# Recap of Previous Lecture

**Topic** — Logistic Regression

**Topic** — Cross Entropy

**Topic** — Log likelihood

**Topic**

**Topic**

# Topics to be Covered

| | |
|---|---|
| **Topic** | Multiclan → logistic regnemion |
| **Topic** | Now to make LR → NL |
| **Topic** | Confusion matrix |
| **Topic** | |
| **Topic** | |

DON'T Complain JUST DOIT

Skill

Naam ← Kaam

Izzat

# About the Faculty

- AIR 1 GATE 2021, 2023 (ECE).
- AIR 3 ESE 2015 ECE.
- M.Tech from IIT Delhi in VLSI.
- Published 2 papers in field of AI-ML.
- Paper 1 : Feature Selection through Minimization of the VC dimension.
- Paper 2 : Learning a hyperplane regressor through a tight bound on the VC dimension.

By- SIDDHARTH SABHARWAL SIR

**Logistic Regression :**

done

$\longrightarrow$ Concave

- Log likelihood

$$\max \sum_{i=1}^{N} y_i \log_e P_i + (1-y_i) \log_e(1-P_i)$$

$$CE = - \sum_{i=1}^{N} \left( \right)$$

minCE

$\longrightarrow$ Convex.

$\Rightarrow$ mean CE

## Logistic Regression : for multiclass case

\* logistic regression is for binary

\* How to use for multi class

Ex: we have data of 3 class.
we create 3 classifier one for each clan.

for Class 1 : Class 1 points '1'
Class 2,3 points 0

one versus rest

2 class data

→ apply logistic regression find $\beta^1$

y

| | |
|---|---|
| 1. | $y_1$ |
| 2. | $y_1$ |
| 3. | $y_2$ |
| 4. | $y_2 y_3 y_1$ |
| . | $y_1$ |
| . | $y_3 y_3$ |

**(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:

Sol.

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

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

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

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

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

Suppose after training our model we get $\vec{\beta} = \begin{bmatrix} -1.2 & -0.005 & 2.5 \end{bmatrix}^T$, 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 $\sigma$. If there is not enough information, write "not enough information".

<span style="color:gold">Sol.</span>

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

   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 $b$.

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.

Sol

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
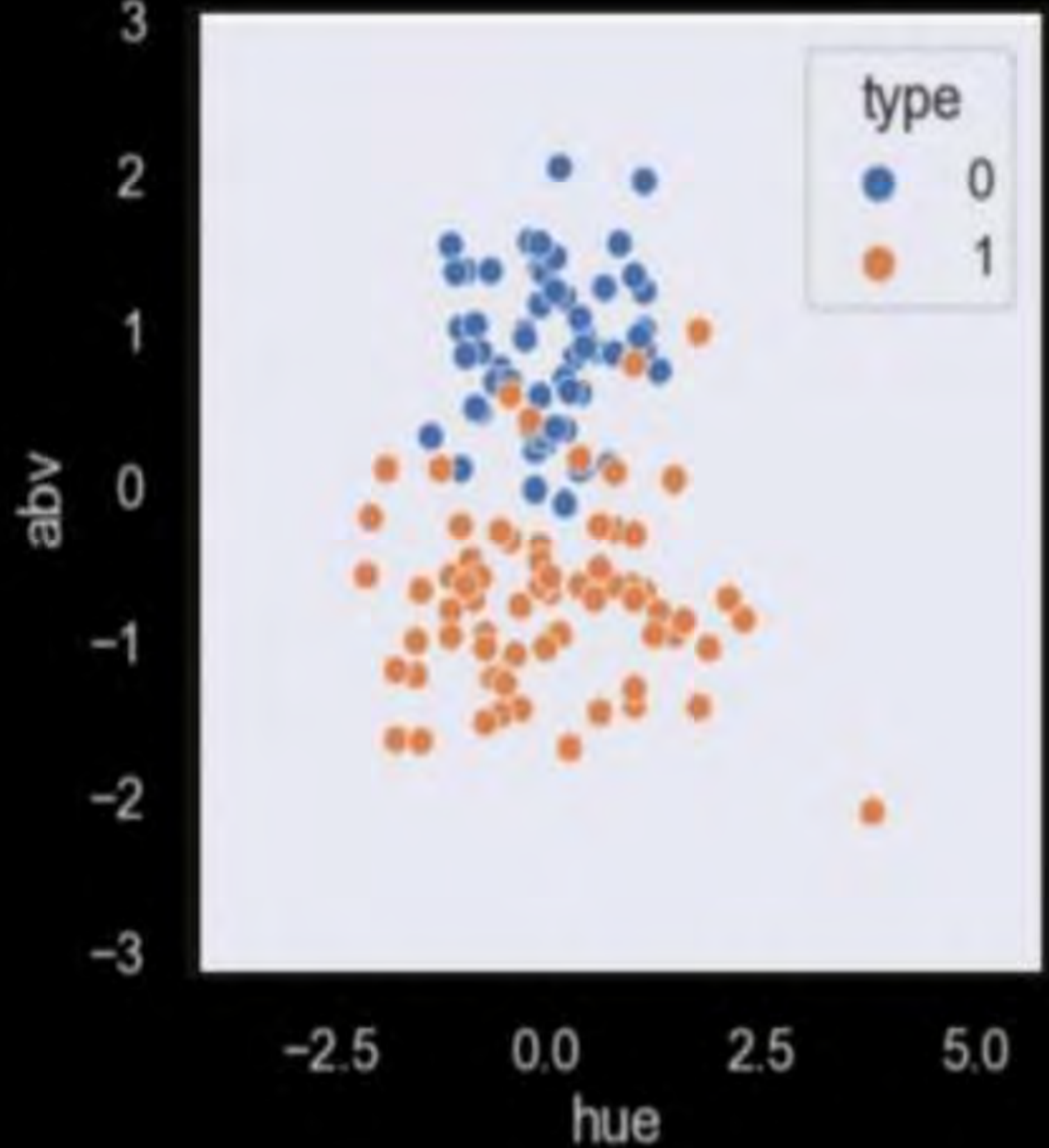
(D) $w_1$ and $w_2$ are negative, $w_3$ is positive

