



# Build It, Break It The Language Edition



Are you a **Builder**?

Do you think your approach to natural language processing problems can **withstand attacks**? Is it **robust** to unexpected inputs, or domain mismatch? Will a **rule-based** approach hold up better than a **learning-based** approach?

Participate in our shared tasks by solving an NLP problem—by rules or by learning—and learn how your approach stacks up against others on **adversarially, linguistically constructed** inputs.

Or maybe you're a **Breaker**? Click to flip...

[read shared task details](#)

[join our mailing list](#) for announcements

[join our workshop](#) at emnlp 2017

designed & implemented by



Emily M. Bender



Hal Daumé III



Allyson Ettinger

# **Sentiment analysis**

## **Logistic Regression**

---

# Logistic Regression (LR)

---

---

# Logistic Regression

---

a neural network can be viewed as a series of logistic regression classifiers stacked on top of each other.

---

# Generative vs. Discriminative

---

---

# NB is a Generative Model

---

$$P(H|E) = \frac{P(E|H)P(H)}{P(E)} = \frac{P(E|H)P(H)}{\sum_{H \in \mathcal{H}} P(E|H)P(H)}$$

If we choose  $H$ ,  $P(H)$ , how can we generate instances of event/evidences  $P(E|H)$   
If we are in  $H_4$ : Winter and in Canada, generate a day  $\rightarrow$  it would be mostly no sun!

---

LR is Discriminative!

---

---

# Logistic Regression

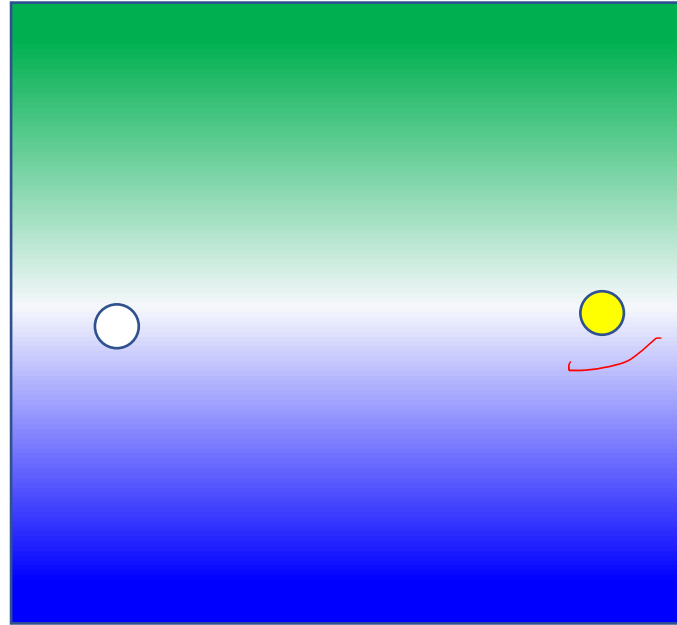
---

Don't care about hypothesis and their probabilities.  
Try to maximize the distance of data from different classes.  
Try to discriminate the most!



# Logistic Regression

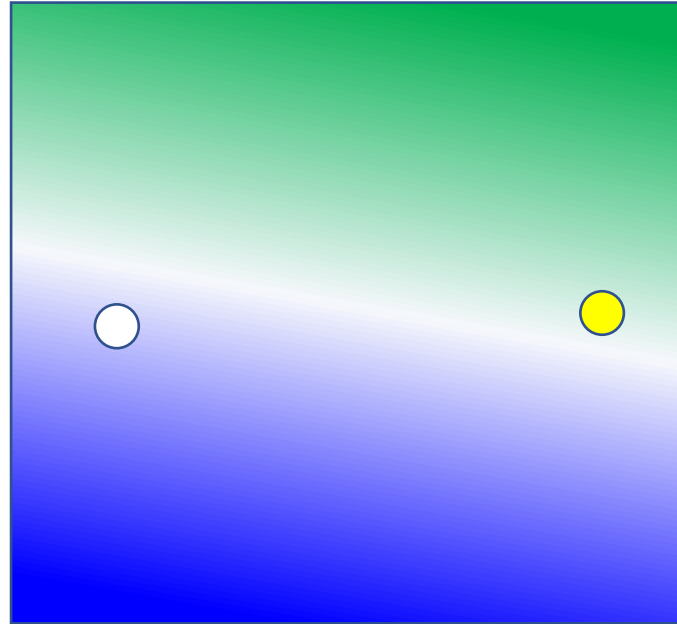
---



$$\left. \begin{array}{l} P(+|x_+) = 0.50 \quad P(-|x_+) = 0.50 \\ P(-|x_-) = 0.50 \quad P(+|x_-) = 0.50 \end{array} \right\}$$

# Logistic Regression

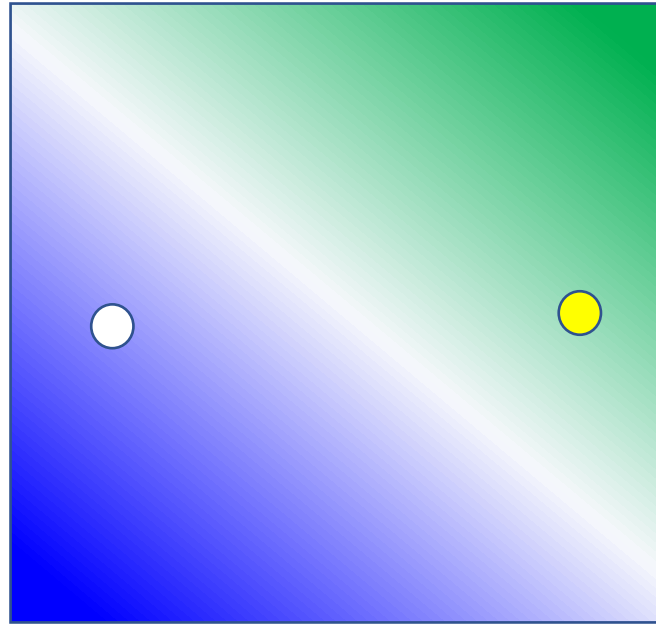
---



$$\begin{aligned} P(+|x_+) &= 0.55 & P(-|x_+) &= 0.45 \\ P(-|x_-) &= 0.55 & P(+|x_-) &= 0.45 \end{aligned}$$

# Logistic Regression

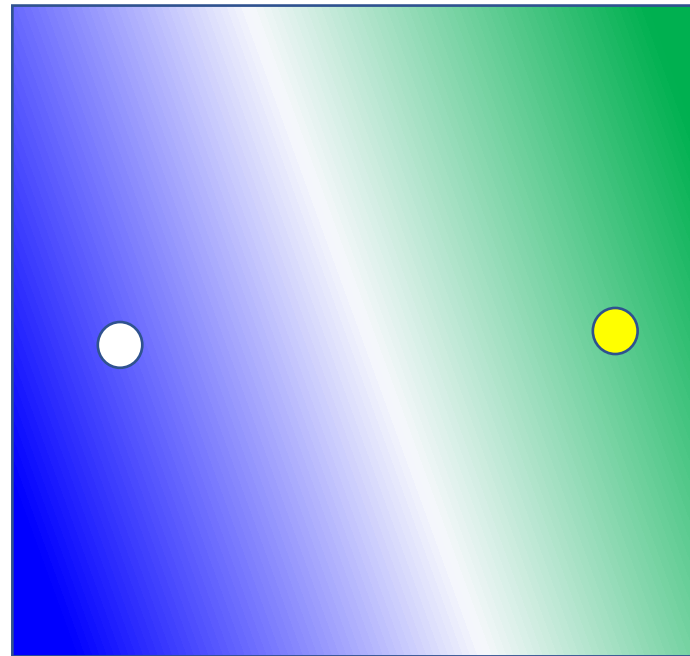
---



$$\begin{aligned} P(+|x_+) &= 0.65 & P(-|x_+) &= 0.35 \\ P(-|x_-) &= 0.65 & P(+|x_-) &= 0.35 \end{aligned}$$

