

# Evaluate ML System

*Ha Noi, 11/12/2022,  
kan pham*

# Content:

1. Accuracy
2. Sensitivity, Specificity and prevalence
3. PPV and NPV
4. Confusion matrix
5. Type error I & II
6. Precision, Recall, and F1-score
7. ROC curve and AUC
8. evaluate multi-classification model

# 1. Accuracy

How good a classification model is?

$$\text{Accuracy} = \frac{\text{Examples of correctly of classified}}{\text{Total number of examples}}$$

# 1. Accuracy

How good a classification model is?



**Positive:** *Disease*

**Negative:** *Normal*

$$\text{Accuracy} = \frac{\text{Examples of correctly of classified}}{\text{Total number of examples}}$$

Ground Truth		
	Positive	Negative
Prediction	20	0
	20	60

*How many Accuracy ?*

# 1. Accuracy

## Accuracy in Terms of Conditional Probability

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= ?\end{aligned}$$



**Positive:** *Disease*

**Negative:** *Normal*

Ground Truth		
	Positive	Negative
Prediction	20	0
	20	60

# 1. Accuracy

## Accuracy in Terms of Conditional Probability

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= P(\text{Correct} \mid \text{Disease}) P(\text{Disease}) + P(\text{Correct} \mid \text{Normal}) P(\text{Normal}) \\ &= \underbrace{P(+ \mid \text{Disease}) P(\text{Disease})}_{\text{Sensitivity (TPR)}} + \underbrace{P(- \mid \text{Normal}) P(\text{Normal})}_{\text{Specificity (TNR)}}\end{aligned}$$

Sensitivity (TPR)

Specificity (TNR)



**Positive:** *Disease*

**Negative:** *Normal*

Ground Truth

	Positive	Negative
Positive	20	0
Negative	20	60

Prediction

# 2. Sensitivity and Specificity

## Accuracy in Terms of Conditional Probability

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= P(\text{Correct} \mid \text{Disease}) P(\text{Disease}) + P(\text{Correct} \mid \text{Normal}) P(\text{Normal}) \\ &= \underbrace{P(+ \mid \text{Disease}) P(\text{Disease})}_{\text{Sensitivity (TPR)}} + \underbrace{P(- \mid \text{Normal}) P(\text{Normal})}_{\text{Specificity (TNR)}}\end{aligned}$$

Sensitivity (TPR)

Specificity (TNR)



**Positive:** *Disease*

**Negative:** *Normal*

Prediction

Ground Truth		
	Positive	Negative
Positive	20	0
Negative	20	60

# 2. Sensitivity and Specificity

## Accuracy in Terms of Conditional Probability



**Positive:** *Disease*  
**Negative:** *Normal*

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= P(\text{Correct} \mid \text{Disease}) P(\text{Disease}) + P(\text{Correct} \mid \text{Normal}) P(\text{Normal}) \\ &= \underbrace{P(+ \mid \text{Disease}) P(\text{Disease})}_{\text{Sensitivity (TPR)}} + \underbrace{P(- \mid \text{Normal}) P(\text{Normal})}_{\text{Specificity (TNR)}}\end{aligned}$$

Sensitivity (TPR)

Specificity (TNR)

*How could you literally define **Sensitivity** and **Specificity** ?*



# 2. Sensitivity and Specificity

## Accuracy in Terms of Conditional Probability



**Positive:** *Disease*  
**Negative:** *Normal*

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= P(\text{Correct} \mid \text{Disease}) P(\text{Disease}) + P(\text{Correct} \mid \text{Normal}) P(\text{Normal}) \\ &= \underbrace{P(+ \mid \text{Disease}) P(\text{Disease})}_{\text{Sensitivity (TPR)}} + \underbrace{P(- \mid \text{Normal}) P(\text{Normal})}_{\text{Specificity (TNR)}}\end{aligned}$$

Sensitivity (TPR)

Specificity (TNR)

*How could you literally define **Sensitivity** and **Specificity** ?*

**$P(+ \mid \text{Disease})$**

If a patient has disease, what is probability that model predict disease?

***Sensitivity***

**$P(- \mid \text{Normal})$**

If a patient is normal, what is probability that model predicts normal?

