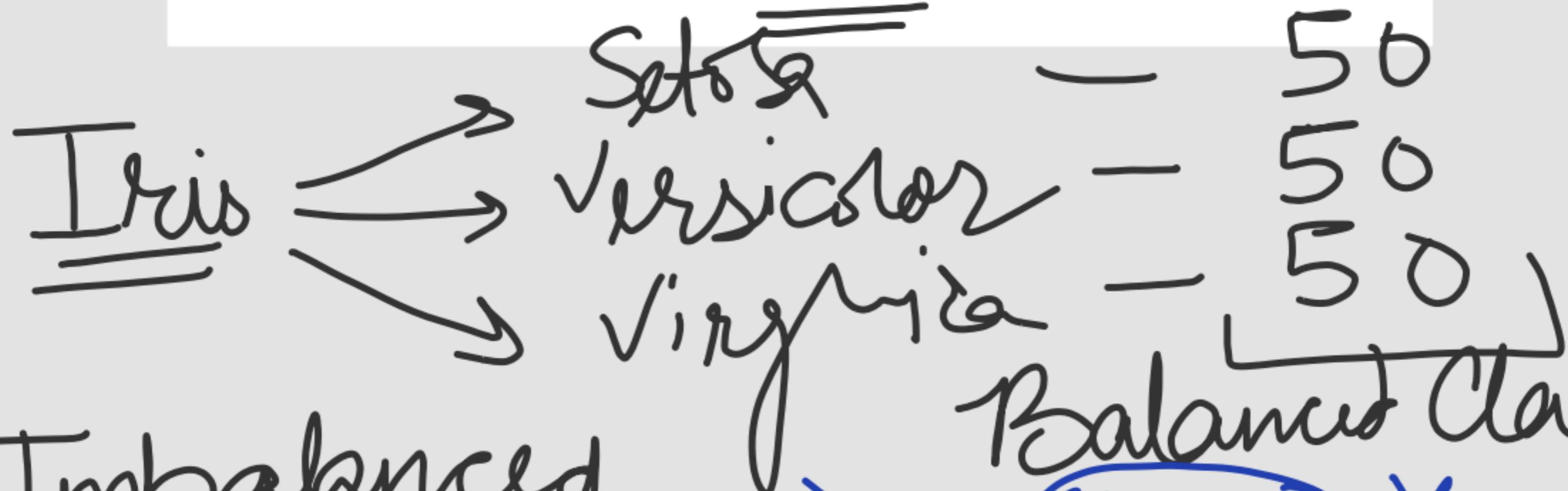


Imbalanced Class Problem



- ① Credit Card Fraud Detection
- ①: Class Y Transactions 100%
②: Legitimate 95%
③: Fraudulent 5%

Class Imbalance Problem

- Lots of classification problems where the classes are skewed (more records from one class than another)
 - Credit card fraud
 - Intrusion detection
 - Defective products in manufacturing assembly line
 - COVID-19 test results on a random sample

Aewma
classed

10 Samples

Test Dataset

Model

MD

0
8
-
0
8
-

Evaluate Performance
Accuracy

Predicted classes

90°

Challenges

- Evaluation measures such as accuracy are not well-suited for imbalanced class
- Detecting the rare class is like finding a needle in a haystack

Confusion Matrix

→ Summarizes the performance of classification

- Confusion Matrix:

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a	b
	Class>No	c	d

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

→ How many instances are correctly classified

Confusion Matrix

□ Confusion Matrix:

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	Positive	Negative
	POSITIVE	a TP	b FN
	Class=No	c FP	d TN
Negative			

a: TP (true positive)

b: FN (false negative)

c: FP (false positive)

d: TN (true negative)

→ Truly classified as +ve

→ Falsey classified as -ve

→ Falsey classified as +ve

→ Truly classified as -ve

		Predicted Class	
		Covid	Non-Covid
Actual Class	Covid	Positive No/0	- Negative
	Non-Covid	1	3
		1	3

Confusion Matrix

$$\text{Accuracy} = \frac{(1+3)}{8} = 50\%$$

		Actual class	
		C	N
R	C	C	N
	N	C	N
R	C	C	N
R	N	N	C
R	C	C	N
R	N	N	C
R	C	C	N
R	N	N	C
R	C	C	N
R	N	N	C

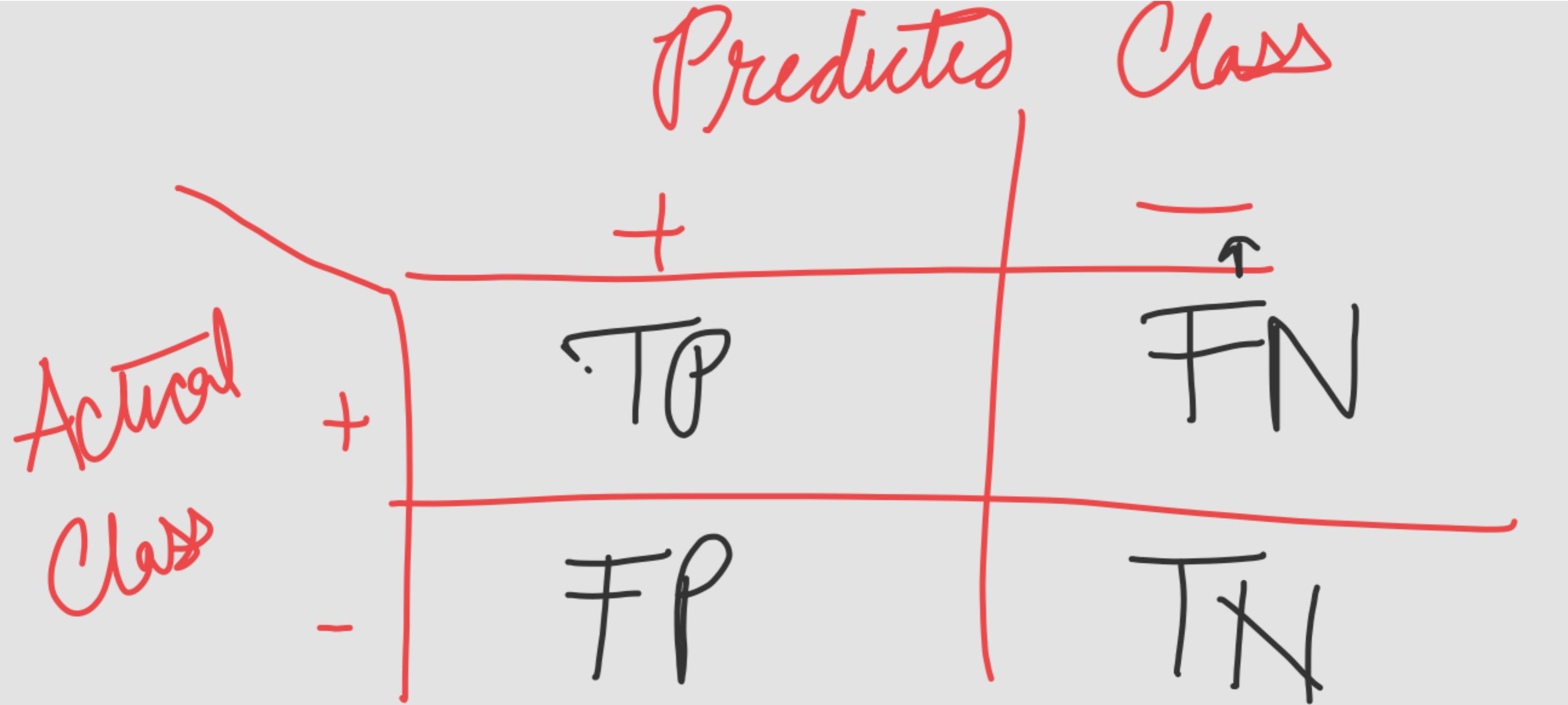
What is the problem with accuracy?
for imbalanced classes.