# Logistic Regression

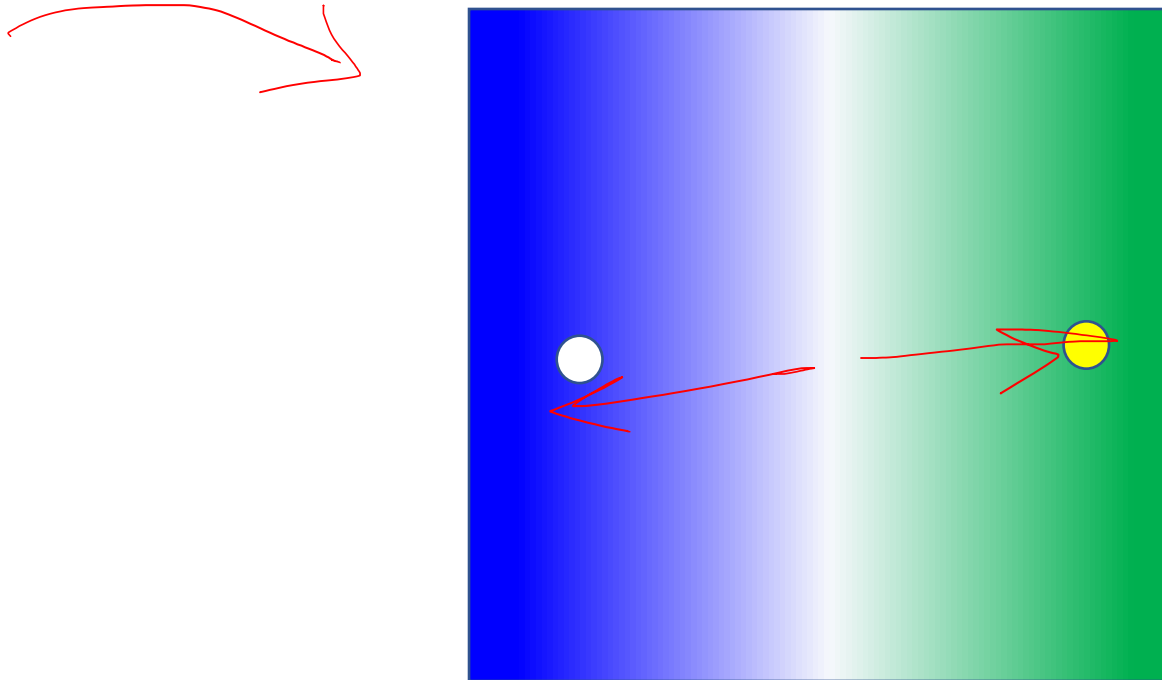
---



$$\begin{aligned} P(+|x_+) &= 0.75 & P(-|x_+) &= 0.25 \\ P(-|x_-) &= 0.75 & P(+|x_-) &= 0.25 \end{aligned}$$

# Logistic Regression

---



$$\begin{aligned}P(+|x_+) &= 0.85 & P(-|x_+) &= 0.15 \\P(-|x_-) &= 0.85 & P(+|x_-) &= 0.15\end{aligned}$$

---

# Logistic Regression

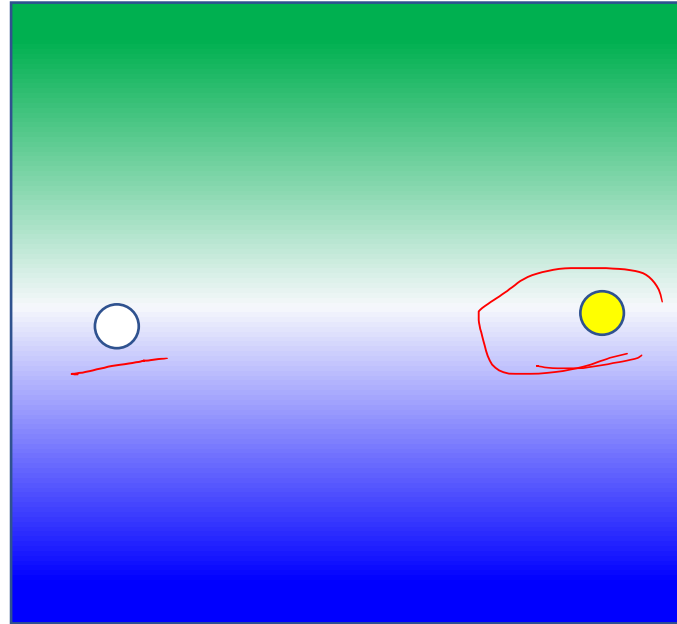
---

- (I) Iterative Process (online vs. offline)
- (II) Optimization = Discriminate classes the most

# Logistic Regression

---

More for  $P(-|x_-) \uparrow$

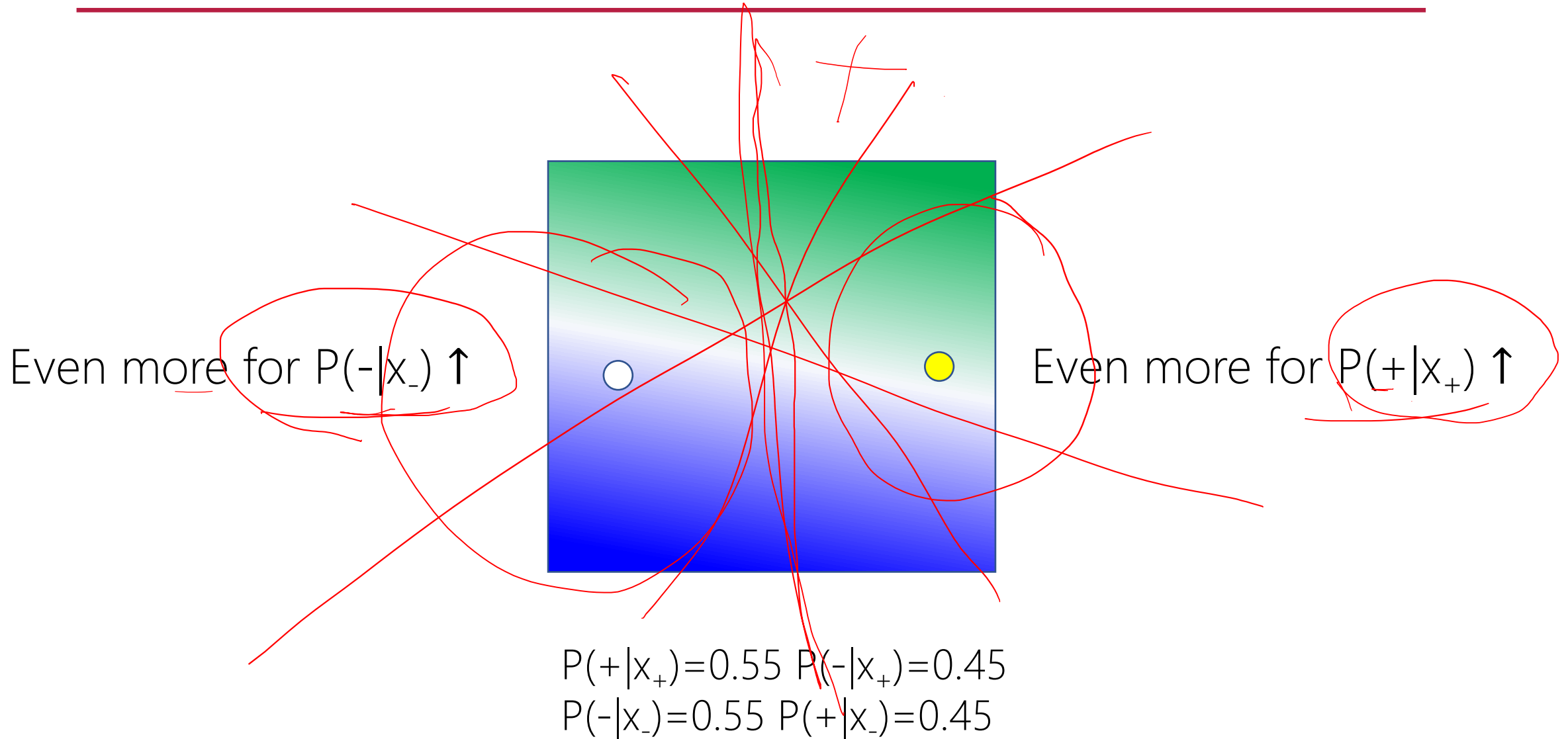


More for  $P(+|x_+) \uparrow$

$$\begin{aligned} P(+|x_+) &= 0.50 & P(-|x_+) &= 0.50 \\ P(-|x_-) &= 0.50 & P(+|x_-) &= 0.50 \end{aligned}$$