***Specificity***

## 2. Sensitivity, Specificity and prevalence

### Accuracy in Terms of Conditional Probability



**Positive:** *Disease*

**Negative:** *Normal*

$$\begin{aligned}\text{Accuracy} &= P(\text{Correct}) \\ &= P(\text{Correct, Disease}) + P(\text{Correct, Normal}) \\ &= P(\text{Correct} \mid \text{Disease}) P(\text{Disease}) + P(\text{Correct} \mid \text{Normal}) P(\text{Normal}) \\ &= \underbrace{P(+ \mid \text{Disease}) P(\text{Disease})}_{\text{Sensitivity (TPR)}} + \underbrace{P(- \mid \text{Normal}) P(\text{Normal})}_{\text{Specificity (TNR)}}\end{aligned}$$

Sensitivity (TPR)

Specificity (TNR)

$$= \text{Sensitivity} * P(\text{Disease}) + \text{Specificity} * P(\text{Normal})$$

# 3. PPV and NPV

$P(+ | \text{Disease})$



$P(\text{Disease} | +)$

If a patient has disease, what is probability that model predict disease?

**Sensitivity**

If a model prediction is positive, what is probability that patient has the disease?

**PPV**  
(positive predictive value)



**Positive:** *Disease*  
**Negative:** *Normal*

Ground Truth

Prediction

	Positive	Negative
Positive	20	0
Negative	20	60

# 3. PPV and NPV

$P(- | \text{Normal})$



$P(\text{Normal} | -)$

If a patient is normal,  
what is probability  
that model predict  
normal?

**Specificity**

If a model prediction is  
negative, what is  
probability that patient  
is normal?

**NPV**  
(negative predictive  
value)



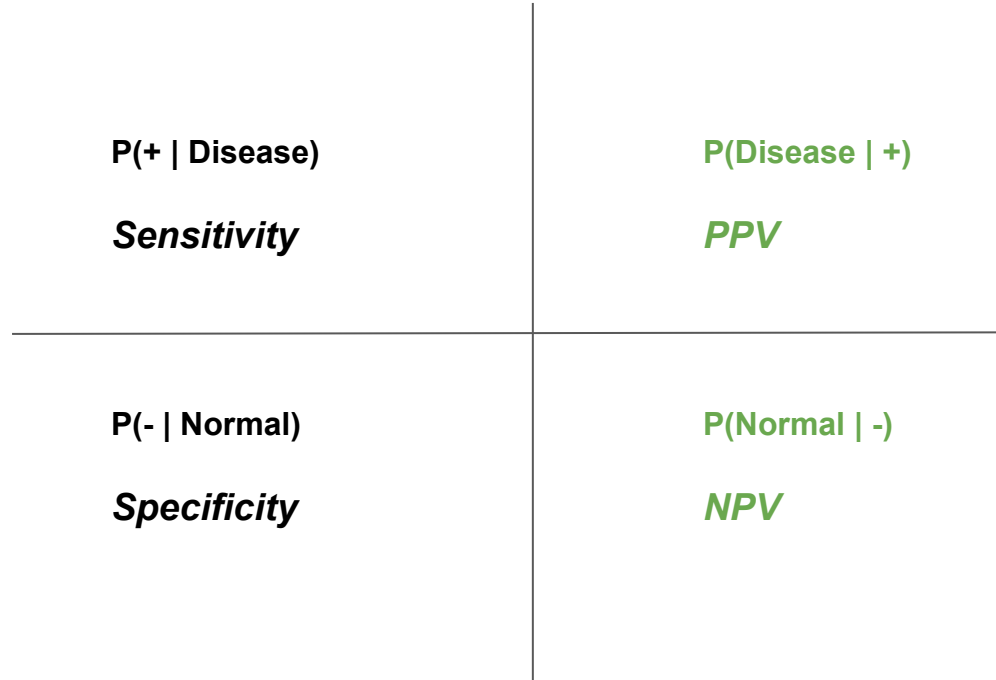
**Positive:** *Disease*  
**Negative:** *Normal*

