
| -0.17 | 0.24 | 0 | 0.45 ✓ |
| -1.18 | 1.61 | 0 | 0.19 ✓ |
| 1.25 | -0.97 | 1 | 0.80 ✓ |

$\sigma(x^T \theta)$

What is the mean cross-entropy loss on just the above three rows of our training data?

☐ $-\frac{1}{3}(\log(0.45) + \log(0.19) + \log(0.20))$

☐ $-\frac{1}{3}(\log(0.55) + \log(0.19) + \log(0.80))$

☐ $-\frac{1}{3}(\log(0.45) + \log(0.81) + \log(0.80))$

☒ $-\frac{1}{3}(\log(0.55) + \log(0.81) + \log(0.80))$

$$-\frac{1}{3} \left[(1-0) \log_e(1-P_1) + (1-0) \log_e(1-P_2) + 1 \log_e P_3 \right]$$

(1 pt) In this question, assume that we are using the logistic regression model $\hat{y} = \sigma(x^T \theta)$.

Suppose we want to modify cross-entropy loss to penalize predictions for observations that are truly positive twice as much as we penalize predictions for observations that are truly negative. Which of the following loss functions could we use? Recall that the average cross-entropy loss is:

CE

$$R(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$$

• ~~#~~ positive class
Points KO
2 times
Penalize

- ☐ $R(\theta) = -\frac{2}{n} \sum_{i=1}^n (y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$
- ☒ $R(\theta) = -\frac{1}{n} \sum_{i=1}^n (2y_i \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$
- ☐ $R(\theta) = -\frac{1}{n} \sum_{i=1}^n (y_i \log(\hat{y}_i) + 2(1 - y_i) \log(1 - \hat{y}_i))$
- ☐ $R(\theta) = -\frac{1}{n} \sum_{i=1}^n ((y_i + 2) \log(\hat{y}_i) + (1 - y_i) \log(1 - \hat{y}_i))$

$$CE = - \left[\sum_{i=1}^N y_i \log p_{i1} + (1-y_i) \log(1-p_{i1}) \right]$$

if $y_i = 1$
Term $\log p_{i1}$

if $y_i = 0$

Term $\log(1-p_{i1})$

CE is loss we want to minimize it

* CE minimize = 0

if for class ^{Point} $\log p_{i1} = 0$

if for class ^{Point} $\log(1-p_{i1}) = 0$

CE = 0

loss = 0

CE min \Rightarrow for class point $p_{i1} = 1$
" " " " $p_{i1} = 0$

$$CE = - \left[\sum_{i=1}^N 2y_i \log p_{i1} + (1-y_i) \log(1-p_{i1}) \right]$$

Classifier biased toward class 1 points.

Suppose after training our model we get $\vec{\beta} = [-1.2 \quad -0.005 \quad 2.5]^T$, where -1.2 is an intercept term, -0.005 is the parameter corresponding to passenger's age, and 2.5 is the parameter corresponding to sex.

- i. [3 Pts] Consider Sīlānah Iskandar Nāsīf Abī Dāghir Yazbak, a 20 year old female. What chance did she have to survive the sinking of the Titanic according to our model? Give your answer as a probability in terms of σ . If there is not enough information, write "not enough information".

$$\frac{1}{1 + e^{-x_i \beta}} \quad \text{p.w.} \quad \text{p.w.}$$

$$P(Y = 1 | \text{age} = 20, \text{female} = 1) = \boxed{}$$

- ii. [3 Pts] Sīlānah Iskandar Nāsīf Abī Dāghir Yazbak actually survived. What is the cross-entropy loss for our prediction in part i? If there is not enough information, write "not enough information."

Suppose you have a logistic regression model for spam detection, using a dataset with a binary outcome that indicates whether an email is spam (1) or not spam (0). The predictor variables x_1 , x_2 , and x_3 are boolean values (0 or 1) that indicate whether the email contains the words "free", "order", and "homework", respectively. The model has four parameters: weights w_1 , w_2 , w_3 , and offset b .