# Logistic Regression

---

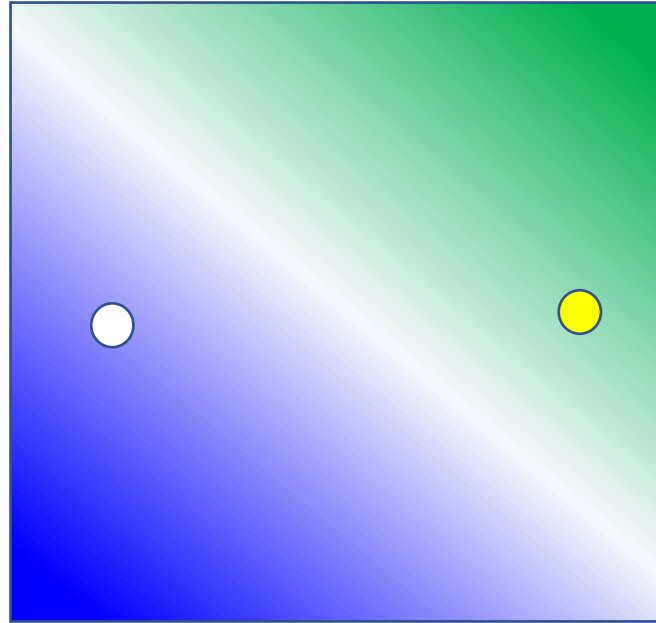




# Logistic Regression

---

Even more for  $P(-|x_-) \uparrow$



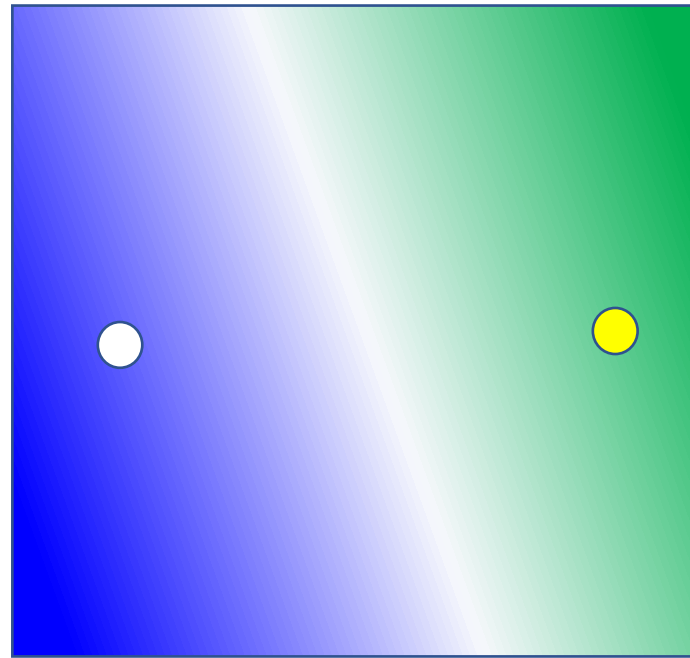
Even more for  $P(+|x_+) \uparrow$

$$\begin{aligned} P(+|x_+) &= 0.65 & P(-|x_+) &= 0.35 \\ P(-|x_-) &= 0.65 & P(+|x_-) &= 0.35 \end{aligned}$$

# Logistic Regression

---

Even more for  $P(-|x_-) \uparrow$

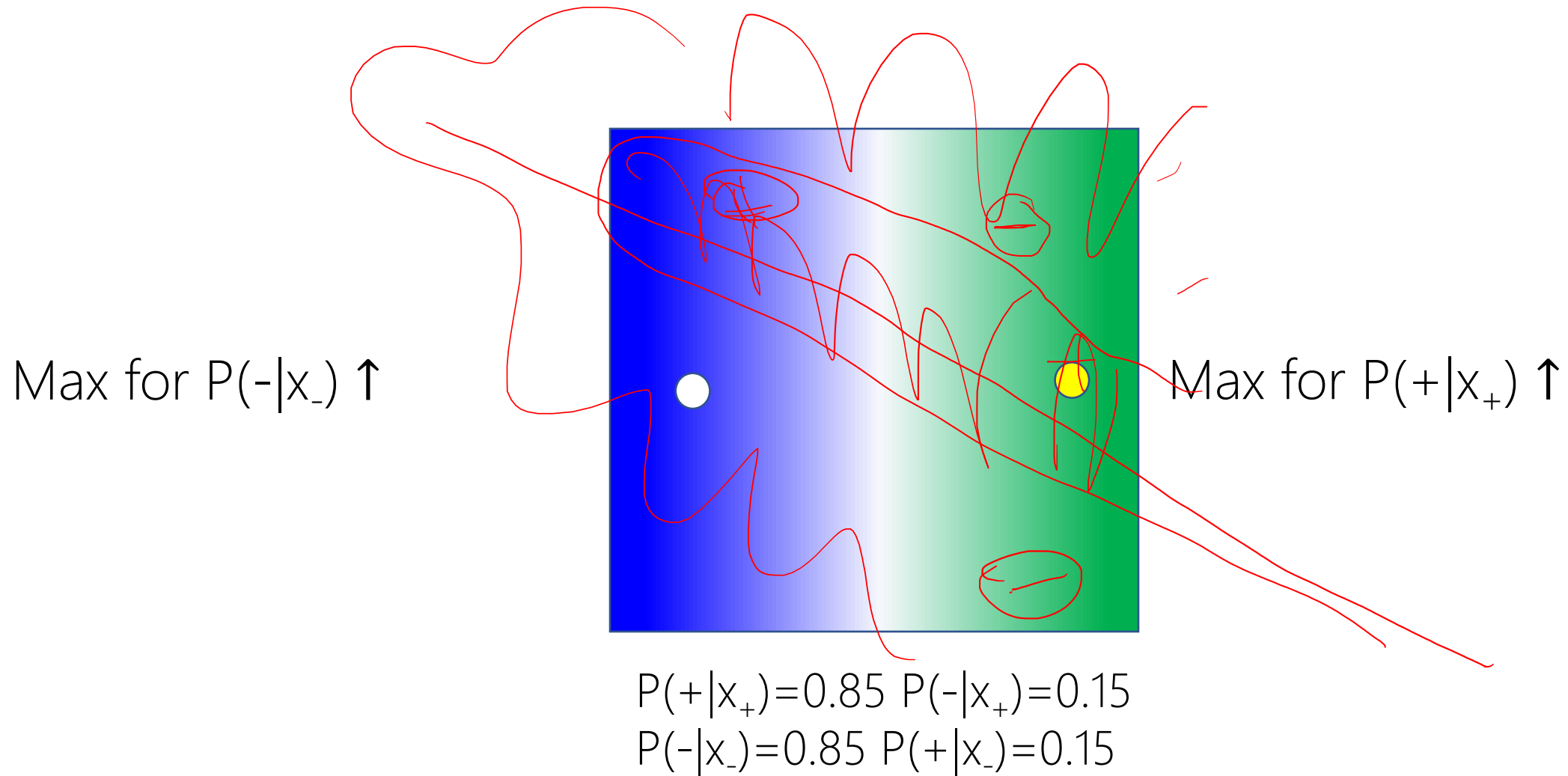


Even more for  $P(+|x_+) \uparrow$

$$\begin{aligned} P(+|x_+) &= 0.75 & P(-|x_+) &= 0.25 \\ P(-|x_-) &= 0.75 & P(+|x_-) &= 0.25 \end{aligned}$$

# Logistic Regression

---



---

# Logistic Regression

---

- (I) Iterative Process: We can change the **line** that discriminate
- (II) Optimization = Discriminate classes the most
- (I) **Maximizing** the  $P(+|x_+) + P(-|x_-)$

---

# Logistic Regression

---

- (I) What is that line?
- (II) What if we have more than two inputs?

$$\begin{aligned}
 &= P(A)P(B) \dots \\
 &P(A, B, C, D \dots) \\
 \hline
 &P(A, B) = P(A)P(B|A) \quad P(A)P(B|\cancel{A})P(C|A, B) \dots \\
 &= P(A)P(B) \quad \text{independence} \\
 \hline
 &\Rightarrow \text{Assumption} \Rightarrow \text{train}
 \end{aligned}$$

# Logistic Regression

(I) Iterative Process: We can change the line that discriminate

(II) Optimization

(I)  $\text{Max} \left( \prod_{x \in \{+\}} P(+|x_+) \prod_{x \in \{-\}} P(-|x_-) \right)$

# Logistic Regression

$$M \Rightarrow \begin{aligned} M(X_1^+) &= \log(0.3) + -2000 \\ M(X_2^+) &= \log(0.8) - 10 \\ M(X_3^+) &= \log(0.5) - 1000 \\ M(X_4^-) &= \log(0.2) - 3000 \end{aligned}$$

$$= 56$$

(I) Iterative Process: We can change the line that discriminate

(II) Optimization

$$(I) \text{ Max } (\sum_{x \in \{+\}} \text{Log } P(+|x_+) + \sum_{x \in \{-\}} \text{Log } P(-|x_-))$$

---

# Logistic Regression

---

(I) Iterative Process: We can change the line that discriminate

(II) Optimization

(I) **Min**  $-(\sum_{x \in \{+\}} \text{Log } P(+|x_+) + \sum_{x \in \{-\}} \text{Log } P(-|x_-))$



---

# Logistic Regression

$$c = \{-, +\} \rightarrow y = \{0, 1\}$$

Handwritten red annotations: The minus sign is circled, the plus sign is boxed, and the mapping to  $y = \{0, 1\}$  is indicated by arrows. To the right, the numbers 2, 3 are handwritten in red.