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.

Sol

☐ 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.

## 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

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

Sol.

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

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:

$$H\ H\ H\ H\ H\ H\ H\ H\ H\ T$$

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 $\frac{9}{10}$

Sol.

There are 5 balls in a bag. Each ball is either red or blue. Let $\theta$ (an integer) be the number of blue balls. We want to estimate $\theta$, 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 $\theta$ is fixed, what is the likelihood of getting exactly that sequence of colors (expressed as a function of $\theta$)?

Sol.

**Logistic Regression**

- **Extending the case for more than 2 class**

for class 2 : Class 2 points $\rightarrow$ '1'

Class 1, 3 points $\rightarrow$ '0'

$\underbrace{\qquad\qquad\qquad}$

2 class data

$\rightarrow$ logistic regression find $\beta^2$

## Logistic Regression

- **Extending the case for more than 2 class**

for class 3: class 3 points → '1'

Class 1, 2 points → '0'

$\underbrace{\hspace{6cm}}_{\text{2 class data}}$

→ logistic regression find $\beta^3$

we create 3 classifier

Now we have a test point

$X_t$

we find $X_t\beta^1, X_t\beta^2, X_t\beta^3$

whichever is max

decide class of $X_t$

**Logistic Regression**

- **What is softmax**

Sigmoid Convert distance $(x\beta)$ into poobability, 2 Class Case.

→ we have multiclass then distance Ko poobability Conversion Softmax. $x_t$ is test point, $\beta^1, \beta^2, \beta^3$ show 3 Classifier

$$P_{class1} = \frac{e^{x_t\beta^1}}{e^{x_t\beta^1} + e^{x_t\beta^2} + e^{x_t\beta^3}}$$

$$P_{class2} = \frac{e^{x_t\beta^2}}{e^{x_t\beta^1} + e^{x_t\beta^2} + e^{x_t\beta^3}}$$

$$\boxed{\text{Sigmoid Convert distance } (x\beta) \text{ into probability, 2 class Case.}}$$

→ we have multiclass then distance ko probability Conversion Softmax . $x_t$ is test point, $\beta^1, \beta^2, \beta^3$ show 3 Classifier

$$P_{class1} = \frac{e^{x_t \beta^1}}{e^{x_t \beta^1} + e^{x_t \beta^2} + e^{x_t \beta^3}}$$

$$P_{class2} = \frac{e^{x_t \beta^2}}{e^{x_t \beta^1} + e^{x_t \beta^2} + e^{x_t \beta^3}}$$

$$P_{class3} = \frac{e^{x_t \beta^3}}{e^{x_t \beta^1} + e^{x_t \beta^2} + e^{x_t \beta^3}}$$

• The max probab decide the predicted class

- **What is softmax loss – categorical cross entropy loss**

$$CE_{loss} \text{ for } 2\text{Class Case} \quad -\sum_{i=1}^{N} Y_i \log P_i + (1-Y_i) \log(1-P_i)$$

$$CE = -\sum_{i=1}^{N} \log P_{y_i}$$

original label

| $y$ | $P_i$ |
|-----|-------|
| 1 | .8 |
| 0 | .3 |
| 1 | .6 |
| 1 | .7 |
| 0 | .2 |

$$CE = -\left[ 1.\log \text{class 1 hone Ki Pred. Probab} + \right.$$
$$1. \log \text{Class 0} \quad " \quad " \quad " \quad " \quad +$$
$$1. \log \text{Clan 1} \quad " \quad " \quad " \quad " \quad +$$
$$1. \log \quad " \quad 1 \quad " \quad " \quad " \quad " \quad +$$
$$\left. 1. \log \quad " \quad 0 \quad " \quad " \quad " \quad \right]$$

So $CE = -\sum_{i=1}^{N} \log P_{y_i}$

This formulae Valid for multiclass Case also

datapoint No.   $y$

| datapoint No | $y$ |
|---|---|
| 1 | $y_1$ |
| 2 | $y_1$ |
| 3 | $y_3$ |
| 4 | $y_3$ |
| 5 | $y_2$ |
| 6 | $y_2$ |
| 7 | $y_2$ |

$$CE = -\left[ \log P_{11} + \log P_{21} + \log P_{33} + \log P_{43} + \log P_{52} + \log P_{62} + \log P_{72} \right]$$

Consider the following image data point, where the model's predicted probabilities for the classes (Dog, Cat, Mountain) are:

| Class | Probability |
|---|---|
| Dog | 0.06 |
| Cat | 0.0 |
| Mountain | 0.94 |

Assuming that the true class of the image is Dog, what is the cross-entropy loss for this data point?