Covid - 10

		Predicted Class	
		C	NC
Actual Class	C	0	10
	NC	990	0

Accuracy = $\frac{990}{1000} \times 100\% = 99\%$

classify everything as NC.



Aim: $\frac{TP + TN}{TP + TN + FP + FN}$

Accuracy

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a (TP)	b (FN)
	Class>No	c (FP)	d (TN)

- Most widely-used metric:

$$\text{Accuracy} = \frac{a + d}{a + b + c + d} = \frac{TP + TN}{TP + TN + FP + FN}$$

Problem with Accuracy

- Consider a 2-class problem
 - Number of Class NO examples = 990
 - Number of Class YES examples = 10

Problem with Accuracy

- Consider a 2-class problem
 - Number of Class NO examples = 990
 - Number of Class YES examples = 10

		PREDICTED CLASS	
		Class=Yes	Class=No
ACTUAL CLASS	Class=Yes	TP	FN
	Class=No	FP	TN

True Positive Rate (TPR)

Out of all +ve samples, how many r correctly classified

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

False Negative Rate (FNR)

Out of all +ve samples, how many r incorrectly classified

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

Problem with Accuracy

- Consider a 2-class problem
 - Number of Class NO examples = 990
 - Number of Class YES examples = 10

- If a model predicts everything to be class NO, accuracy is $990/1000 = 99\%$
 - This is misleading because the model does not detect any class YES example
 - Detecting the rare class is usually more interesting (e.g., frauds, intrusions, defects, etc)

Accuracy

Which model is better?

A

		PREDICTED	
		Class=Yes	Class>No
ACTUAL	Class=Yes	0	10
	Class>No	0	990

99%

B

		PREDICTED	
		Class=Yes	Class>No
ACTUAL	Class=Yes	10	0
	Class>No	90	900

91%

Precision

$$\frac{TP}{TP+FP}$$

Actual
Spam Non-Spam

+ Non-Spam

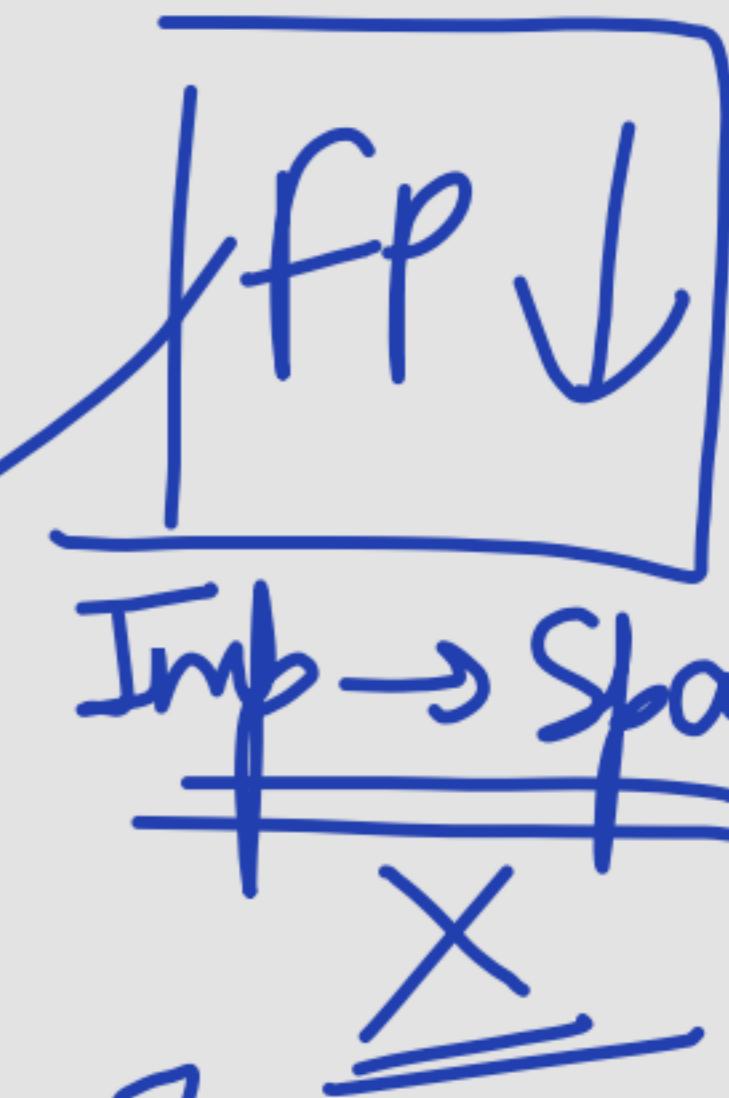
Fraction of records which are actually +Ve out of total records which are classified or predicted as +Ve.

Dont want to miss info.