(I) Iterative Process: We can change the line that discriminate

(II) Optimization

(I) Min  $-\sum_{(x,y) \in D} [(y) \text{Log } P(y|x_y) + (1-y) \text{Log } P(y|x_y)]$

Handwritten red annotations: The word "Min" is highlighted in yellow. The entire expression is crossed out with a large red 'X'. Above the first term,  $P(1|x_1)$  is written. Above the second term,  $P(0|x_0)$  is written. Below the first term,  $1 \times \log(1/x_1)$  is written. Below the second term,  $1 \times \log(1/x_0)$  is written.

---

# Logistic Regression

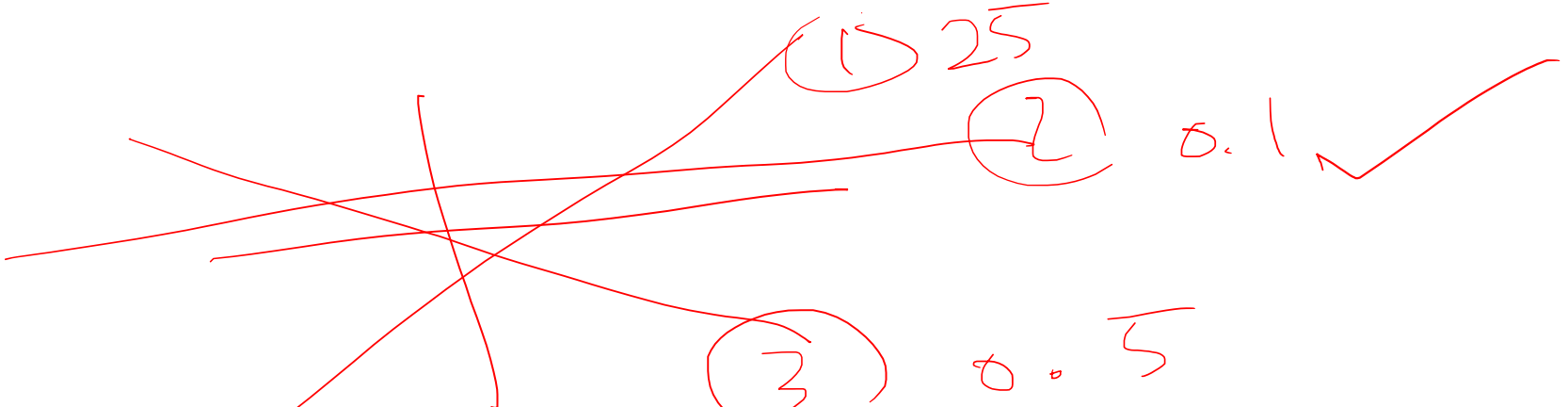
$c = \{-, +\} \rightarrow y = \{0, 1\}$

---

(I) Iterative Process: We can change the line that discriminate

(II) Optimization

(I)  $\text{Min} - \sum_{(x,y) \in D} \text{Log } P(y|x_y)$   $\rightarrow 0$



# Logistic Regression

$C = \{-, +\} \rightarrow y = \{0, 1\}$

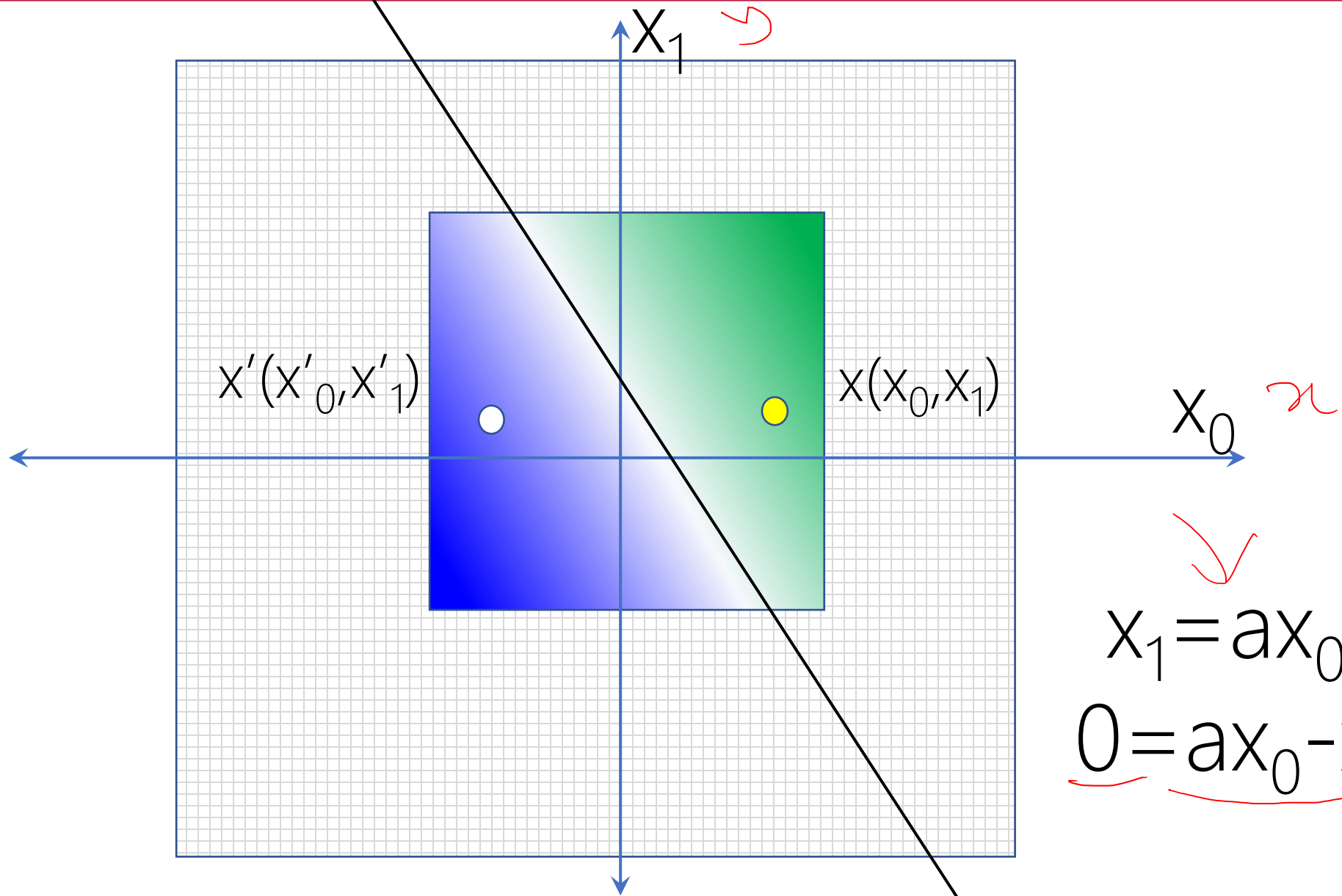
---

(I) Iterative Process: We can change the **line** that discriminate

(II) Optimization

(I)  $\text{Min} - \sum_{(x,y) \in D} [(y) \text{Log } P(+|x_+) + (1-y) \text{Log } P(-|x_-)]$