Ground Truth

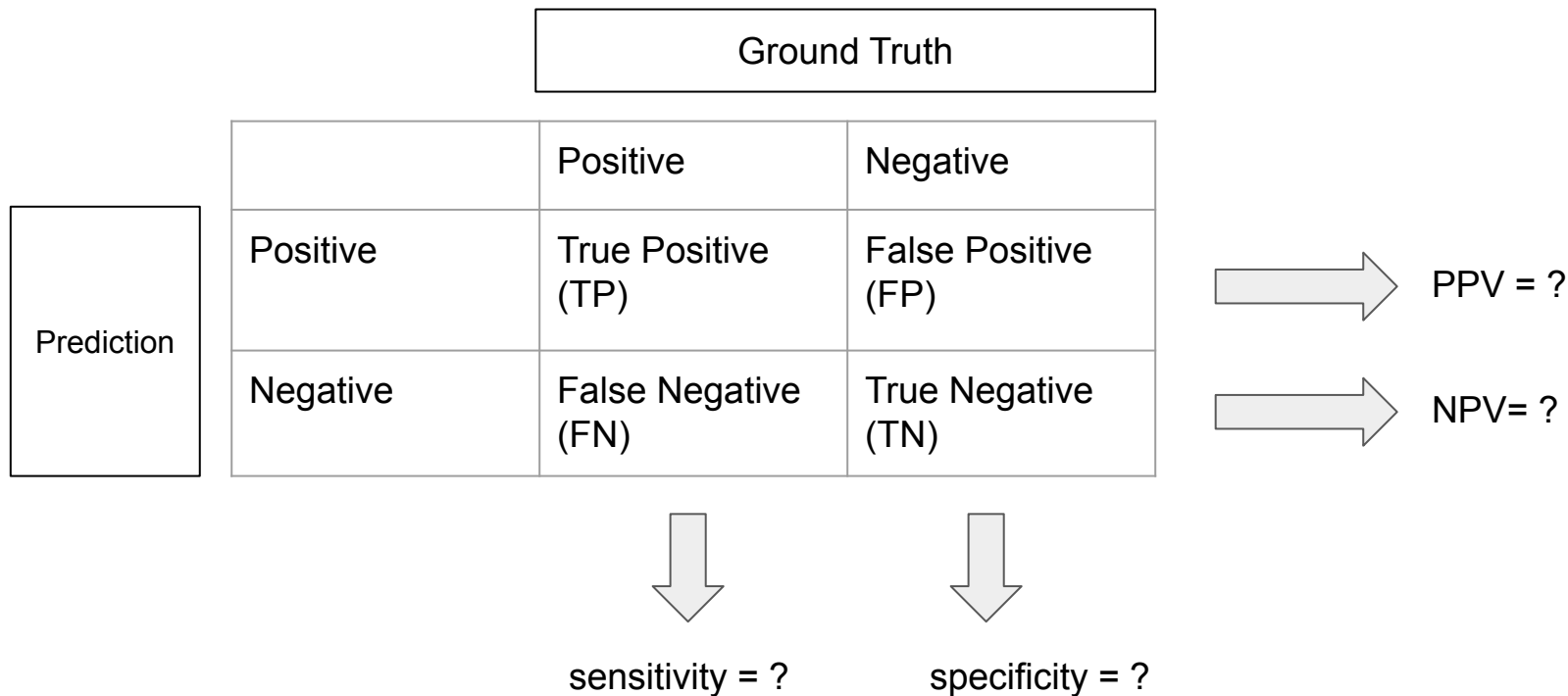
	Positive	Negative
Positive	20	0
Negative	20	60