+P.W

You find that emails containing the words "free" and "order" have a higher probability of being spam, while emails containing the word "homework" have a lower probability of being spam.

Given this information, which of the following signs is most likely for the weights w_1 , w_2 , and w_3 ?

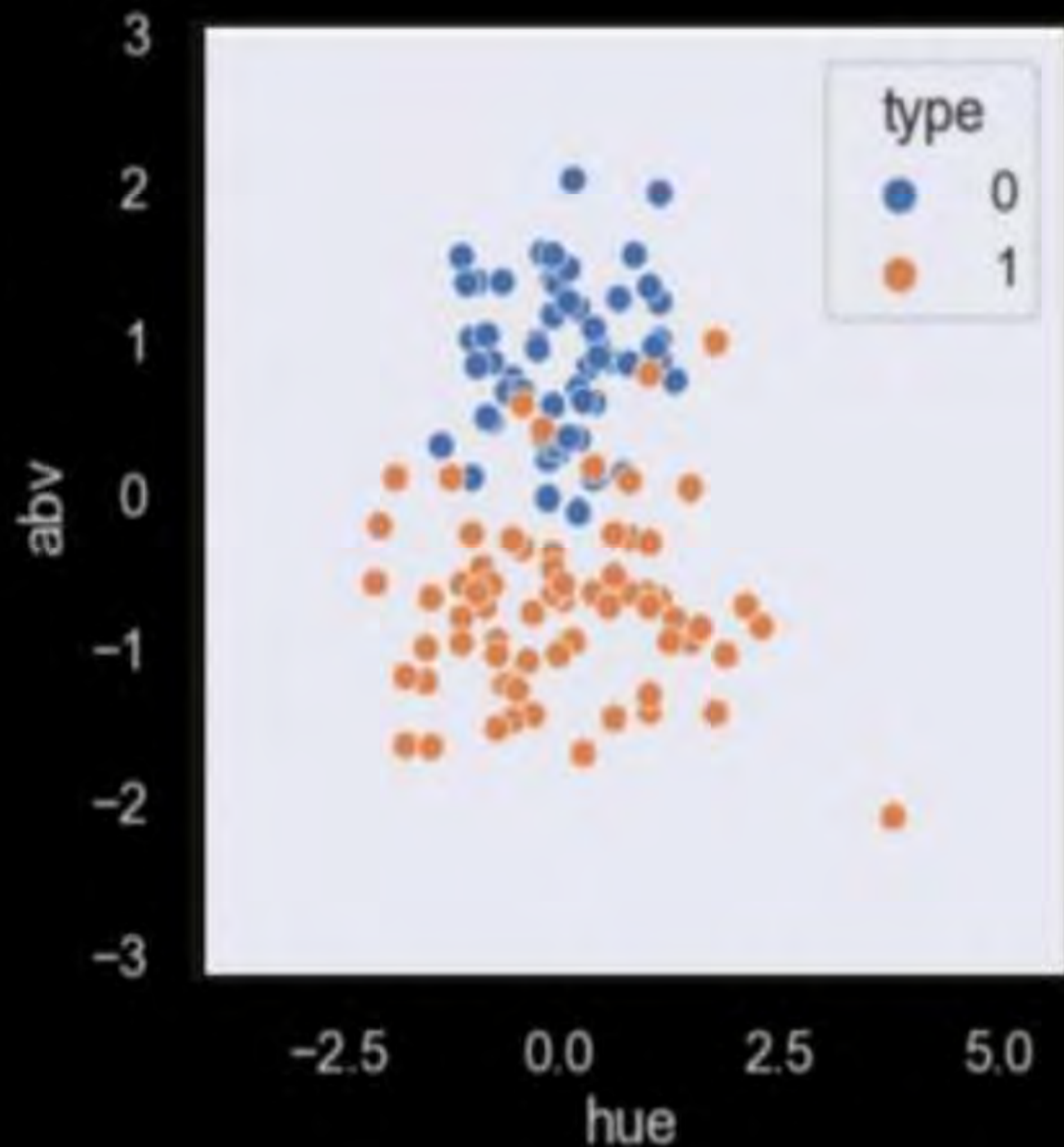
- (A) All positive
- (B) All negative
- (C) w_1 and w_2 are positive, w_3 is negative
- (D) w_1 and w_2 are negative, w_3 is positive

Consider the following scatter plot of our two (standardized) features.