# Logistic Regression



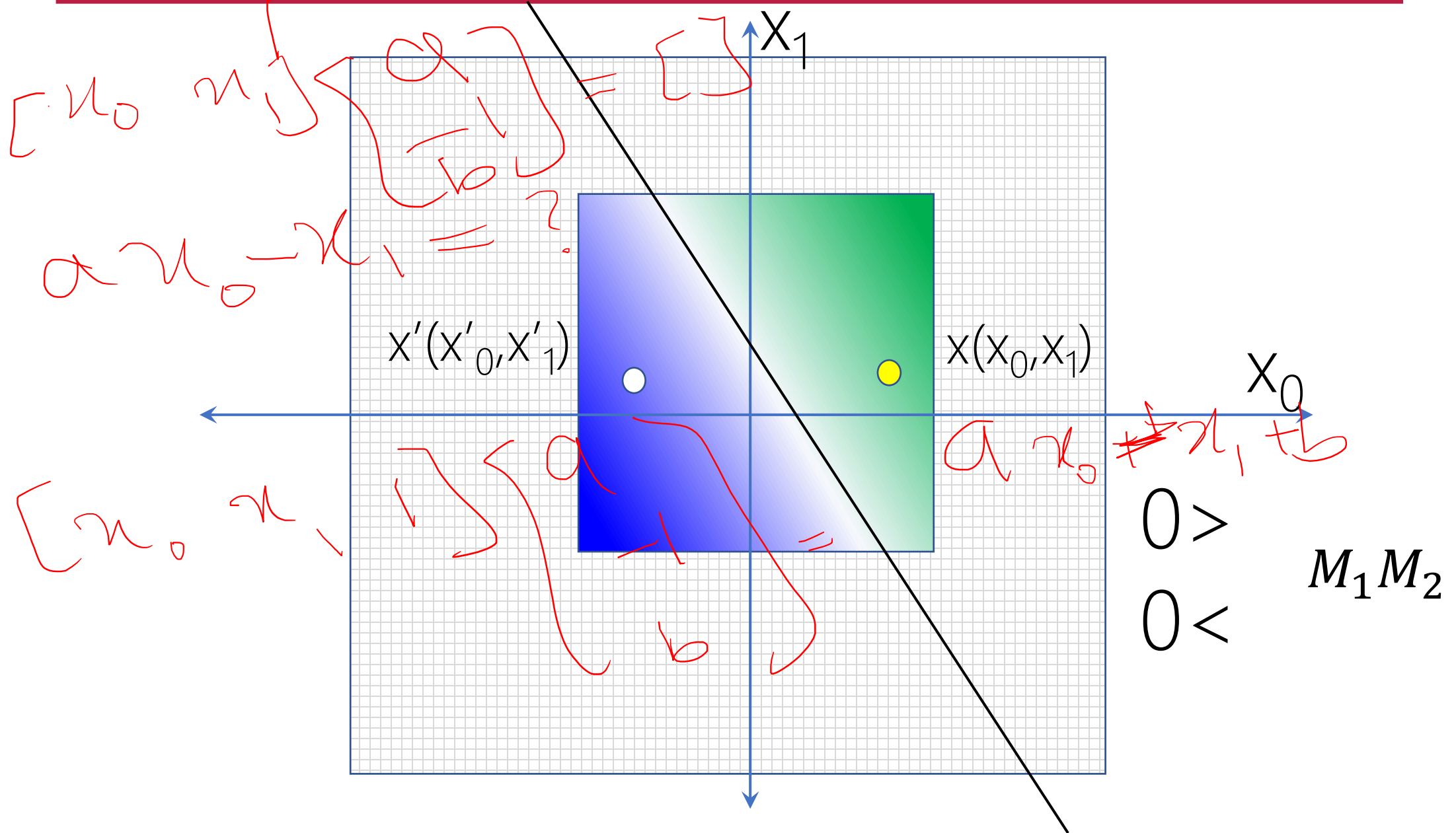
$$x_1 = ax_0 + b$$

$$0 = ax_0 - x_1 - b$$

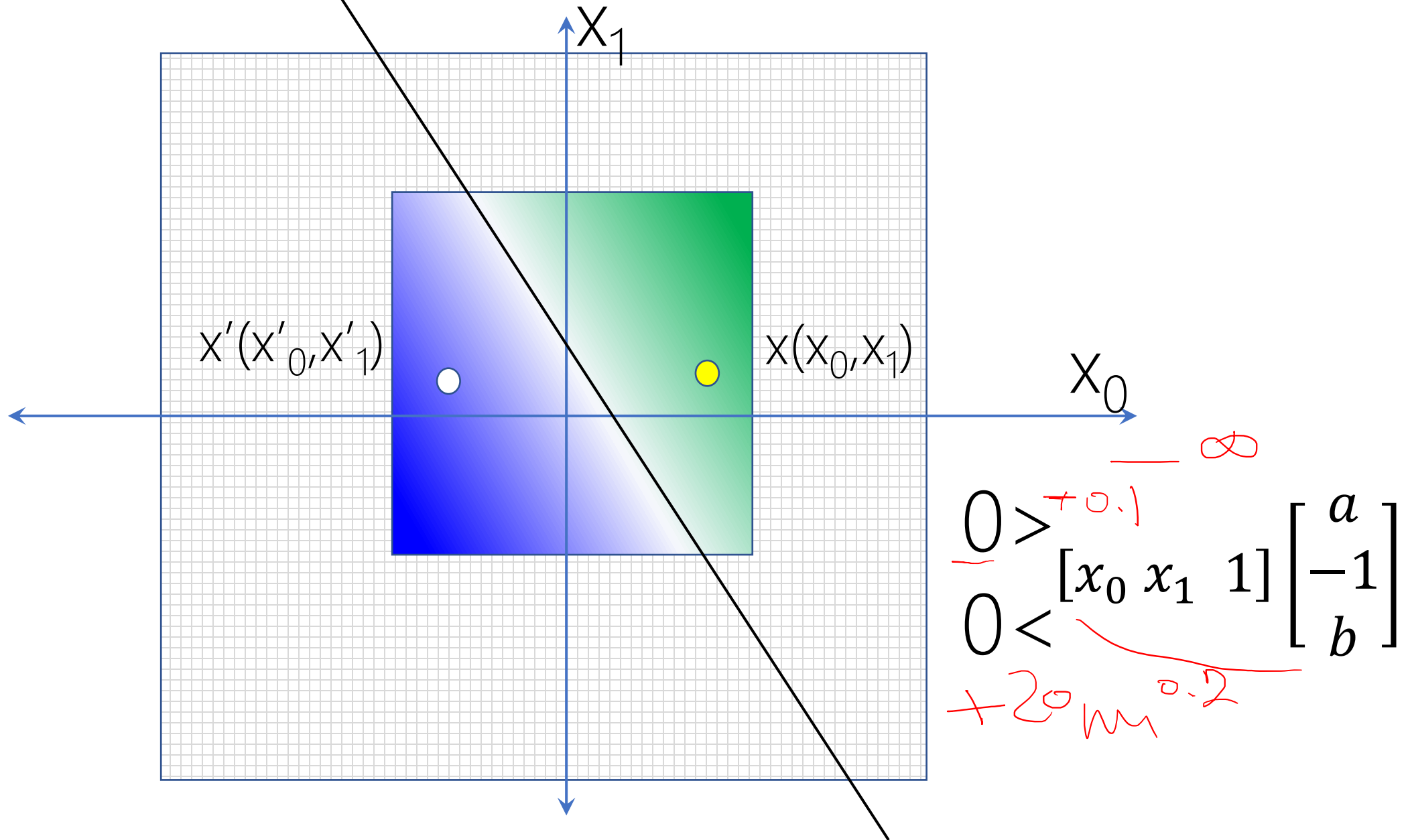
A 2D plot on a grid showing a blue square region and a green square region. The blue region is labeled  $x'(x'_0, x'_1)$  and contains a white dot. The green region is labeled  $x(x_0, x_1)$  and contains a yellow dot. A black line passes through the regions. Red handwritten annotations include  $x$ ,  $x_0$ ,  $x_1$ , and  $x$ .

$$x_1 = ax_0 + b$$
$$0 > ax_0 - x_1 + b$$
$$0 < ax_0 - x_1 + b$$

# Logistic Regression



# Logistic Regression



# Logistic Regression

---

$0 >$

$0 <$

$$\begin{bmatrix} x_0 & x_1 & \dots & x_d & 1 \end{bmatrix} \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_{d+1} \end{bmatrix} = [?]$$

$x_0 \ x_1 \ x_2$

$x_0 \ x_1 \ \dots \ x_n$

$(n)$

$d = 1 \rightarrow$  Line in 2-dimension

$d = 2 \rightarrow$  Plane in 3 dimension

$d = n \rightarrow$  Hyperplane in  $(n+1)$  dimension



---

# Parametric vs. non-Parametric

---

LR vs. Naïve Bayes

(True Bayesian Inference is Parametric, Why?)