A) $-\log(0.06)$

B) $-\log(0.94)$

C) $-0.6\log(0.6)$

*(handwritten annotations)*

3 class

Single data point

3 class

data belong to Clan DOG

$$CE = -\sum_{i=1}^{N} \log P_{yi}$$

$$CE = -\log \cdot 06$$

Assume our prediction and ground truth for the three classes for $i^{th}$ point is:

Probab of 3 clames.

$$CE =$$

$$-\log(P_{y_i})$$

$$-\log(.8)$$

$$\hat{y}_i = \begin{bmatrix} 0.1 \\ 0.8 \\ 0.1 \end{bmatrix} = \begin{bmatrix} \hat{y}_i^1 \\ \hat{y}_i^2 \\ \hat{y}_i^3 \end{bmatrix}$$

original label

$$y_i = \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} = \begin{bmatrix} y_i^1 \\ y_i^2 \\ y_i^3 \end{bmatrix}$$

## Logistic Regression

# How to classify this ?

2D data

Not linearly classifiable $\Rightarrow$ logictic Reg and linear Classification Cannot work on data

$x_2$

$(0,1)$

$(1,1)$

Each point has 2 Values $(x^1, x^2)$

$(0,0)$

$(1,0)$

- But we want to use logistic regression only.

$x_1$

$(0,1) \rightarrow (1,0,1)$  $(1,1) \rightarrow (2,1,1)$

$(1,0) \rightarrow (1,1,0)$  $(0,0) \rightarrow (0,0,0)$

2D data

$(x_1, x_2)$ $\xrightarrow{\text{data transform}}$ 3D. space

each data has
2 values

$$\begin{bmatrix} x_1+x_2, & x^1, & x^2 \end{bmatrix}$$

5D

$$\begin{bmatrix} x_1^2, & x_2^2, & x_1 x_2, & x_1^3, & x_2^3 \end{bmatrix}$$

$$2D \longrightarrow 5D$$

If after transformation we get linearly separable data → Now logistic regression is possible on this new (new/old) of 5D (2D/5D)

And in

In this case, @ higher dimension 3D, Classifier is ~linear/hyperplane~,

But in 2D lower dimension classifier is ~Nonlinear~.

Logistic Regression

No linearly seperable

$2D \longrightarrow 3D$

$x^2$

$(x^1-2)^2 + (x^2-2)^2 = 3^2$

$\left[ x^1, x^2, (x^1-2)^2 + (x^2-2)^2 - 9 \right]$

$(2,2)$

$3$

$x^2$

$x^3 > 0$

$x^3 < 0$

$x^1$

$x^2$

$x^3$

But $x^3 = 0$ here is

$(x_2^2 - 2)^2 + (x_1^2 - 2)^2 - 9 = 0$.

$x^2$

Jo ki Circle hai

No Linearly Seperable

$(x^1 - 2)^2 + (x^2 - 2)^2 = 3^2$

$(2,2)$  3

$x^2$

$X$  $X$  $X$  $X$  $X$  $X$

Classification Rule
Classification Plane
$x^3 = 0$ Plane   $x^3 > 0 \Rightarrow$ class 1
                  $x^3 < 0 \Rightarrow$ class 0

$2D \longrightarrow 3D$

$\left[ x^1, x^2, (x^1 - 2)^2 + (x^2 - 2)^2 - 9 \right]$

$\wedge x^3$

$x^3 > 0$

$x^2$

$x^1$

$x^3 < 0$

- Original data
not linearly seperable

higher dimension
mein le gaye

Linearly seperable
Plane that clarify
data

Linear
Plane

The Linear
Classifier in
higher dimension =
NL clasifier in original
Plane

**Confusion Matrix**

Confusion Matrix

To measure accuracy of Classifier

Predicted

Actual

| | P | N | |
|---|---|---|---|
| P | 50 | 10 | → Pred P = 60 |
| N | 20 | 80 | → Pred N = 100 |