Actual		Predicted	
Spam	Non-Spam	Spam	Non-Spam
TP	FN	FP	TN
Non-Spam			IF

Actual		Predicted		NS	
S	NS	Spam	Non-Spam	NS	IF
4	4	4	2	6	
NS					

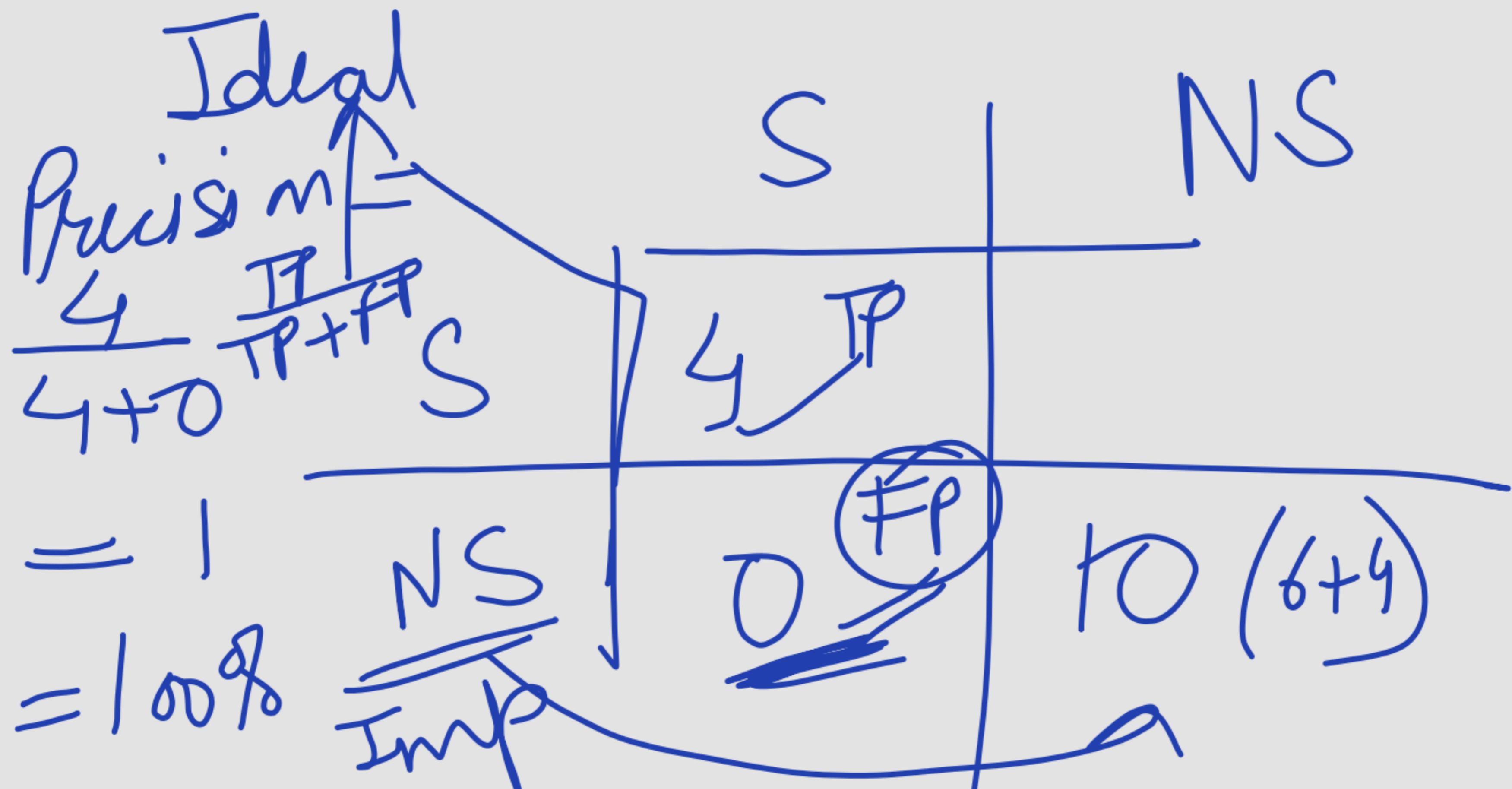
Classified as Spam = 8
actual Spams = 6



Precision

$$\frac{TP}{TP+FP} = \frac{4}{4+4} = 50\%$$

High Precision



		Predict	Recall
		Cancer	Healthy
Sensitivity	Actual	TP	FN
	Cancer	FP	TN
		TP + FN	TP + TN
		TP	TP + TN

I Out of all total +ve records how many are correctly classified as +ve

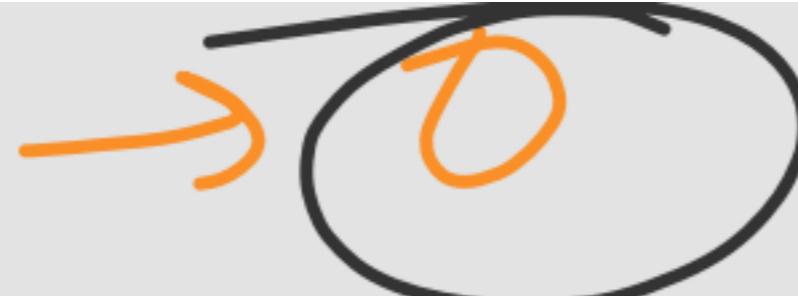
		Cancer	Healthy
		0	10
Cancer	Actual	0	10
	Predict	0	10