# Logistic Regression

---

$$0 > f(X) \rightarrow -\infty$$

$$0 < f(X) \rightarrow +\infty$$

$d = 1 \rightarrow$  Line in 2-dimension

$d = 2 \rightarrow$  Plane in 3 dimension

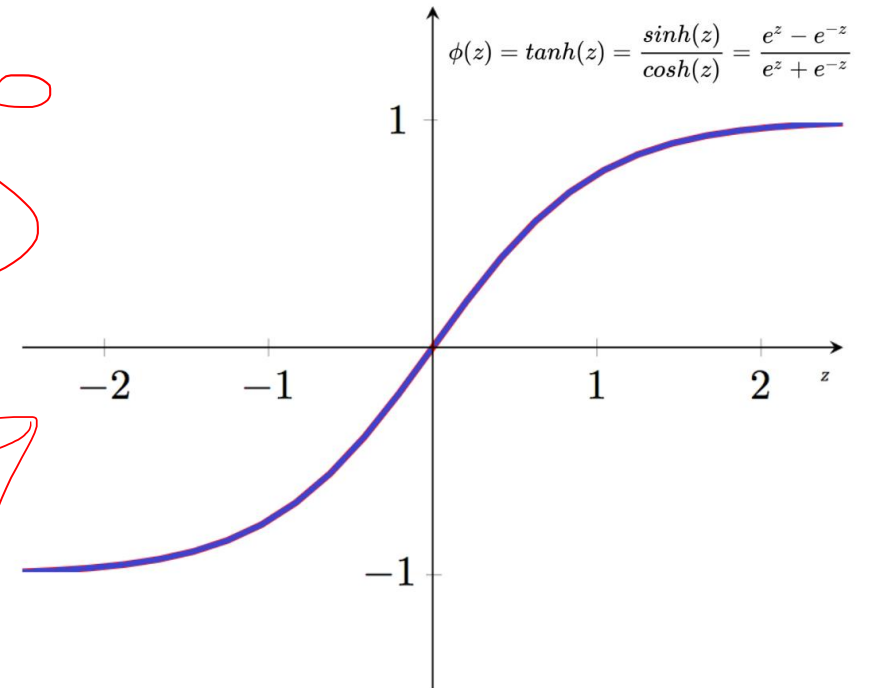
$d = n \rightarrow$  Hyperplane in  $(n+1)$  dimension

# Logistic Regression: Squish by Tanh

$[-1, +1]$

$$0 > f(X) \rightarrow -1$$

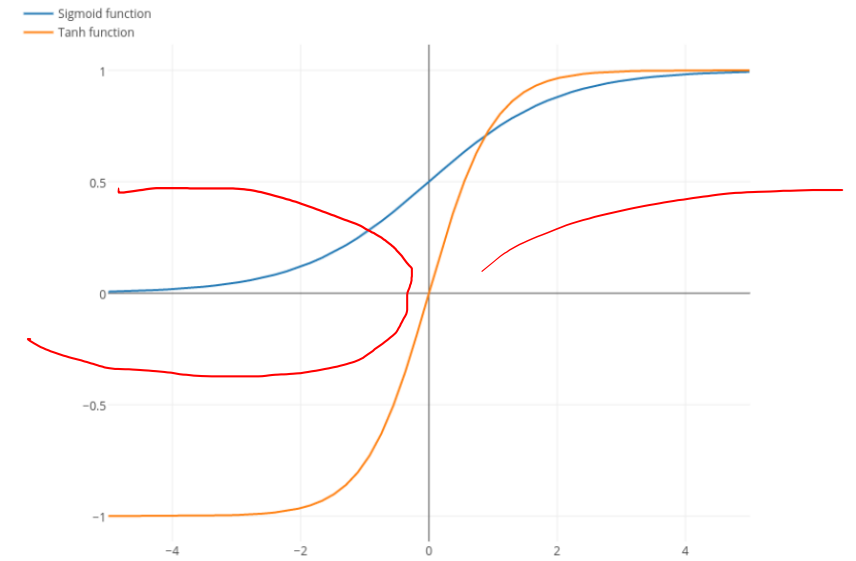
$$0 < f(X) \rightarrow +1$$



# Logistic Regression: Squish by Sigmoid

$$f(X) > 0 \rightarrow 0$$

$$f(X) < 0 \rightarrow +1$$



$$s(f(X)) = \frac{1}{1+e^{-f(X)}}$$

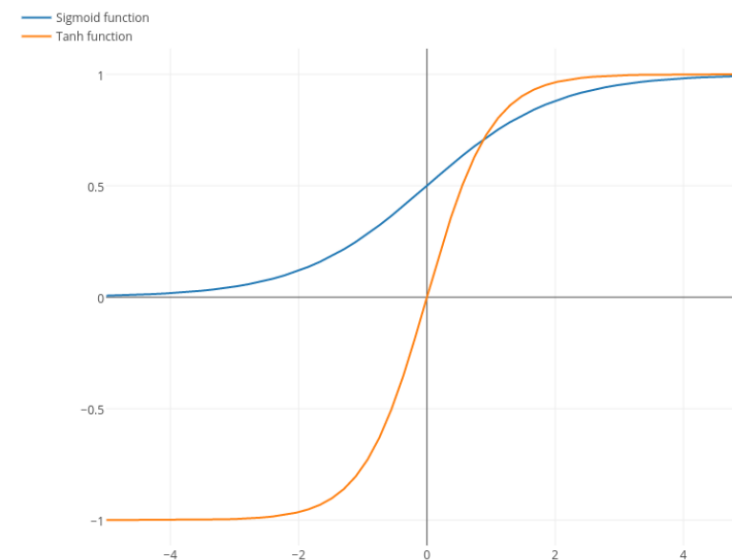
Becomes very similar to probability values!

# Logistic Regression: Squish

---

$$f(X) > 0 \rightarrow 0$$

$$f(X) < 0 \rightarrow +1$$



$$s(f(X)) = \frac{1}{1+e^{-f(X)}}$$

Becomes very similar to probability values!

But only for positive class (+)

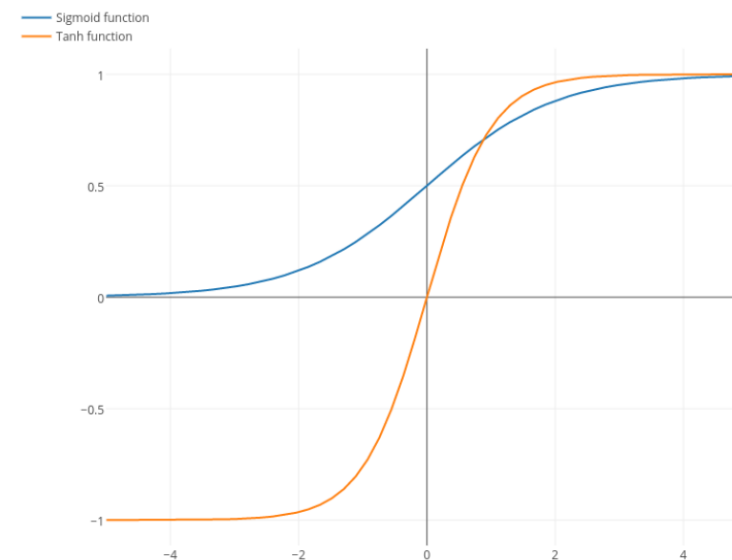
$$y=1 \rightarrow P(y|x) = P(+|x) = \text{Sigmoid}(x)$$

# Logistic Regression: Squish

---

$$f(X) > 0 \rightarrow 1$$

$$f(X) < 0 \rightarrow 0$$



$$s(f(X)) = \frac{1}{1+e^{-f(X)}}$$

Becomes very similar to probability values!

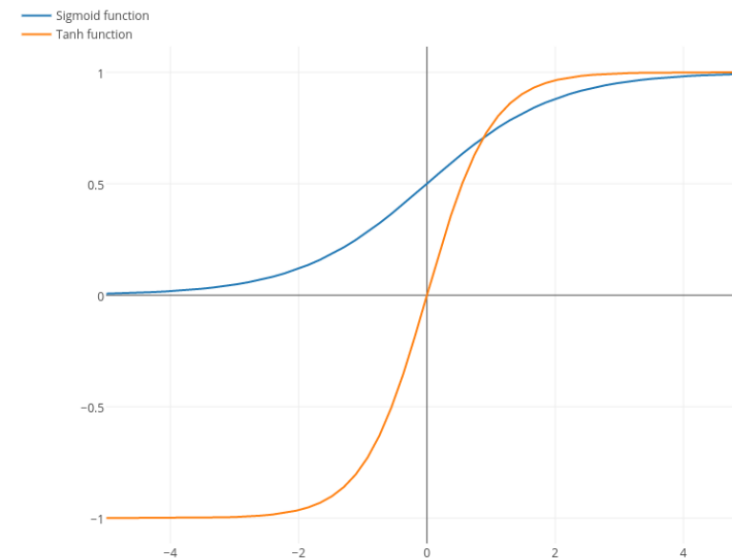
For negative class (-)?