Prediction

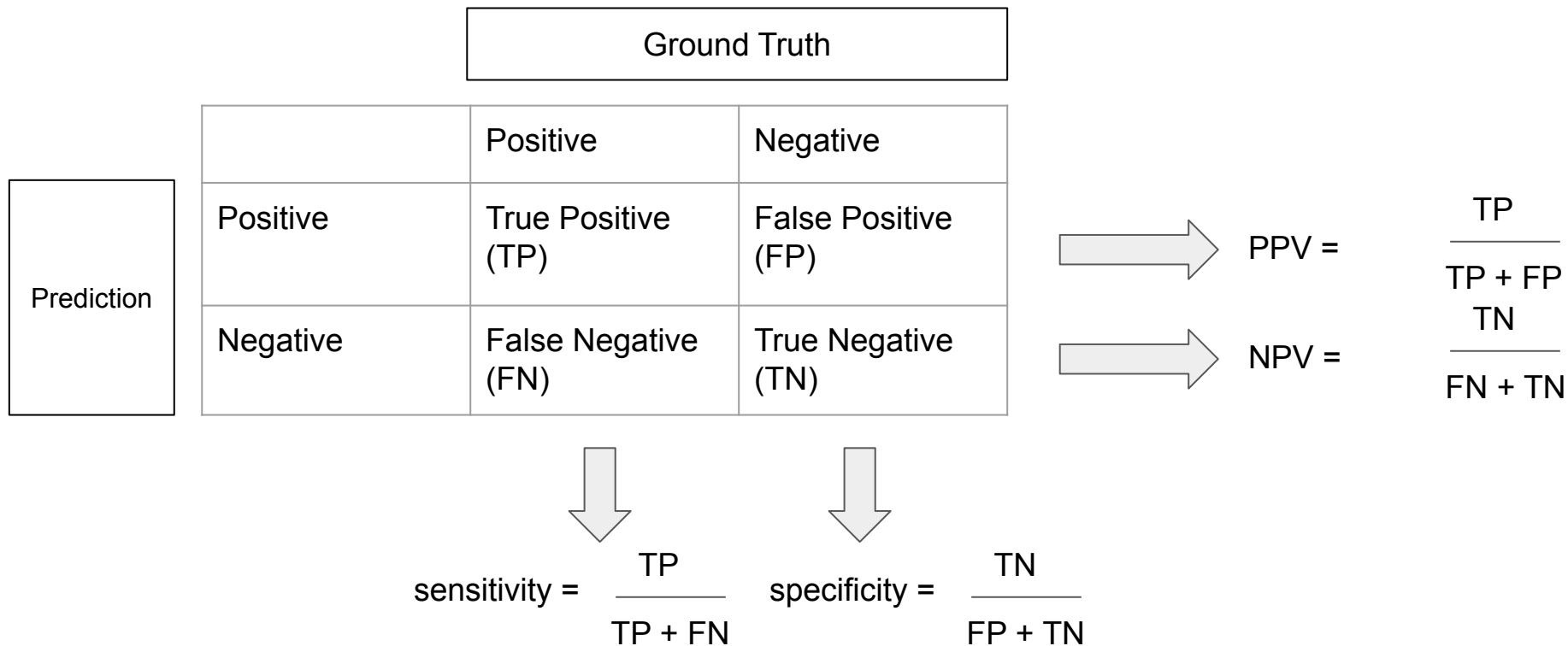
# 3. PPV and NPV



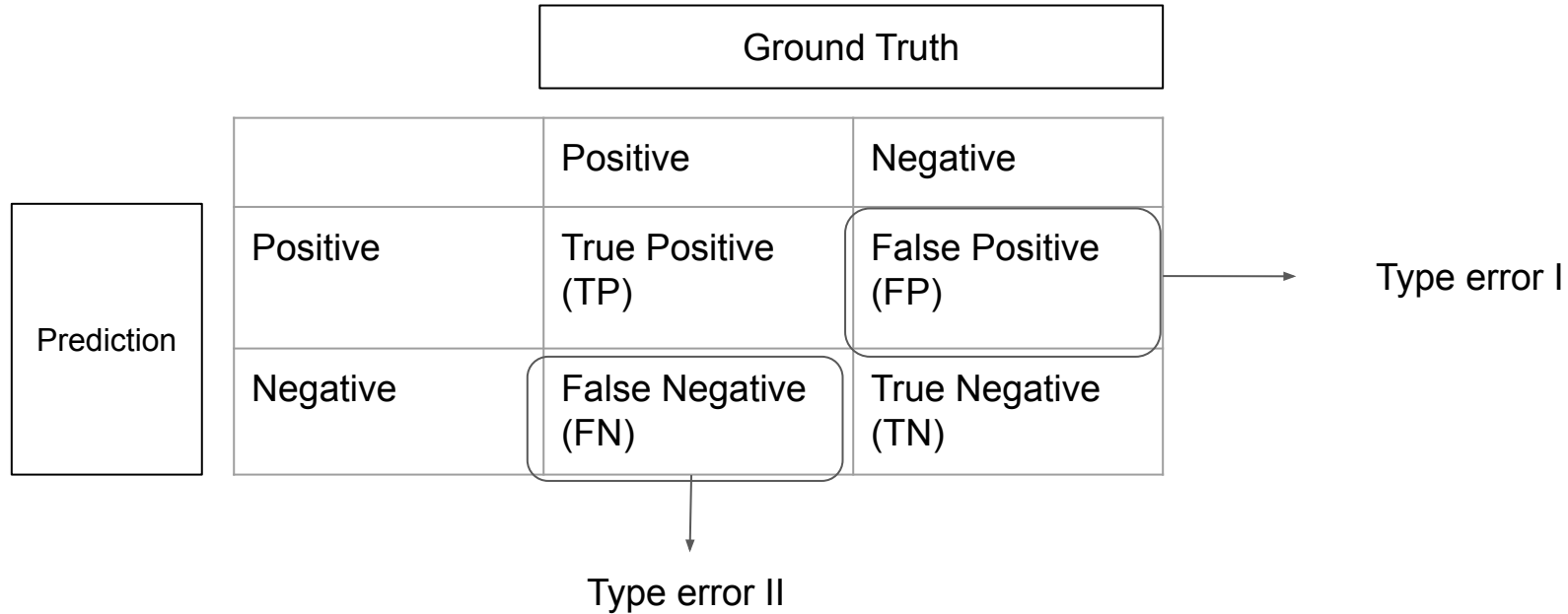
# 4. Confusion matrix



# 4. Confusion matrix



# 5. Type error I & II



What error is more serious?



# 6. f1-score, precision and recall

Model 1 and model 2 is binary classification model  
have the same accuracy is 80%.

**Positive:** *Disease*

**Negative:** *Normal*

**Model 1:**

Ground Truth			
Prediction		Positive	Negative
	Positive	0	0
	Negative	20	80

**Model 2:**

Ground Truth			
Prediction		Positive	Negative
	Positive	20	0
	Negative	20	60

What is the better model?

# 6. f1-score, precision and recall

Model 1 and model 2 is binary classification model  
have the same accuracy is 80%.

**Positive:** *Disease*

**Negative:** *Normal*

**Model 1:**

		Ground Truth	
		Positive	Negative
Prediction	Positive	0	0
	Negative	20	80

**Model 2:**

		Ground Truth	
		Positive	Negative
Prediction	Positive	20	0
	Negative	20	60

What is the better model?

⇒ Model 2 is better because it predict right 50% patients that get diseases.

# 6. f1-score, precision and recall

Model 1 and model 2 is binary classification model  
have the same accuracy is 80%.

**Positive:** *Disease*  
**Negative:** *Normal*

**Model 1:**

Ground Truth			
Prediction		Positive	Negative
	Positive	0	0
	Negative	20	80

**Model 2:**

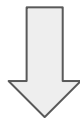
Ground Truth			
Prediction		Positive	Negative
	Positive	20	0
	Negative	20	60

What is the better model?

⇒ Model 2 is better because it predict right 50% patients that get diseases.  
Thus, we need to change accuracy metric to another metric.

# 6. f1-score, precision and recall

		Ground Truth	
		Positive	Negative
Prediction	Positive	True Positive (TP)	False Positive (FP)
	Negative	False Negative (FN)	True Negative (TN)



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



$$\text{Precision} = \frac{\text{TP}}{\text{TP} + \text{FP}}$$

$$\text{f1-score} = \frac{2}{1/\text{precision} + 1/\text{recall}}$$

# 6. f1-score, precision and recall

Model 1 and model 2 is binary classification model  
have the same accuracy is 80%.

**Positive:** *Disease*

**Negative:** *Normal*

**Model 1:**

		Ground Truth	
		Positive	Negative
Prediction	Positive	0	0
	Negative	20	80

Precision = ?

Recall = ?

f1-score = ?