$$= 10$$

$$\text{Recall} = \frac{0}{10}$$

$$\text{Recall} = \frac{0}{10} = 0\%$$

$$\text{Recall} = \frac{10}{10} = 100\%$$

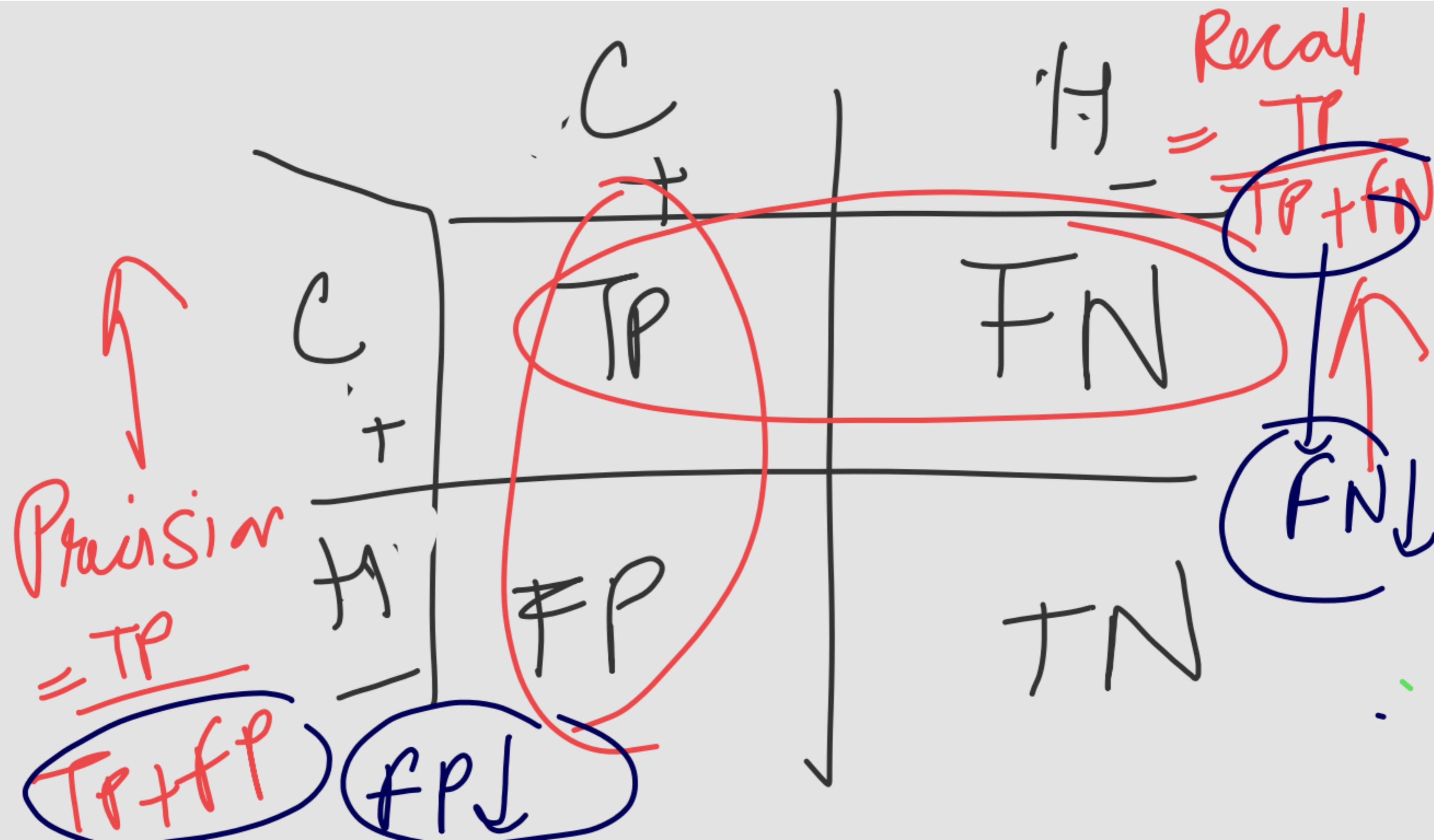
Recall → 

To Cancer → classified as Healthy

If a person has Cancer,
→ if model says ~~No~~

↳ There is a great risk

Recall ↑



For following confusion metrics, find accuracy, precision & recall

~~Model A~~

		-	+	
		9	1	$\sum = 10$
-	+	1	9	
	-	990	10	$\sum = 1000$

Acc	Pr	Re
99.1%	$\frac{1}{1+9}$ $\frac{1}{10}$	$\frac{1}{1+9}$ $\frac{1}{10}$

high precision & recall

~~Model B~~

		-	+	
		0	10	$\sum = 10$
-	+	10	90	
	-	90	900	$\sum = 990$

$\frac{10}{10+90} \times 100\% = 10\%$	$\frac{10}{10+0} \times 100\% = 100\%$

low precision
4

High Pr & High Recall $\rightarrow 100\%$

100%

$\frac{10}{10+90} \rightarrow FP$

$\frac{90}{90+0} \rightarrow FN$

High recall

Model	Precision	Recall	Measure
	$\frac{10}{10}$	$\frac{10}{10}$	$\frac{10+10}{2}$
2.	$\frac{10}{10}$	$\frac{10}{10}$	$\frac{10+10}{2} = 10$
3.	$\frac{10}{10}$	$\frac{10}{10}$	$\frac{10+10}{2} = 10$

As per arithmetic mean
 go for model 3

Let's break till
11:56 pm

Ques How to compare
precision & recall?

	Precision (P)	Recall (R)	Average ($\frac{P+R}{2}$)	CM
Model 1	0.5	0.4	0.45 ✓	+ - 0 0
Model 2	0.7	0.1	0.4	+ - 0 0
Model 3	0.02	1.0	0.51 ↗ <i>(Highest for model 3)</i>	+ - 0 0

Ques Not a good classifier
Predict y=1 all the time //

Ques 980

Let's break till
11:56 pm

Ques How to compare precision & recall?

F-Score
C measure.
Precision
Recall
F-Score
Mean

	Precision (P)	Recall (R)	Average ($\frac{P+R}{2}$)	F-Score
Model 1	0.5	0.4	0.45	0.44
Model 2	0.7	0.1	0.4	0.175
Model 3	0.02	1.0	0.51	0.53

(Not a good classifier
Predict y=1 all the time)

0 → 1 → 0.5 → 0

Harmonic Mean (#F Score)

$$= \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

$$P = \frac{TP}{TP+FP}$$

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

$$= \frac{2TP}{2TP+FP+FN}$$

$$\frac{2PR}{P+R}$$

Alternative Measures

		PREDICTED CLASS	
		Class=Yes	Class>No
ACTUAL CLASS	Class=Yes	a	b
	Class>No	c	d

$$\text{Precision (p)} = \frac{a}{a + c}$$

$$\text{Recall (r)} = \frac{a}{a + b}$$

$$\text{F - measure (F)} = \frac{2rp}{r + p} = \frac{2a}{2a + b + c}$$

Measures of Classification Performance

		PREDICTED CLASS	
		Yes	No
ACTUAL CLASS	Yes	TP	FN
	No	FP	TN

α is the probability that we reject the null hypothesis when it is true. This is a Type I error or a false positive (FP).

β is the probability that we accept the null hypothesis when it is false. This is a Type II error or a false negative (FN).

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

$$\text{ErrorRate} = 1 - \text{accuracy}$$

$$\text{Precision} = \text{Positive Predictive Value} = \frac{TP}{TP + FP}$$

$$\text{Recall} = \text{Sensitivity} = \text{TP Rate} = \frac{TP}{TP + FN}$$

$$\text{Specificity} = \text{TN Rate} = \frac{TN}{TN + FP}$$

$$\text{FP Rate} = \alpha = \frac{FP}{TN + FP} = 1 - \text{specificity}$$

$$\text{FN Rate} = \beta = \frac{FN}{FN + TP} = 1 - \text{sensitivity}$$

$$\text{Power} = \text{sensitivity} = 1 - \beta$$

$p=0$

$r=1$

$f\text{-measure} = 0$

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class>No
	Class=Yes	10	40
	Class>No	10	40

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class>No
	Class=Yes	25	25
	Class>No	25	25

		PREDICTED CLASS	
ACTUAL CLASS		Class=Yes	Class>No
	Class=Yes	40	10
	Class>No	40	10