Predicted

Actual P = 70    Actual N = 90

Actual

Predicted

→ Predicted P
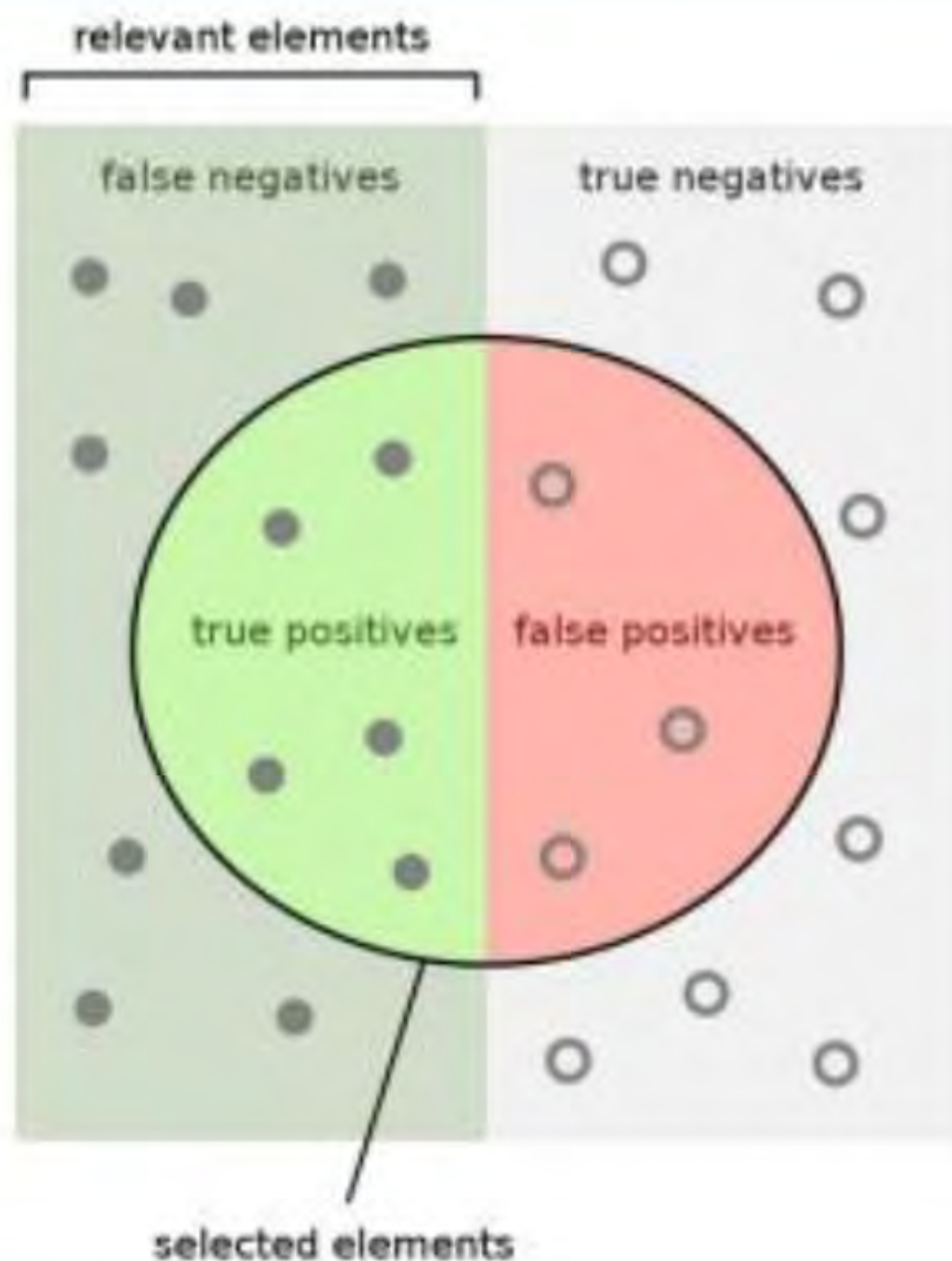
→ Predicted N

actual P    actual N

## Confusion Matrix

- Recall
- Precision

## 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)



relevant elements

false negatives | true negatives

true positives | false positives

selected elements

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

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

Sensitivity =

Specificity =