# Logistic Regression: Squish

---

$$f(X) > 0 \rightarrow 1$$

$$f(X) < 0 \rightarrow 0$$



$$s(f(X)) = \frac{1}{1+e^{-f(X)}}$$

Becomes very similar to probability values!

For negative class (-)

$$P(+|x) + P(-|x) = 1 \rightarrow P(-|x) = 1 - P(+|x)$$

# Logistic Regression

$C = \{-, +\} \rightarrow y = \{0, 1\}$

- (I) Iterative Process: We can change the  $f$  that discriminate
- (II) Optimization
- (I)  $\text{Min} - \sum_{(x,y) \in D} \text{Log } P(y|x_y) \rightarrow \begin{cases} y = 1: P(+|x_+) = \text{Sigmoid}(f(x_+)) \\ y = 0: P(-|x_-) = 1 - \text{Sigmoid}(f(x_+)) \end{cases}$



---

# Logistic Regression

$$P \sim f$$

---

(I) Iterative Process: We can change the  $f$  that discriminate

(II) Optimization

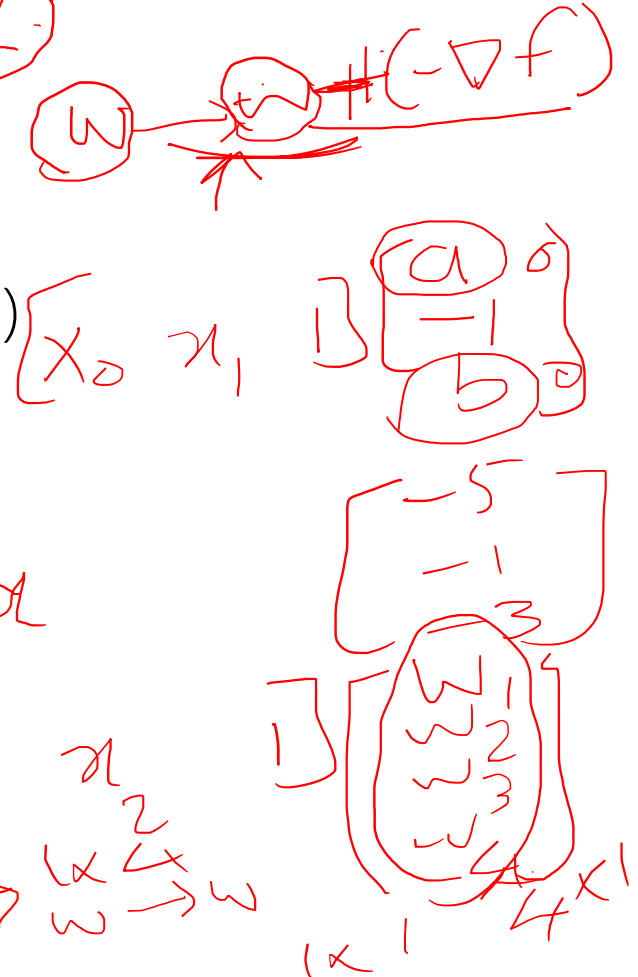
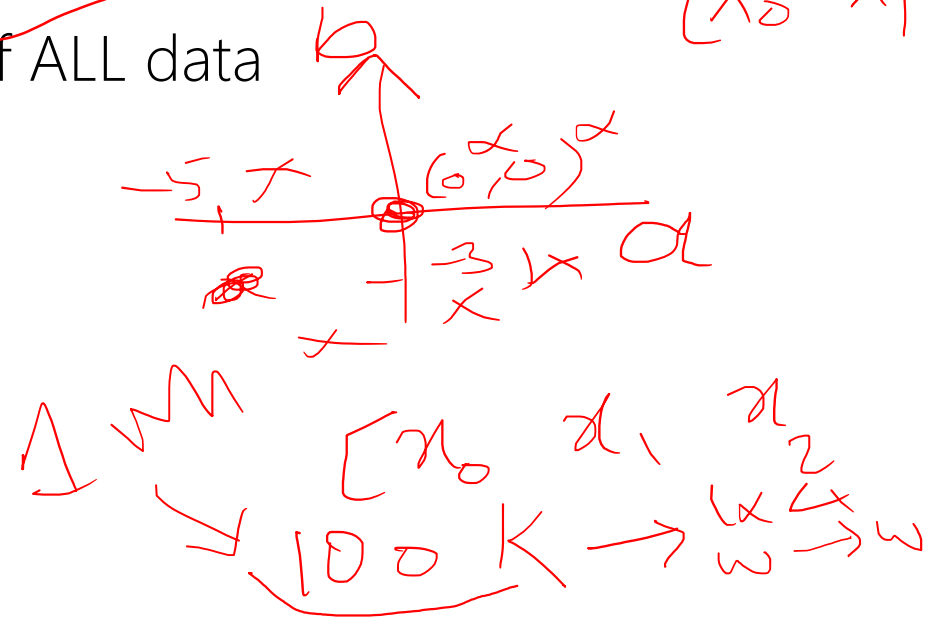
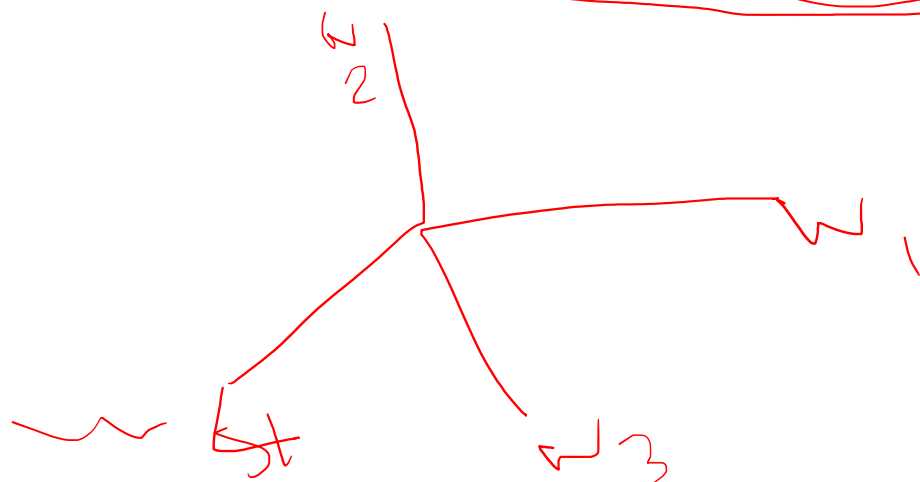
$$(I) \quad \text{Min} -\sum_{(x,y) \in D} \text{Log } P(y|x_y)] \rightarrow \begin{cases} y = 1: P(+|x_+) = \text{Sigmoid}(f(x_+)) \\ y = 0: P(-|x_-) = 1 - \text{Sigmoid}(f(x_+)) \end{cases}$$

# Logistic Regression

Optimization:  $\text{Min } -\sum_{(x,y) \in D} \text{Log } P(y|x_y) \rightarrow \begin{cases} y = 1: P(+|x_+) = \text{Sigmoid}(f(x_+)) \\ y = 0: P(-|x_-) = 1 - \text{Sigmoid}(f(x_+)) \end{cases}$

## Derivatives per function weights (Gradients)

- Update on each input data
- Update on batches of input data (Stochastic Gradient Descent)
- Update on multiple rounds (epoch) of ALL data

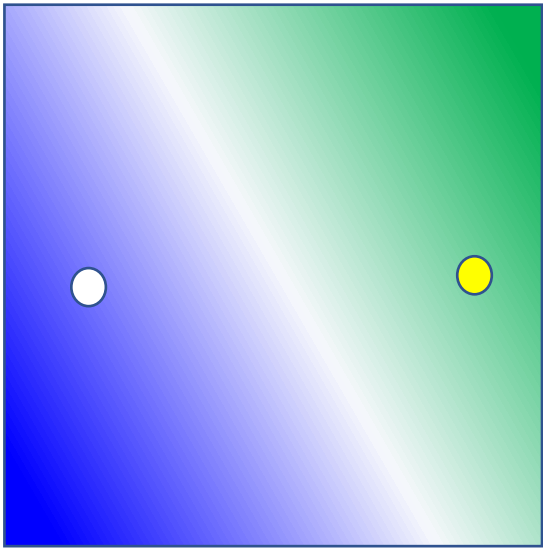


# Logistic Regression

---

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$   $\rightarrow \begin{cases} y = 1: P(+|x_+) = \text{Sigmoid}(f(x_+)) \\ y = 0: P(-|x_-) = 1 - \text{Sigmoid}(f(x_+)) \end{cases}$