$p=1$

$r=1$

$f\text{-measure} = 1$

Precision (p) = 0.5

TPR = Recall (r) = 0.2

FPR = 0.2

F – measure = 0.28

is closer + to lower values

Precision (p) = 0.5

TPR = Recall (r) = 0.5

FPR = 0.5

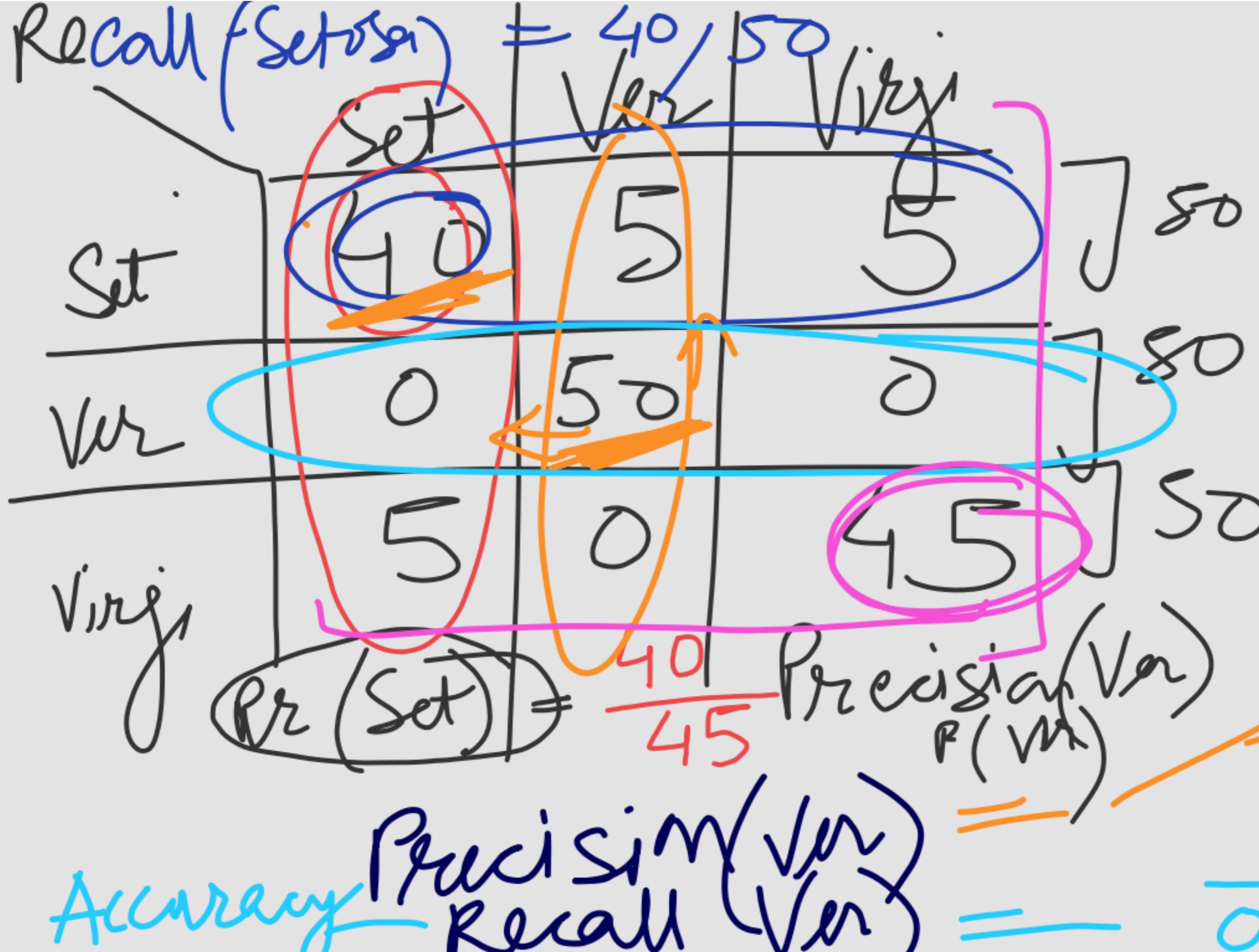
F – measure = 0.5

Precision (p) = 0.5

TPR = Recall (r) = 0.8

FPR = 0.8

F – measure = 0.61



$$Pr(\text{Set}) = \frac{40}{45}$$

$$\text{Precision}(\text{Ver}) = \frac{50}{5+80+0}$$

$$\text{Accuracy} = \frac{50}{0+50+0}$$

$$\text{Precision}(\text{Virg}) = \frac{45}{5+0+0}$$

$$\text{Recall}(\text{Ver}) = \frac{50}{50+0+0}$$

ROC (Receiver Operating Characteristic)

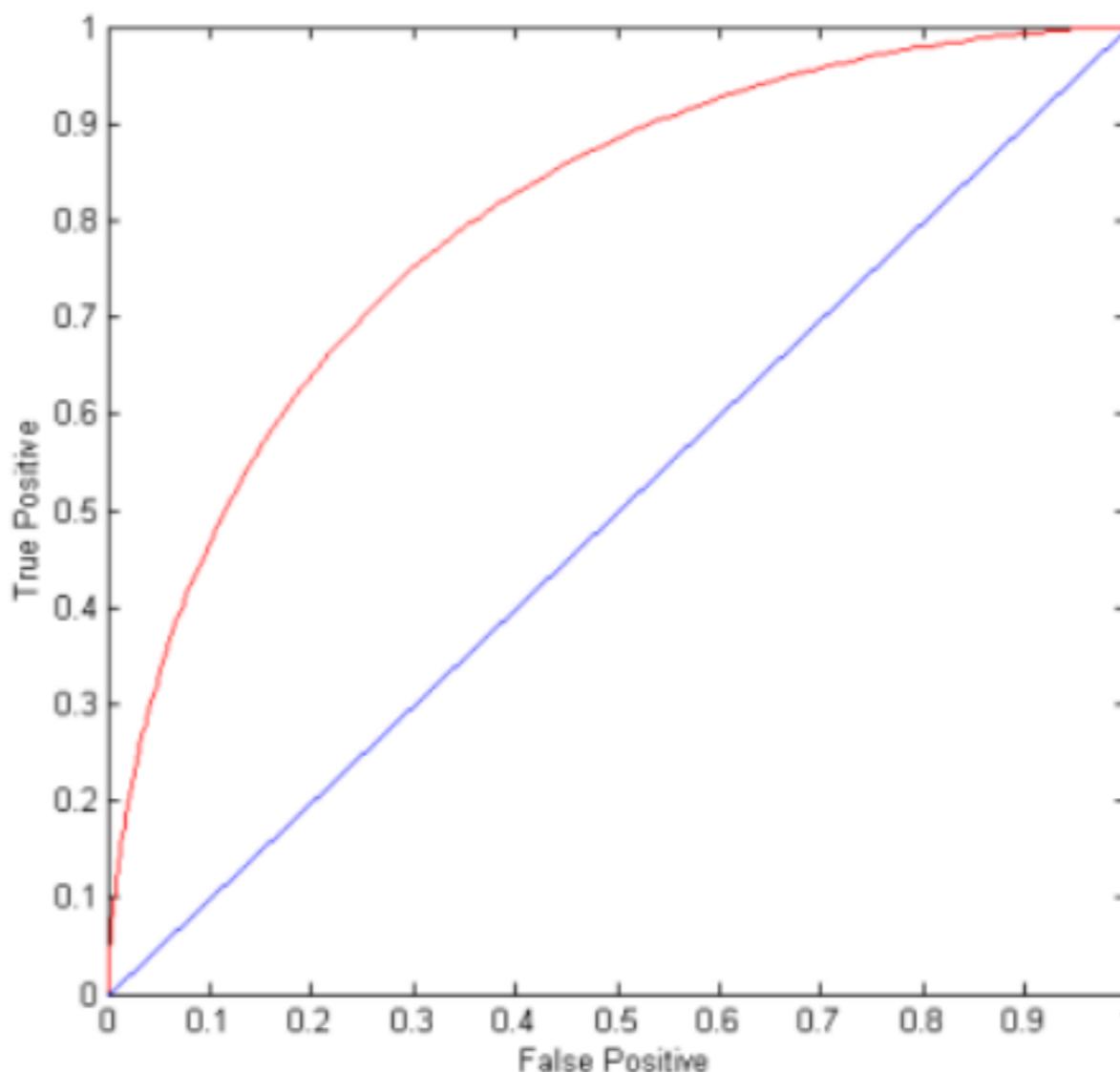
- A graphical approach for displaying trade-off between detection rate and false alarm rate
- Developed in 1950s for signal detection theory to analyze noisy signals
- ROC curve plots TPR against FPR
 - Performance of a model represented as a point in an ROC curve
 - Changing the threshold parameter of classifier changes the location of the point

ROC Curve

(TPR,FPR):

- (0,0): declare everything to be negative class
- (1,1): declare everything to be positive class
- (1,0): ideal

- Diagonal line:
 - Random guessing
 - Below diagonal line:
 - ◆ prediction is opposite of the true class



ROC (Receiver Operating Characteristic)

- To draw ROC curve, classifier must produce continuous-valued output
 - Outputs are used to rank test records, from the most likely positive class record to the least likely positive class record
- Many classifiers produce only discrete outputs (i.e., predicted class)
 - How to get continuous-valued outputs?
 - ◆ Decision trees, rule-based classifiers, neural networks, Bayesian classifiers, k-nearest neighbors, SVM