## What is ROC curve (receiver operating characteristic curve)

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

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

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

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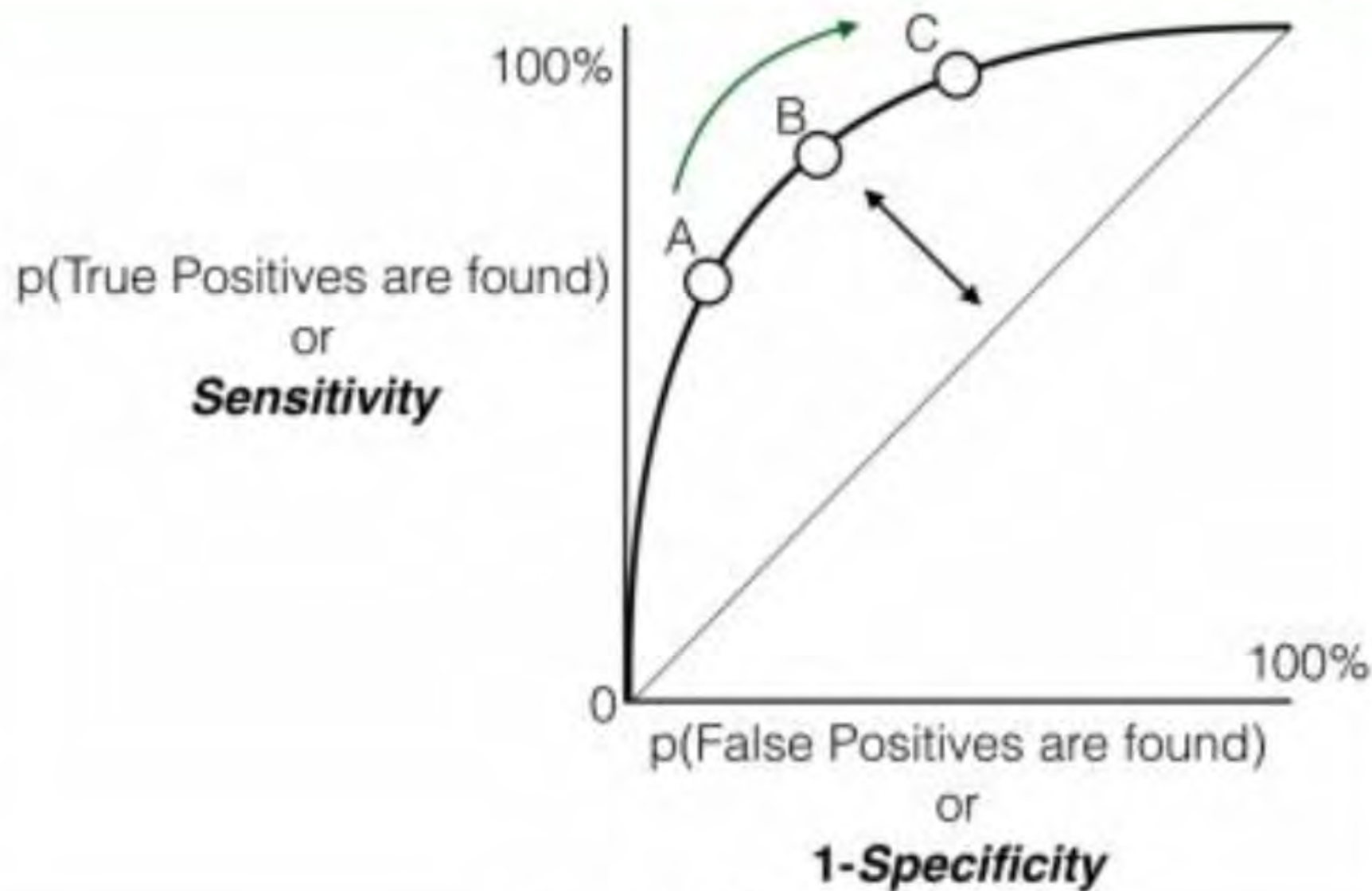
$$TPR = \frac{TP}{TP + FN}$$

False Positive Rate (FPR) is defined as follows:

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

## 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.
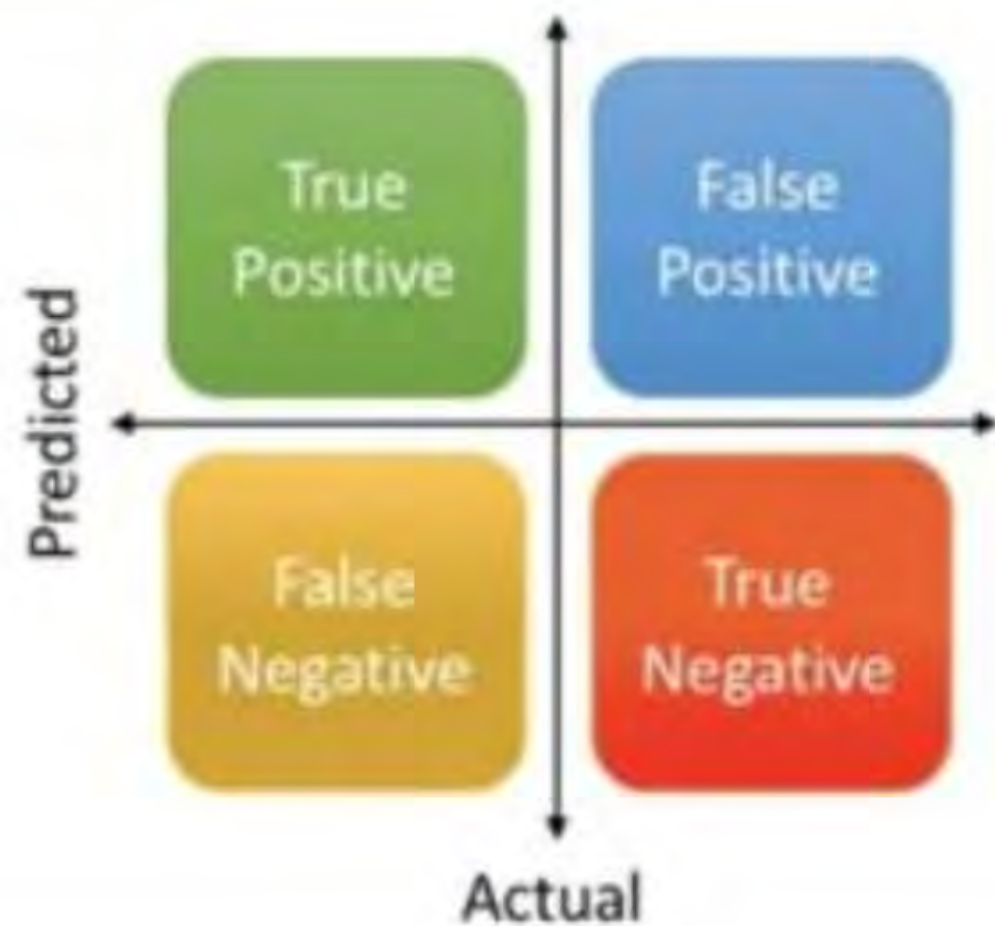
## What is Recall and Precision