**Model 2:**

		Ground Truth	
		Positive	Negative
Prediction	Positive	20	0
	Negative	20	60

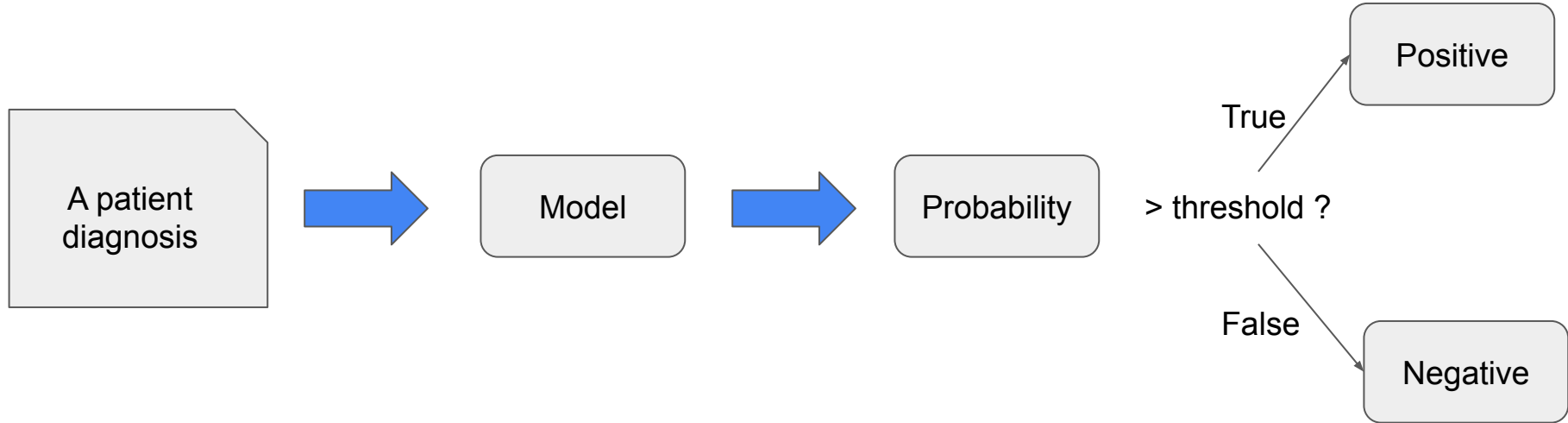
Precision = ?

Recall = ?

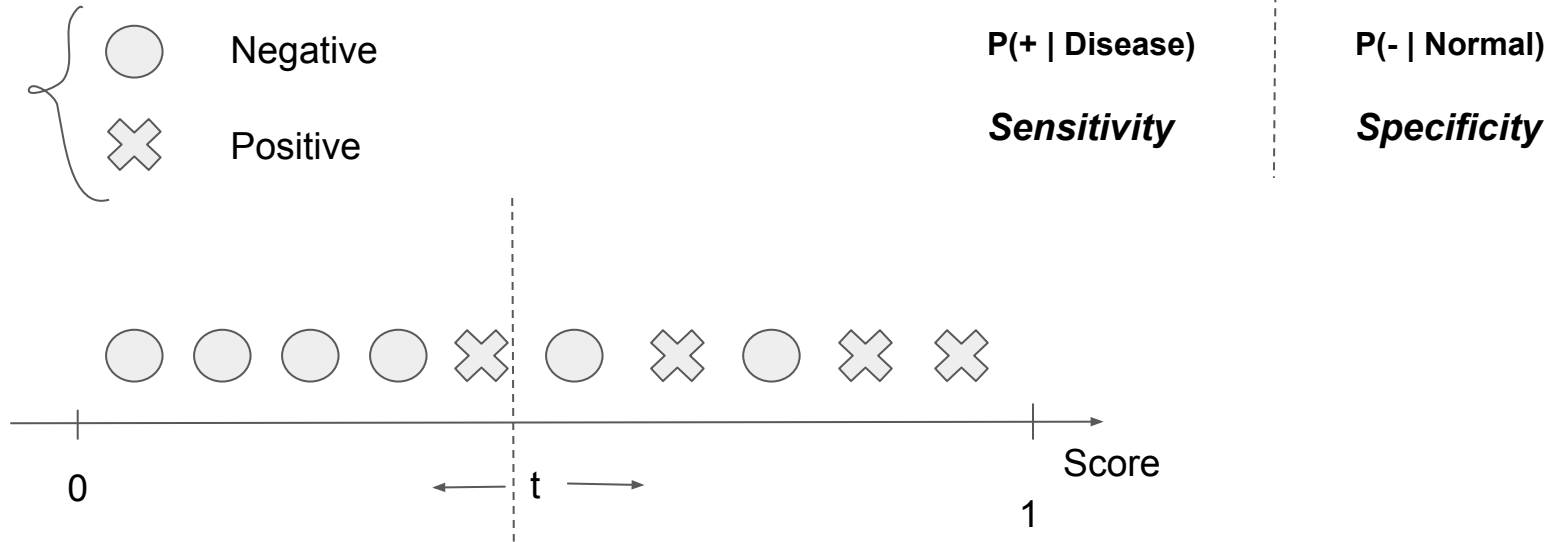
f1-score = ?

# 7. ROC curve and threshold

What is threshold ?