Evaluating the Performance of a Classifier

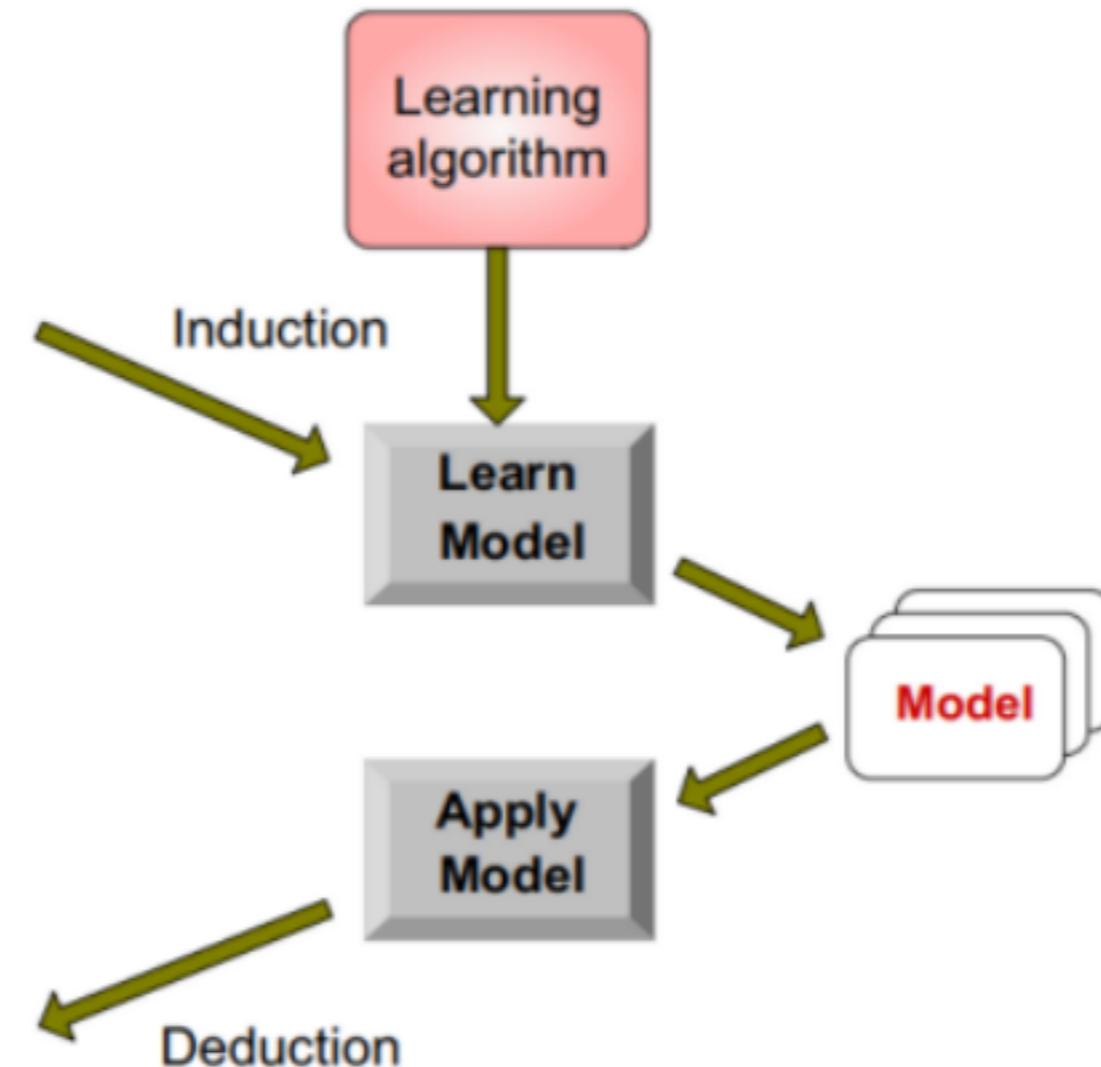
- Often useful to measure the performance of the model on the test set because such a measure provides an unbiased estimate of its generalization error.
- The accuracy or error rate computed from the test set can also be used to compare the relative performance of different classifiers on the same domain.
- Methods for evaluating classifier performance:
 1. Holdout Method
 2. Random Subsampling
 3. Cross-Validation
 4. Bootstrap

Tid	Attrib1	Attrib2	Attrib3	Class
1	Yes	Large	125K	No
2	No	Medium	100K	No
3	No	Small	70K	No
4	Yes	Medium	120K	No
5	No	Large	95K	Yes
6	No	Medium	60K	No
7	Yes	Large	220K	No
8	No	Small	85K	Yes
9	No	Medium	75K	No
10	No	Small	90K	Yes

Training Set

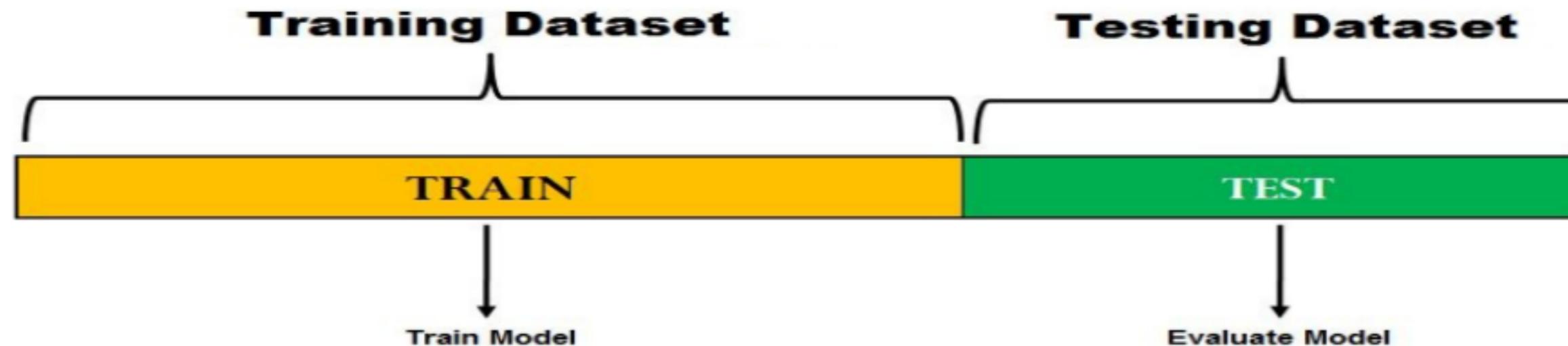
Tid	Attrib1	Attrib2	Attrib3	Class
11	No	Small	55K	?
12	Yes	Medium	80K	?
13	Yes	Large	110K	?
14	No	Small	95K	?
15	No	Large	67K	?

Test Set



1. Holdout Set

- In the holdout method, the original data with labeled examples is partitioned into two disjoint sets, called the training and the test sets, respectively.
- A classification model is then induced from the training set and its performance is evaluated on the test set.
- The proportion of data reserved for training and for testing is typically at the discretion of the analysts (e.g., 50-50 or two-thirds for training and one-third for testing).
- The accuracy of the classifier can be estimated based on the accuracy of the induced model on the test set.