$$\text{Precision} = \frac{\text{True Positive}}{\text{Actual Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$
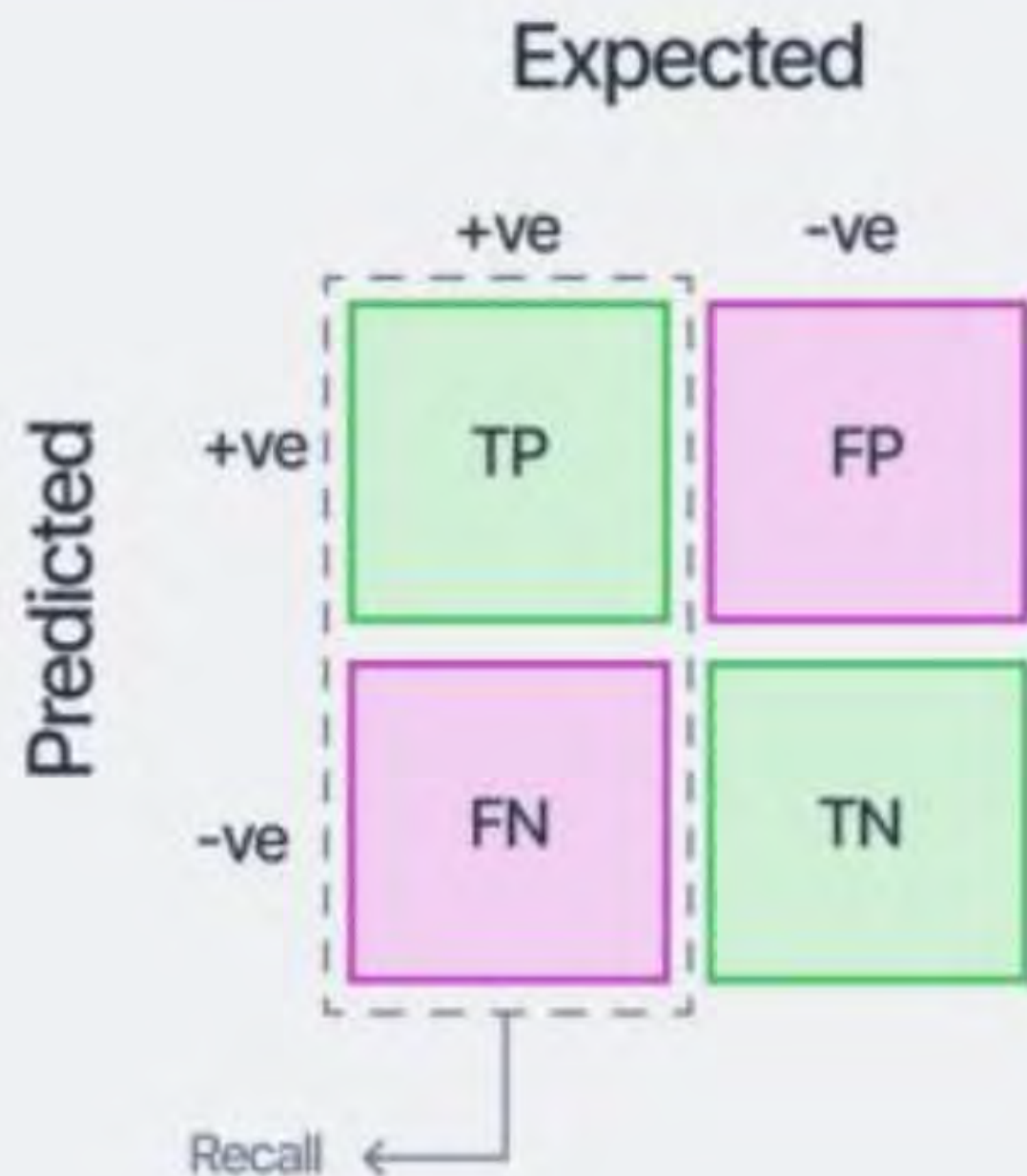