Which of the following statements are true about an unregularized logistic regression model fit on the above data? Select all that apply.

f.w.

- ☐ After performing logistic regression, the weight for the hue feature will very likely have a negative sign.
- ☐ After performing logistic regression, the weight for the abv feature will very likely have a negative sign.

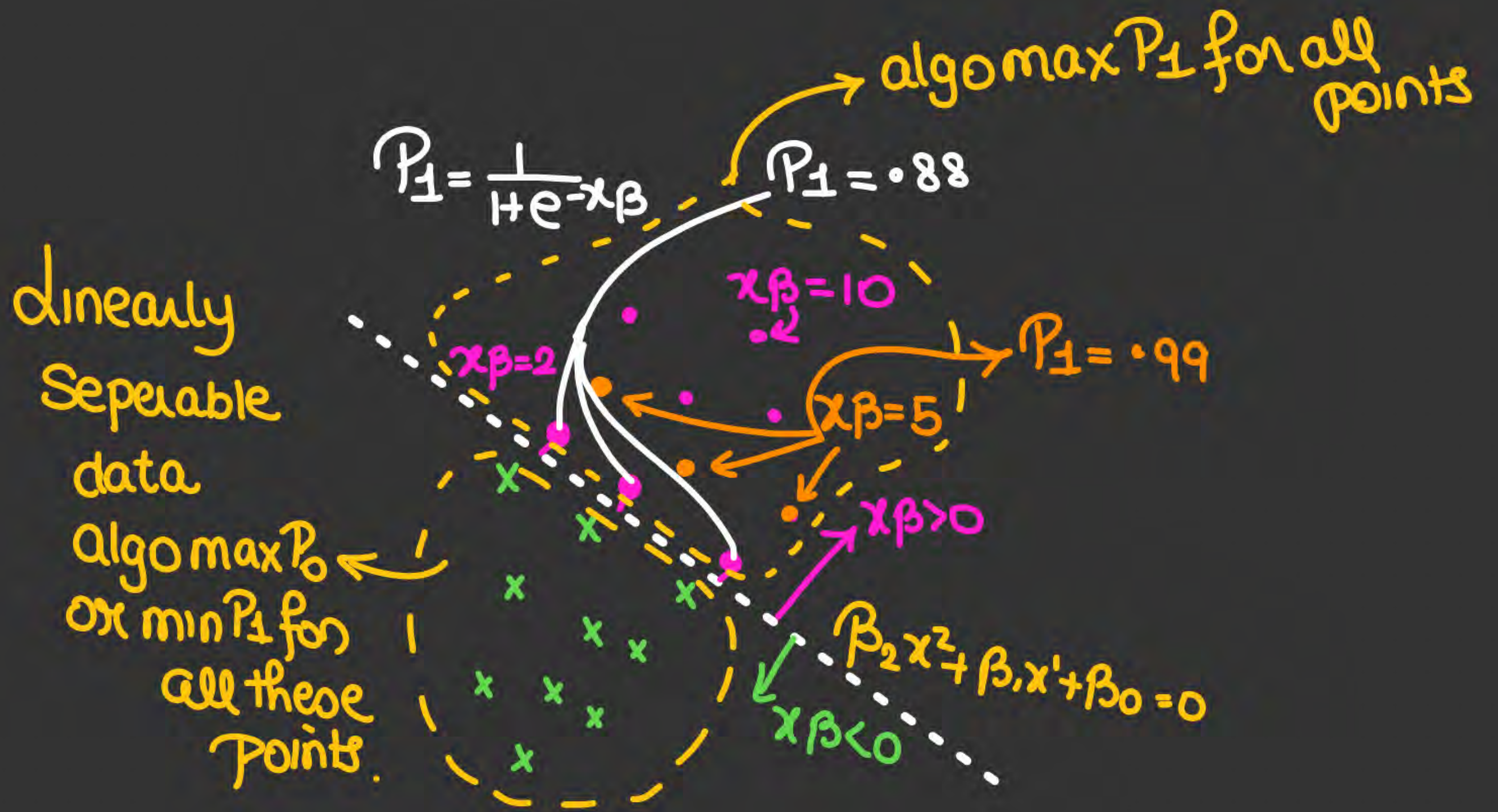




Logistic Regression

- Extending the case for more than 2 classes... (not imp)

$$\begin{aligned} \Rightarrow 4x^2 + 8x' + 10 &= 0 \\ \Rightarrow 8x^2 + 16x' + 20 &= 0 \\ \Rightarrow 400x^2 + 800x' + 10000 &= 0 \end{aligned} \quad \left. \vphantom{\begin{aligned} \Rightarrow 4x^2 + 8x' + 10 &= 0 \\ \Rightarrow 8x^2 + 16x' + 20 &= 0 \\ \Rightarrow 400x^2 + 800x' + 10000 &= 0 \end{aligned}} \right\} \text{Same classifier}$$



logistic Reg { overfit } large β 's } Unstable model.

Regularisation Solution \Rightarrow

Blaise
 Overfit, Probab of class 1
 $\frac{1}{0}$ $x_i \beta > 0$
 $x_i \beta < 0$

Loss $f(x) \Rightarrow$

$$-\sum_{i=1}^N y_i \log_e p_i + (1-y_i) \log_e (1-p_i) + \sum_{i=1}^p \beta_i^2$$