Issues with the “Holdout” Method

- Separate sets for training and validation/testing
 - 1. Fewer labeled examples are available for training because some records are held out for testing
 - 2. Model may be highly dependent on composition of training and testing sets
 - ▣ *Small training sets will have greater variance*
 - ▣ *Small testing sets will be less reliable (will have wider confidence intervals)*

2. Random Sampling

- The holdout method can be repeated several times to improve the estimation of a classifier's performance.

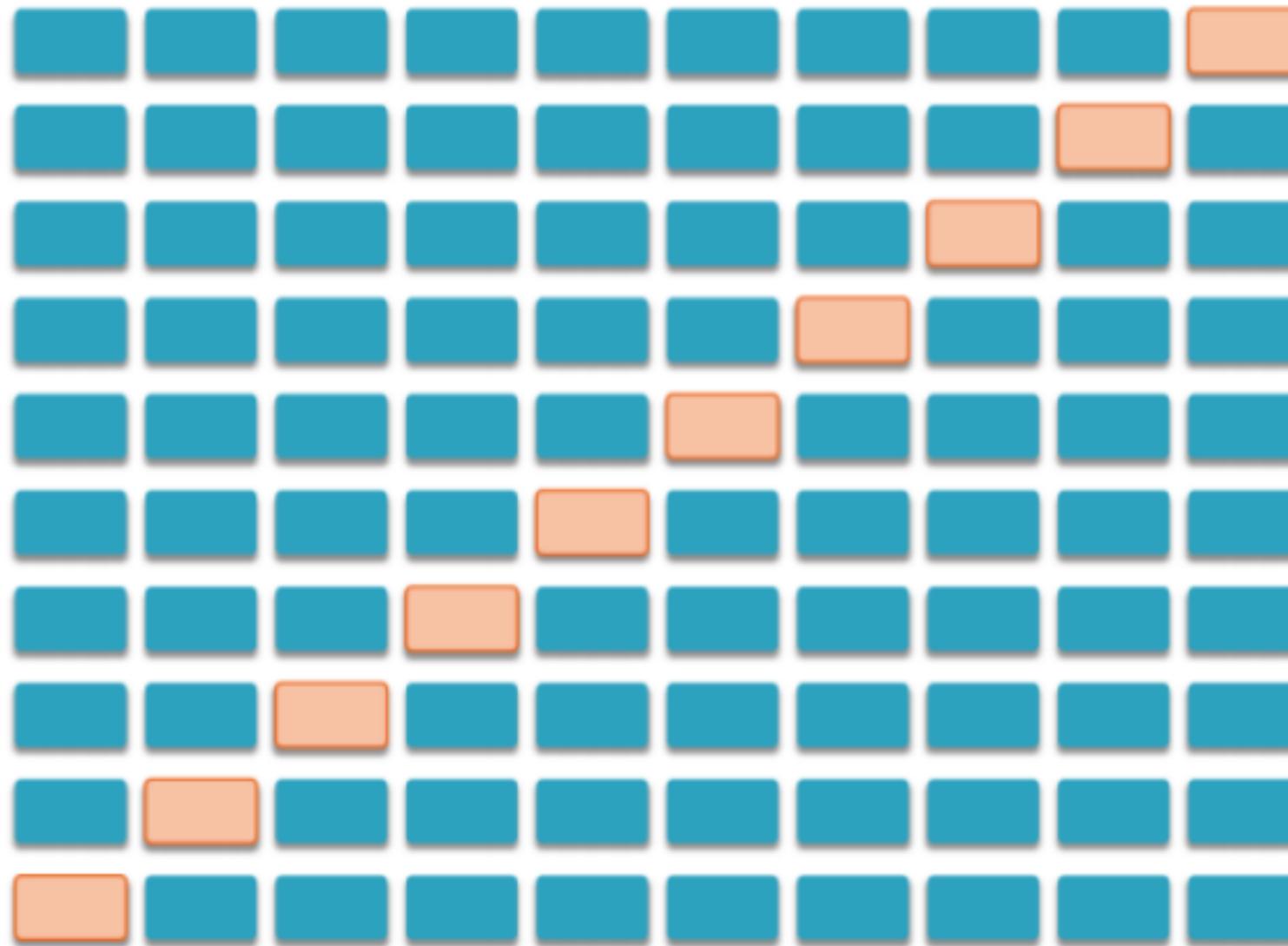
Let acc_i be the model accuracy during the i^{th} iteration. The overall accuracy is given by $acc_{\text{sub}} = \sum_{i=1}^k acc_i/k$. Random subsampling still encounters some

- ☹ No control over the number of times each record is used for training or testing.

3. Cross-Validation

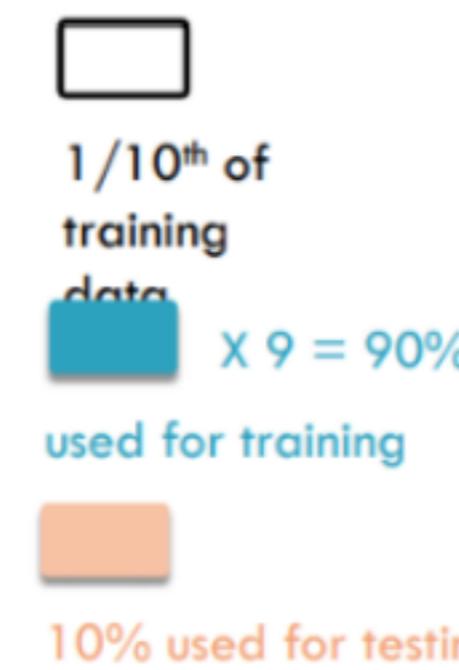
- Each record is used the same number of times for training and exactly once for testing. To illustrate this method, suppose we partition the data into two equal-sized subsets. First, we choose one of the subsets for training and the other for testing.
 - k-fold Cross Validation
 - $k = \text{number of folds (integer)}$
 - $k = 10$ is common
 - More computationally expensive