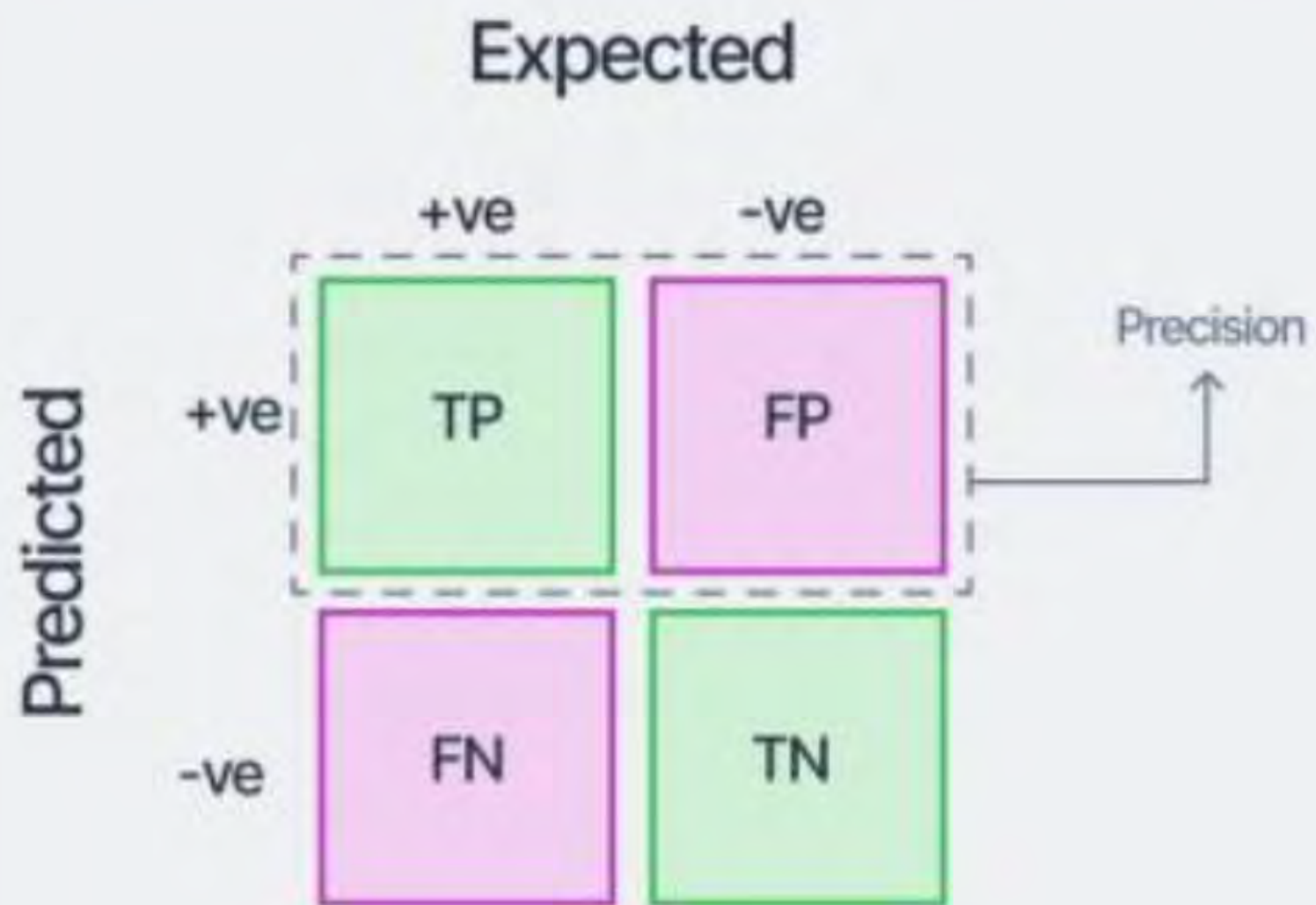
$$\text{Recall} = \frac{\text{True Positive}}{\text{Predicted Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