Linear Classification



Logistic Regression

- **Why we need regularisation in Logistic regression**

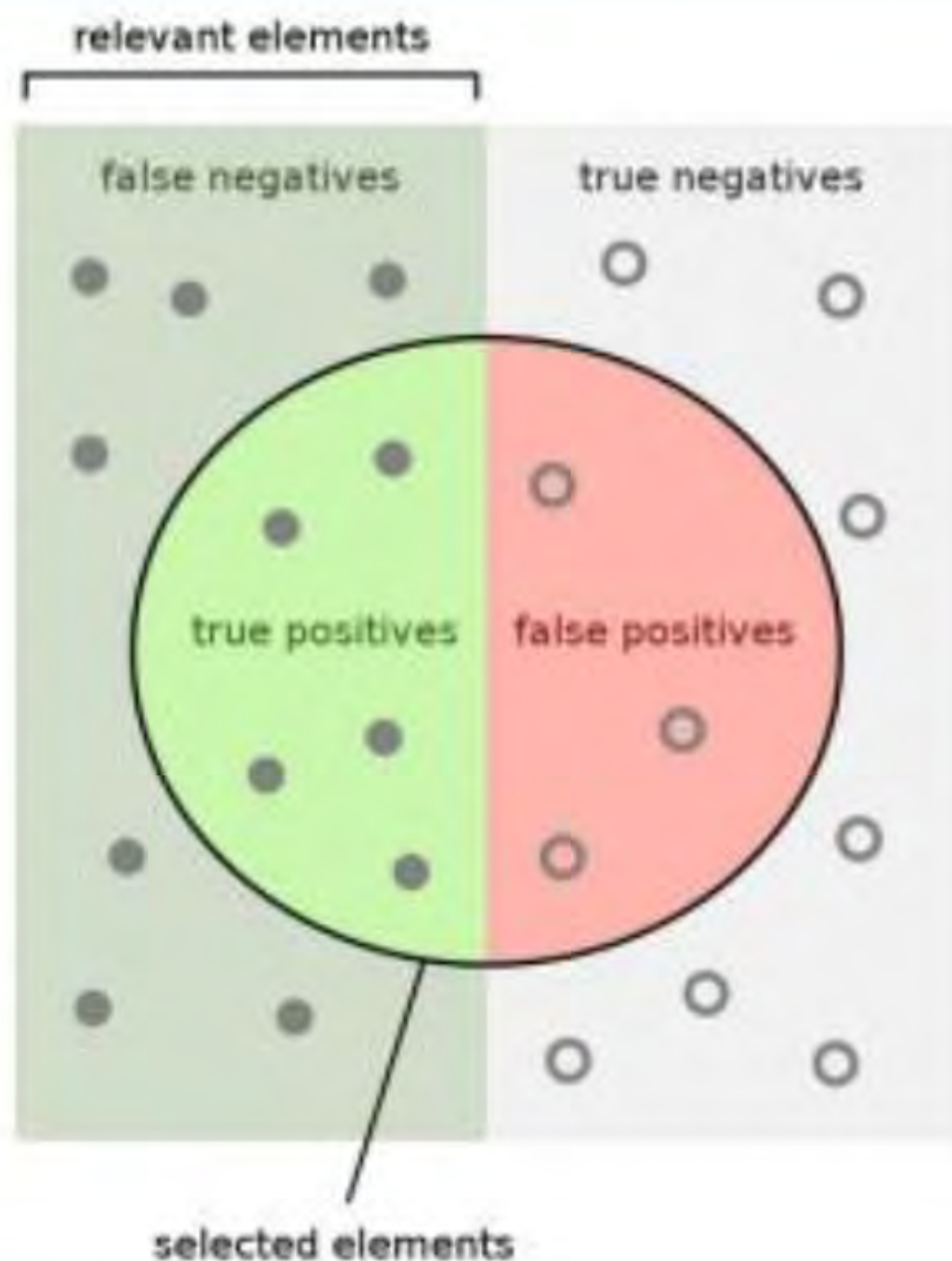


What is ROC curve (receiver operating characteristic curve)

- A receiver operating characteristic curve, or ROC curve, is a graphical plot that illustrates the performance of a binary classifier model (can be used for multi class classification as well) at varying threshold values.
- The ROC curve is the plot of the true positive rate (TPR) against the false positive rate (FPR) at each threshold setting



What is ROC curve (receiver operating characteristic curve)



How many relevant items are selected?
e.g. How many sick people are correctly identified as having the condition.

$$\text{Sensitivity} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

How many negative selected elements are truly negative?
e.g. How many healthy people are identified as not having the condition.

$$\text{Specificity} = \frac{\text{true negatives}}{\text{true negatives} + \text{false positives}}$$



What is ROC curve (receiver operating characteristic curve)

- **Sensitivity is a measure of how well a test can identify true positives**
- **Specificity is a measure of how well a test can identify true negatives:**

$$\text{sensitivity} = \frac{\text{number of true positives}}{\text{number of true positives} + \text{number of false negatives}}$$

$$\text{specificity} = \frac{\text{number of true negatives}}{\text{number of true negatives} + \text{number of false positives}}$$



What is ROC curve (receiver operating characteristic curve)

- What is TPR and FPR ?

True Positive Rate (TPR) is a synonym for recall and is therefore defined as follows:

$$TPR = \frac{TP}{TP + FN}$$

False Positive Rate (FPR) is defined as follows:

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

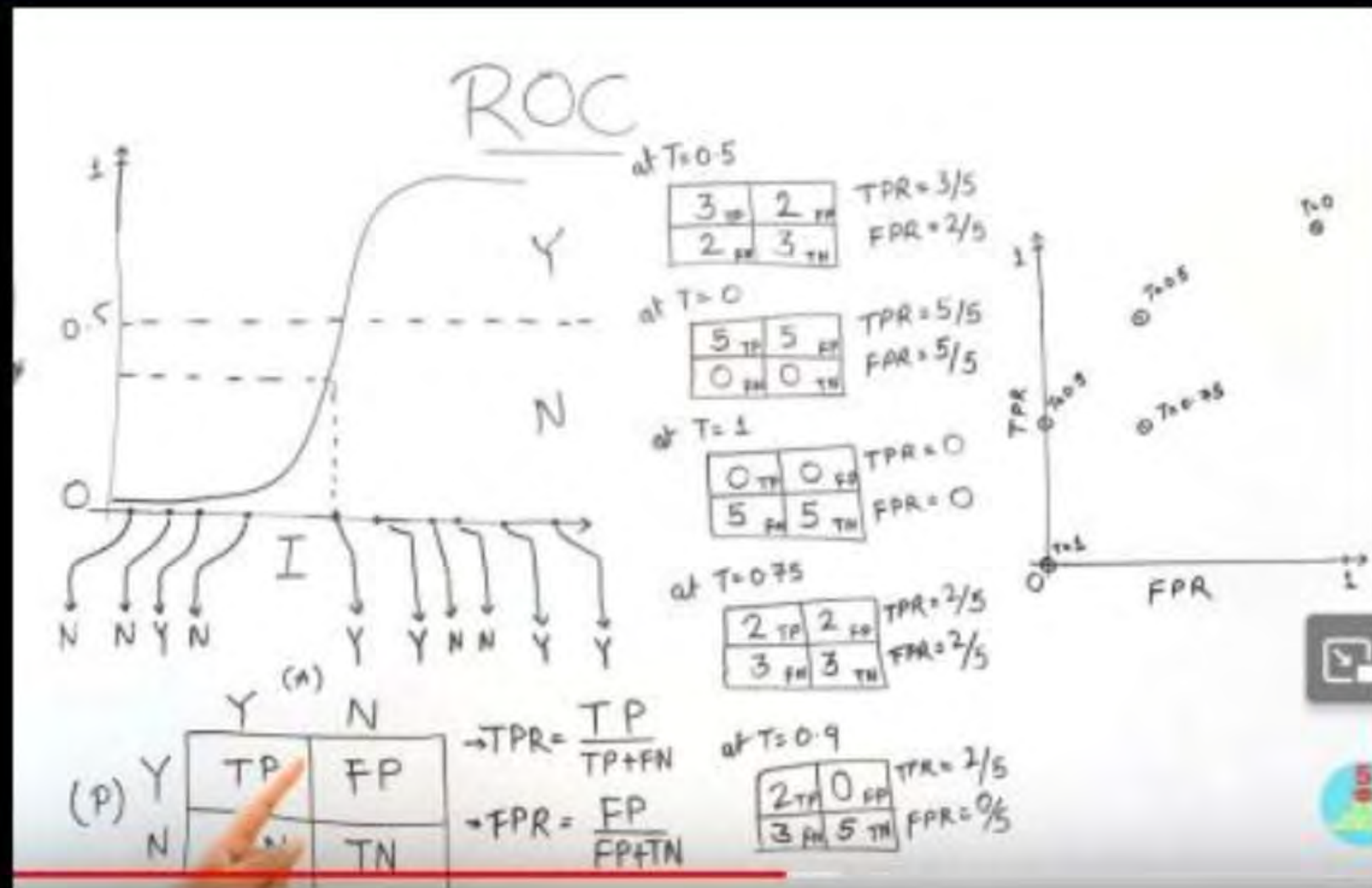


Linear Classification



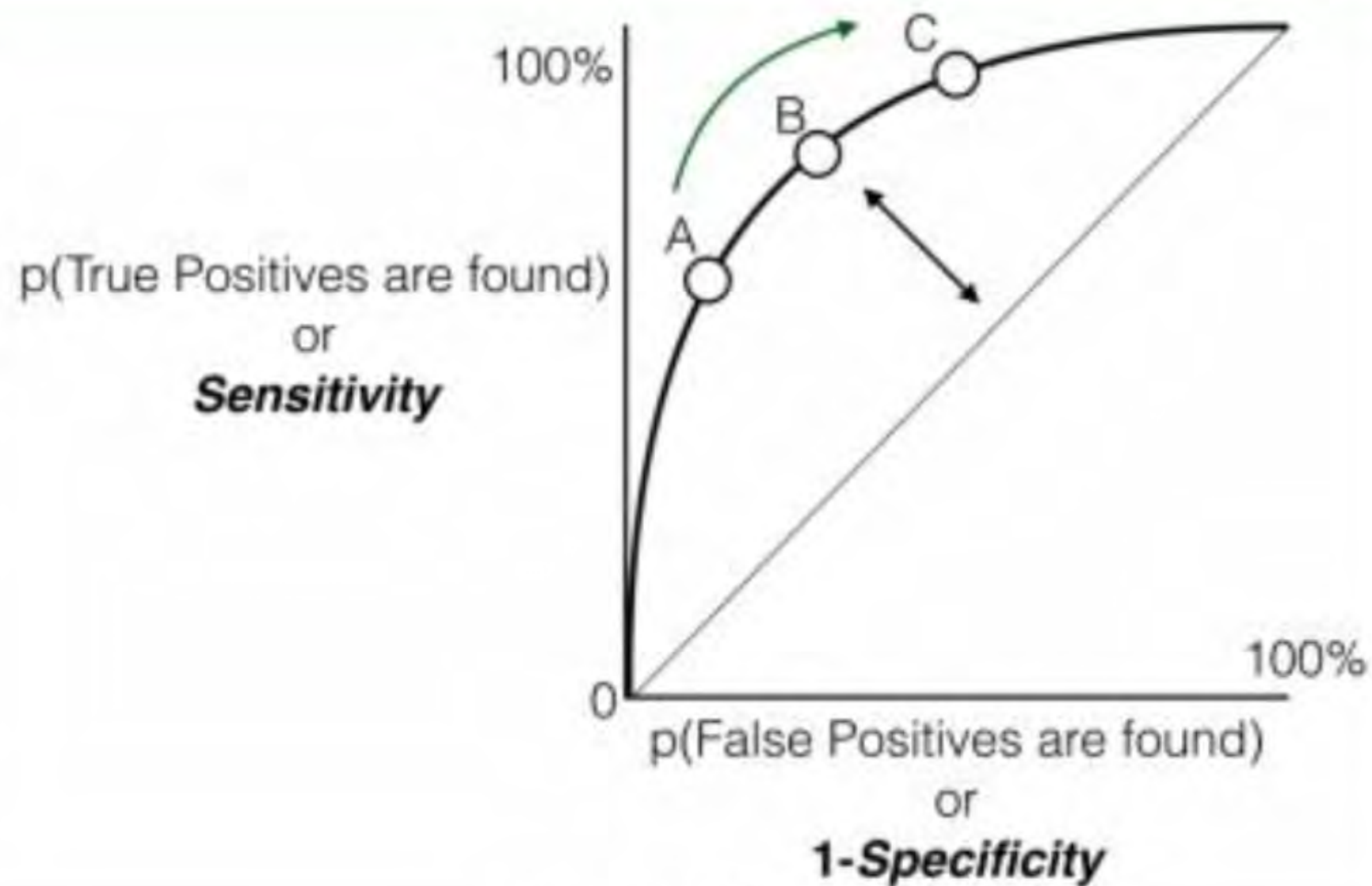
What is ROC curve (receiver operating characteristic curve) an example

The curve between TPR and FPR.





What is ROC curve (receiver operating characteristic curve) an example



Sensitivity versus False Positive Rate plot



What is AUC (Area under the curve)

- AUC stands for the Area Under the Curve, and the AUC curve represents the area under the ROC curve.
- It measures the overall performance of the binary classification model.
- The area will always lie between 0 and 1,
- A greater value of AUC denotes better model performance.
- Our main goal is to maximize this area in order to have the highest TPR and lowest FPR at the given threshold.
- The AUC measures the probability that the model will assign a randomly chosen positive instance a higher predicted probability compared to a randomly chosen negative instance.

THANK - YOU