10-fold Cross Validation

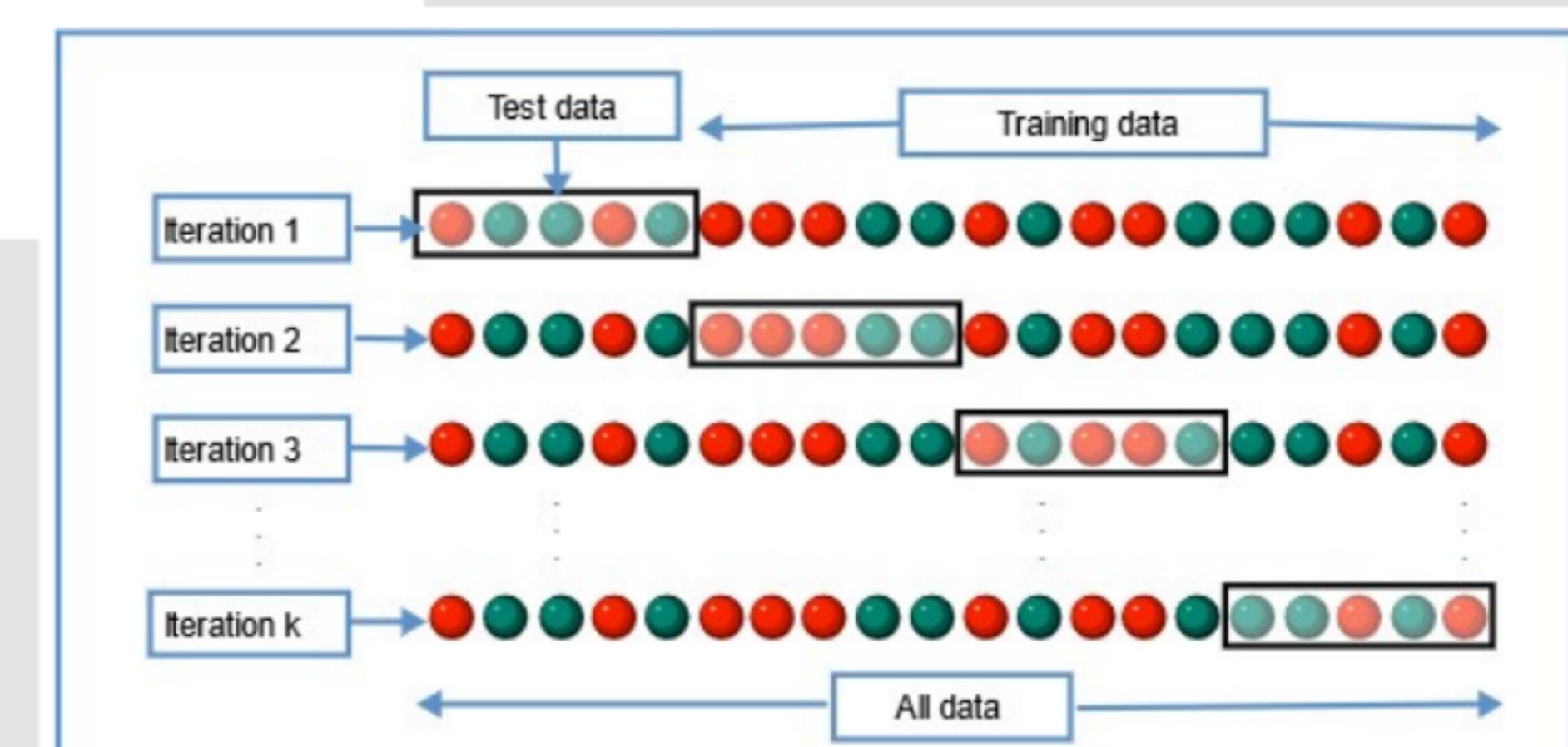


09/21/2020

Introduction to Data Mining, 2nd Edition



- Repeat k times
- Average results
- Each instance will be used once in testing



$$CV_k = \frac{1}{k} \sum_{i=1}^k ErrorRate_i$$

Average the error rate for each fold.

Leave-One-Out Cross-Validation

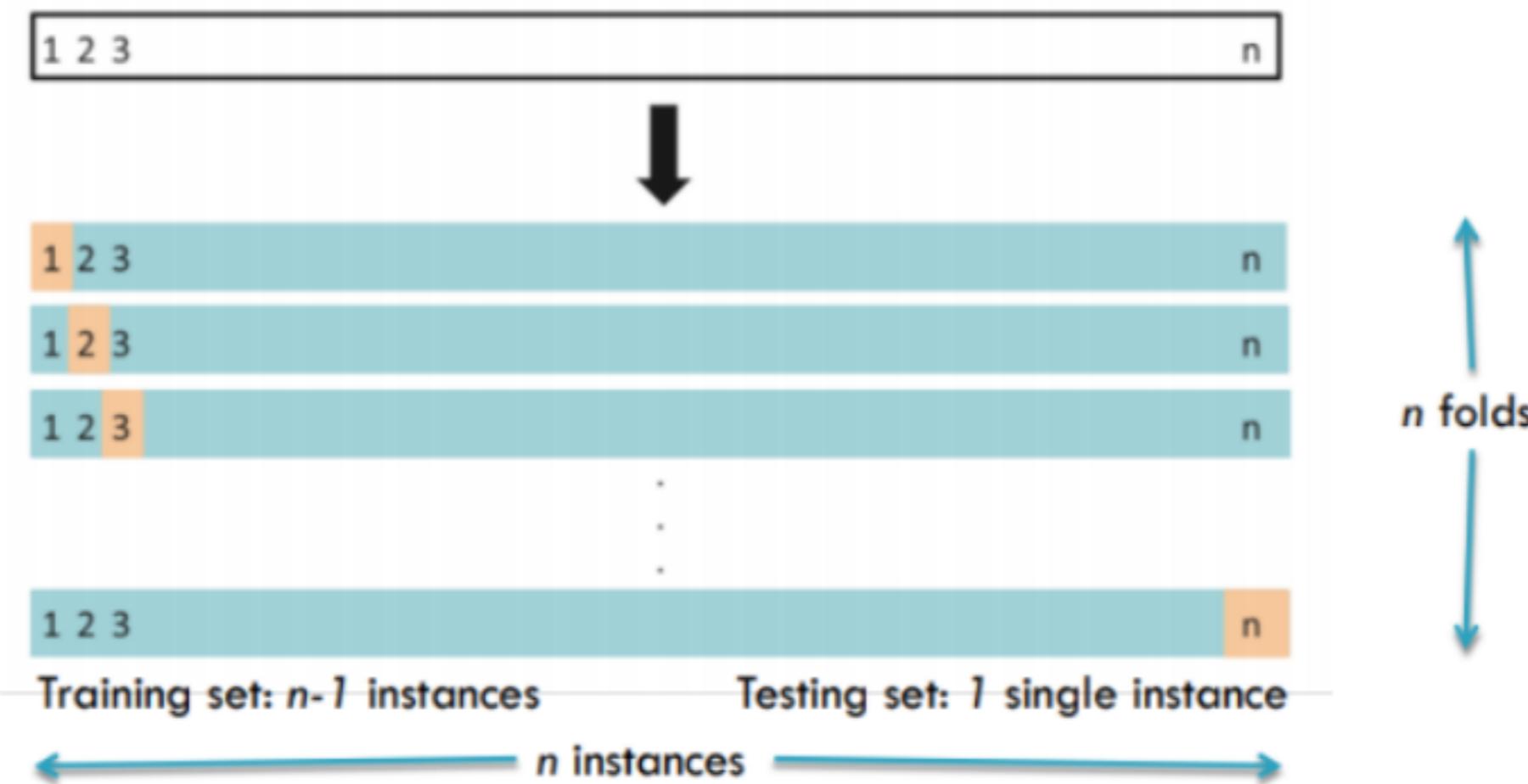
- k-fold Cross-Validation

- $k = \text{number of folds (integer)}$
- $k = 10$ is common

- Leave-One-Out Cross-Validation

- □ *Extreme:* $k = n$, where there are n observations in training+validation set
- Significantly more computationally expensive

Leave-One-Out Cross-Validation



Stratified K-Fold Cross-Validation