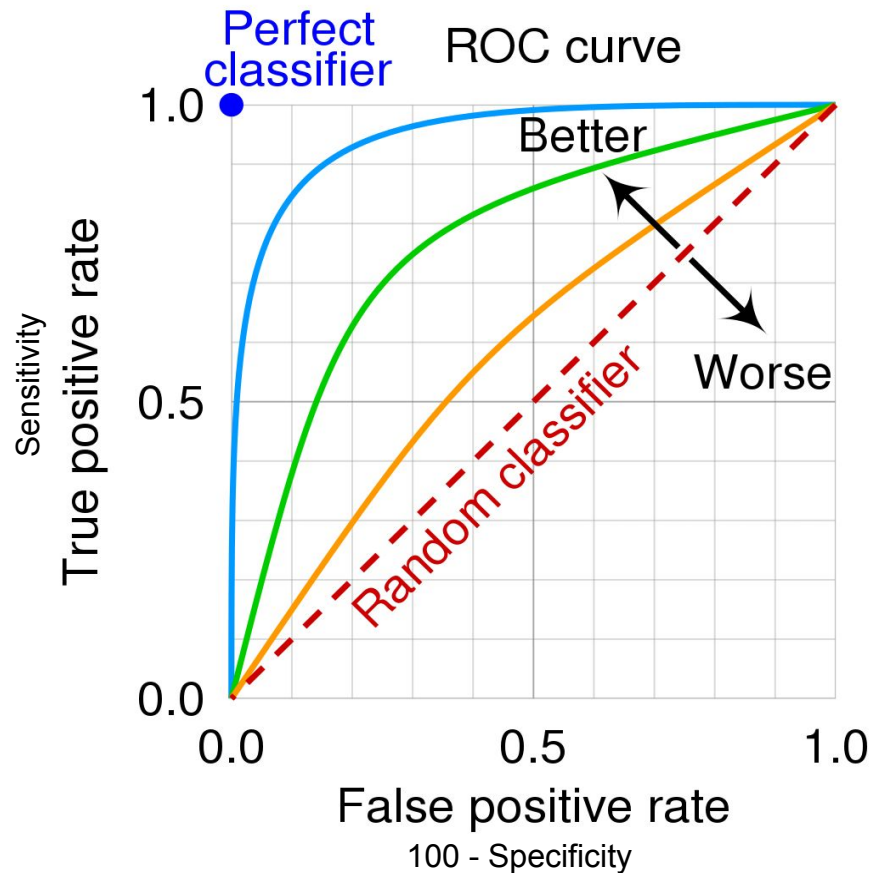


# 7. ROC curve and threshold



If we change threshold  $t$ , how does Sensitivity and Specificity change?