$$\text{Accuracy} = \frac{\text{True Positive} + \text{True Negative}}{\text{Total}}$$

# What is Recall and Precision

## What is Recall and Precision

Both precision and recall may be useful in cases where there is imbalanced data.

It may be valuable to prioritize one over the other in cases where the outcome of a false positive or false negative is costly.
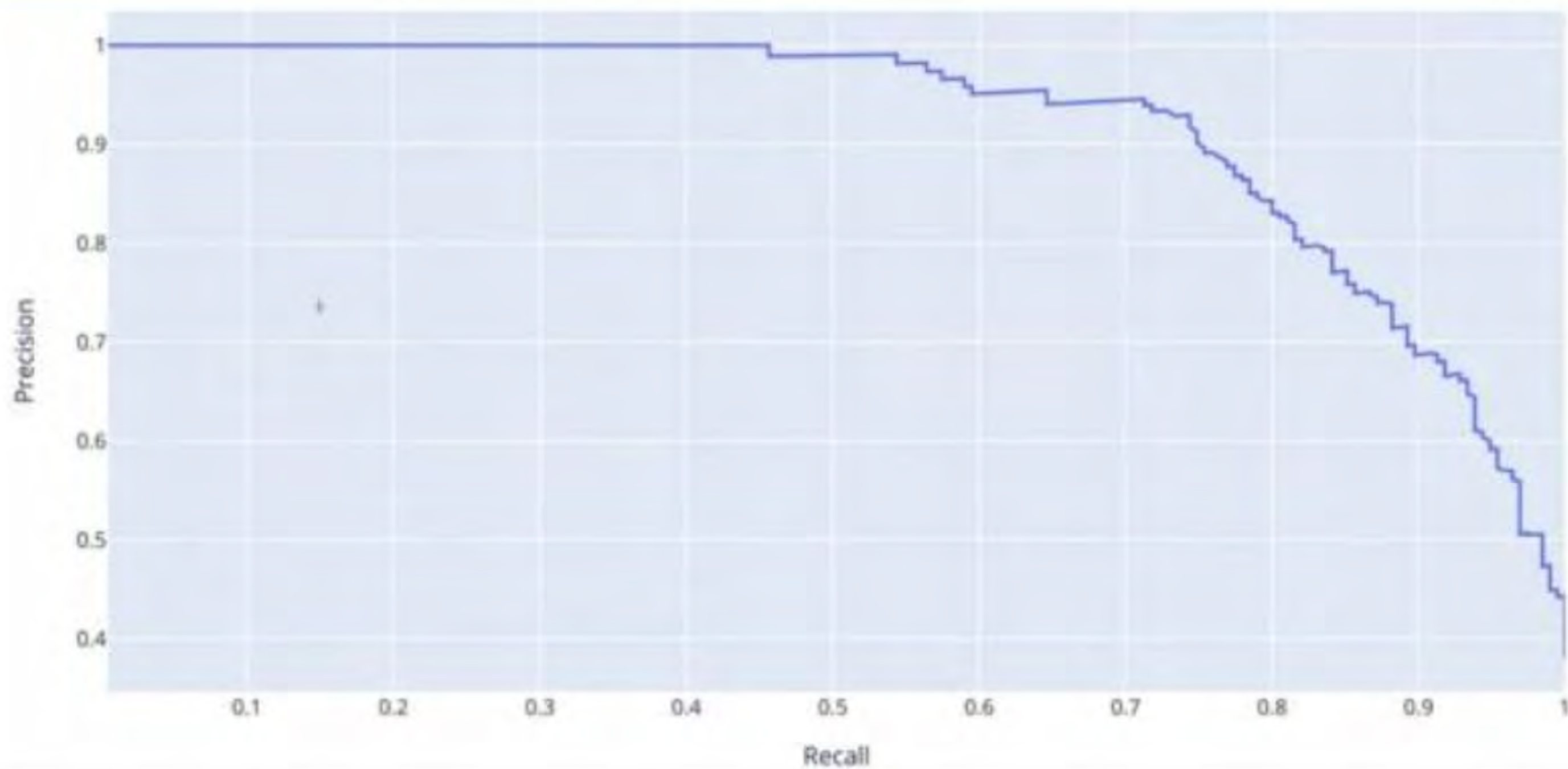
For example, in medical diagnosis, a false positive test can lead to unnecessary treatment and expenses.

In this situation, it is useful to value precision over recall. In other cases, the cost of a false negative is high.

For instance, the cost of a false negative in fraud detection is high, as failing to detect a fraudulent transaction can result in significant financial loss.

## What is Recall and Precision

## What is F-1 Score

In most problems, you could either give a higher priority to maximizing precision, or recall, depending upon the problem you are trying to solve. But in general, there is a simpler metric which takes into account both precision and recall, and therefore, you can aim to maximize this number to make your model better. This metric is known as F1-score, which is simply the harmonic mean of precision and recall.

$$\text{F1 Score} = 2 * \frac{\text{Precision} * \text{Recall}}{\text{Precision} + \text{Recall}}$$

## Practise

The confusion matrix visualizes the ____ of a classifier by comparing the actual and predicted classes.

- ⦿ Accuracy
- ⦿ Stability
- ⦿ Connectivity
- ⦿ Comparativity

## Practise

**From the above Table**

| n=200 | Prediction=NO | Prediction = YES |
|---|---|---|
| Actual = NO | 60 | 10 |
| Actual = YES | 5 | 125 |

○ In reality, there are totally 135 accounts who have a balance more than $1000 and 70 accounts with balance - less than $1000

○ In reality, there are totally 60 accounts who have a balance more than $1000 and 70 accounts with balance - less than $1000

○ In reality, there are totally 125 accounts who have a balance more than $1000 and 10 accounts with balance - less than $1000

○ In reality, there are totally 130 accounts who have a balance more than $1000 and 70 accounts with balance - less than $1000

**Practise**

For the below confusion matrix, what is the recall?

|         | Not 5 | 5    |
|---------|-------|------|
| Not 5   | 53272 | 1307 |
| 5       | 1077  | 4344 |

○ 0.7

○ 0.8

○ 0.9

○ 0.95

For the below confusion matrix, what is the precision?

|  | Not 5 | 5 |
|---|---|---|
| Not 5 | 53272 | 1307 |
| 5 | 1077 | 4344 |

○ 0.73

○ 0.76

○ 0.78

○ 0.82

## What is F-1 Score

F1 score is:

- ○ absolute mean of precision and recall

- ○ harmonic mean of precision and recall

- ○ squared mean of precision and recall

## What is F-1 Score

For the below confusion matrix, what is the F1 score?

|  | Not 5 | 5 |
| --- | --- | --- |
| Not 5 | 53272 | 1307 |
| 5 | 1077 | 4344 |

- ○ 0.72
- ○ 0.784
- ○ 0.82
- ○ 0.84

## What is F-1 Score

For a model to detect videos that are unsafe for kids, we need (safe video = postive class)

- ○ High precision, low recall

- ○ High recall, low precision

THANK - YOU