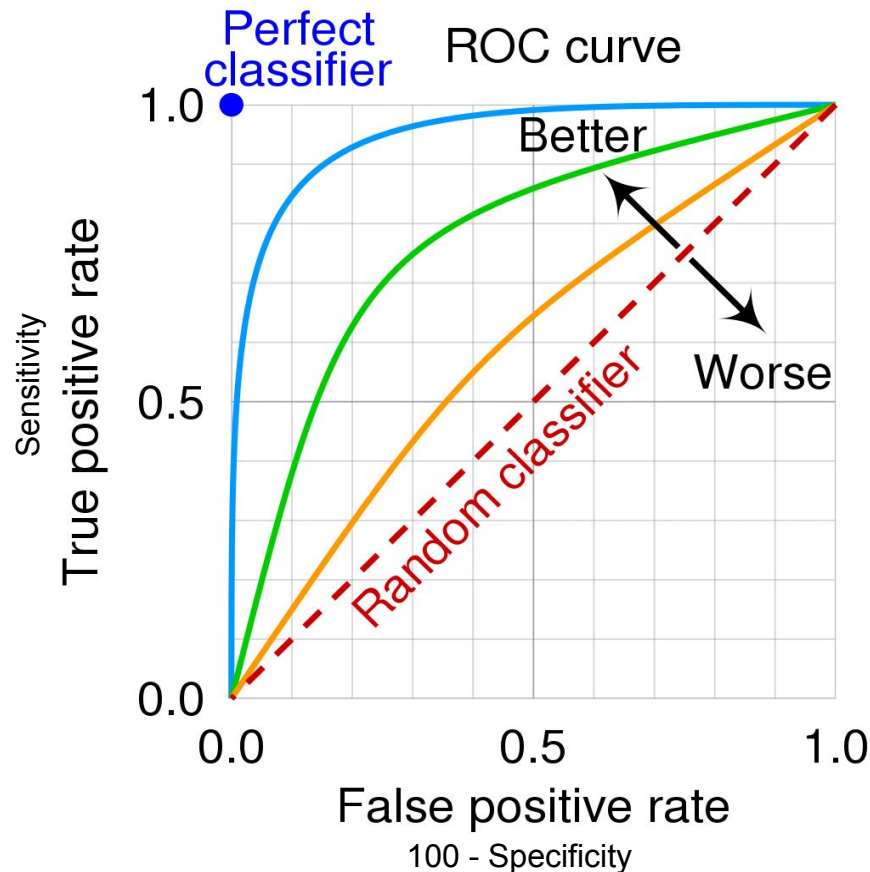
# 7. ROC and AUC





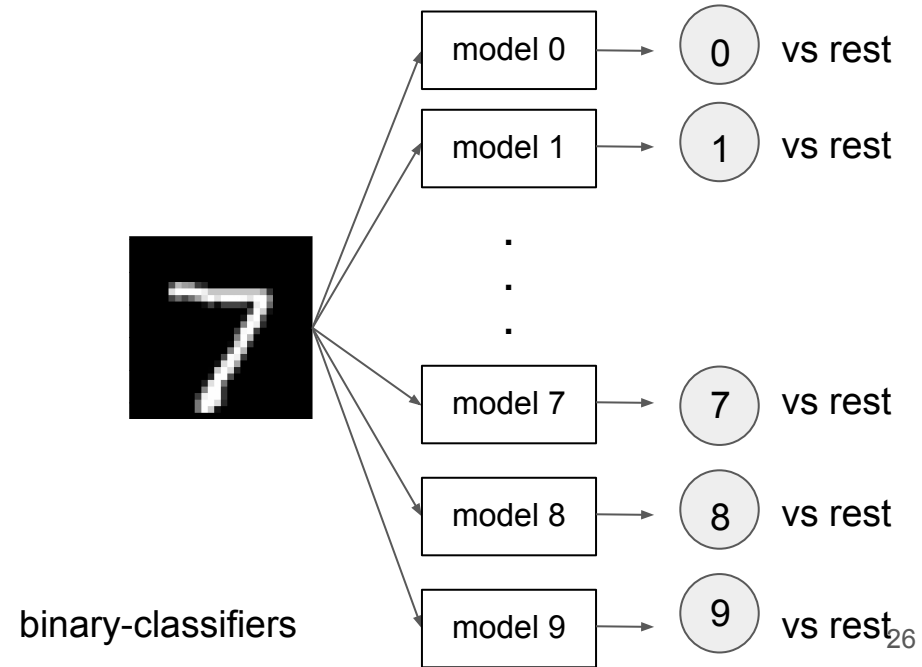
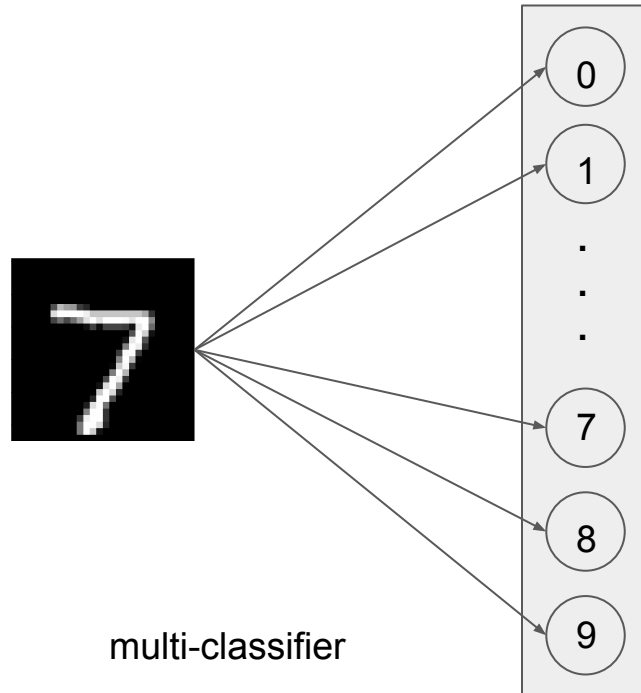
# 7. ROC and AUC

- *AUC = Area under the curve ROC*
- *AUC in  $[0, 1.0]$*
- *Bigger AUC is greater classifier*



# 8. Evaluate multi-classification model

- Multi-classifier literally is a group of binary classifier
- One-vs-rest method:



# 8. Evaluate multi-classification model

- Evaluate class  $i$  based on binary-classifier *model  $i$*