- Function  $f$  is linear function of weights



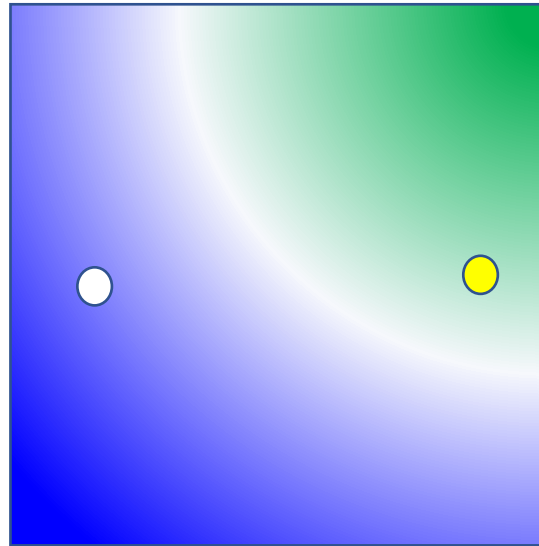
$$[x_0 \ x_1 \ \dots \ x_d \ 1] \begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_{d+1} \end{bmatrix} = [?]$$

# Logistic Regression

---

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)] \rightarrow \begin{cases} y = 1: P(+|x_+) = \text{Sigmoid}(f(x_+)) \\ y = 0: P(-|x_-) = 1 - \text{Sigmoid}(f(x_+)) \end{cases}$

- Can we have  $f$  as a non-linear function of weights?



---

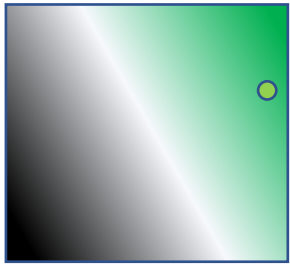
# Multiclass with Logistic Regression

---

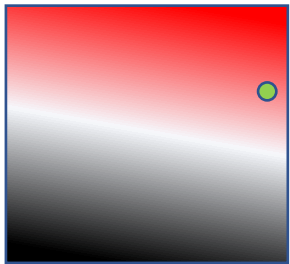
# Multi-target (Multi-label)

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$

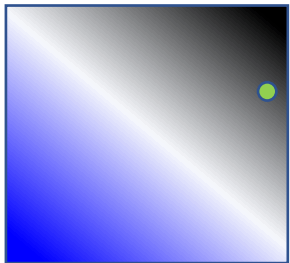
- Data point can be member of more than one class



green class vs. else:  $[x_0 \ x_1 \ \dots \ x_d \ 1]$   $\begin{bmatrix} w_1 \\ w_2 \\ \dots \\ w_{d+1} \end{bmatrix} = [?]$



red class vs. else:  $[x_0 \ x_1 \ \dots \ x_d \ 1]$   $\begin{bmatrix} w'_1 \\ w'_2 \\ \dots \\ w'_{d+1} \end{bmatrix} = [?]$



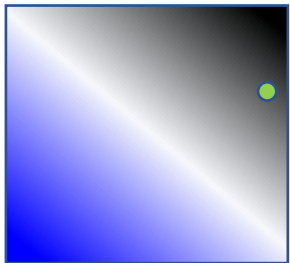
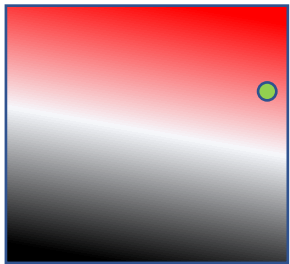
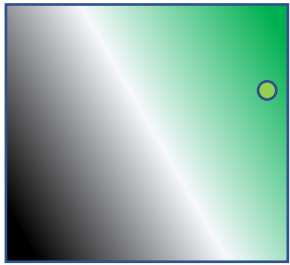
blue class vs. else:  $[x_0 \ x_1 \ \dots \ x_d \ 1]$   $\begin{bmatrix} w''_1 \\ w''_2 \\ \dots \\ w''_{d+1} \end{bmatrix} = [?]$

Multiple  
Binary  
Classification

# Multi-target (Multi-label)

Optimization:  $\text{Min } -\sum_{(x,y) \in D} \text{Log } P(y|x_y)$

- Data point can be member of more than one class



$$\begin{matrix} \text{green} \\ \text{no + red} \\ \text{not blue} \end{matrix} \begin{bmatrix} x_0 & x_1 & \dots & x_d & 1 \end{bmatrix} \begin{bmatrix} w_1 & w'_1 & w''_1 \\ w_2 & w'_2 & w''_2 \\ \dots & \dots & \dots \\ w_{d+1} & w'_{d+1} & w''_{d+1} \end{bmatrix} = [\text{?}, \text{?}, \text{?}] = [1, 0, 1]$$

[ ] [ ] [ ]  
Si -

$([0, 1], [0, 1], [0, 1])$

---

## Multinomial Logistic Regression

---

Logistic Regression  $\rightarrow$  Softmax Regression



# Softmax Regression

$$\text{softmax}(\mathbf{z} = [z_1, z_2, \dots, z_d]) = \frac{1}{\sum_{i=1}^d e^{z_i}} [e^{z_1}, e^{z_2}, \dots, e^{z_d}]$$

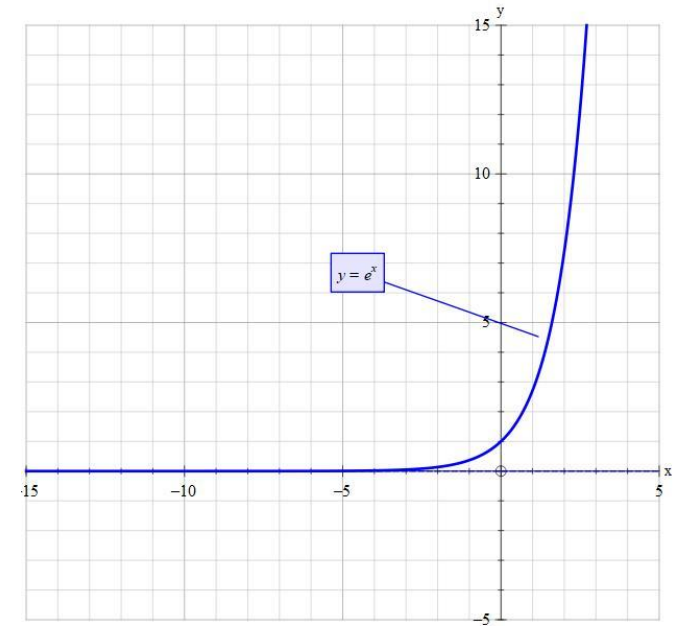
Handwritten notes:  $e^z \Rightarrow \frac{e^z}{\sum e^i}$

$$0 \leq \frac{e^{z_j}}{\sum_{i=1}^d e^{z_i}} \leq 1$$

$$\sum_{j=1}^d \frac{e^{z_j}}{\sum_{i=1}^d e^{z_i}} = 1$$

facts  
linke

```
in[2]: import torch
in[3]: sample_vector = torch.rand*(1,10)
in[4]: sample_vector
Out[4]:
tensor([[0.6084, 0.8744, 0.8193, 0.8784, 0.2906, 0.7169, 0.0126, 0.8783, 0.3822,
         0.3883]])
in[5]: torch.softmax(sample_vector, dim=1)
Out[5]:
tensor([[0.0985, 0.1285, 0.1216, 0.1290, 0.0717, 0.1098, 0.0543, 0.1290, 0.0786,
         0.0790]])
in[6]: torch.sum(torch.softmax(sample_vector, dim=1))
Out[6]: tensor(1.0000)
```



# Multi-target (Multi-label)

Optimization:  $\text{Min } -\sum_{(x,y) \in D} \text{Log } P(y|x_y)$

$$[x_0 \ x_1 \ \dots \ x_d \ 1] \begin{bmatrix} w_1 & w'_1 & w''_1 \\ w_2 & w'_2 & w''_2 \\ \dots & \dots & \dots \\ w_{d+1} & w'_{d+1} & w''_{d+1} \end{bmatrix} = [1.5, -1.4, 10] \Rightarrow \text{Softmax}$$

$\Rightarrow [2.0342e-04, 1.1193e-05, 9.9979e-01]$

```
in[12]: torch.softmax(torch.as_tensor([1.5, -1.4, 10.0]).view(-1), dim=0)
Out[12]: tensor([2.0342e-04, 1.1193e-05, 9.9979e-01])
```

$$[x_0 \ x_1 \ \dots \ x_d \ 1] \begin{bmatrix} w_1 & w'_1 & w''_1 \\ w_2 & w'_2 & w''_2 \\ \dots & \dots & \dots \\ w_{d+1} & w'_{d+1} & w''_{d+1} \end{bmatrix} = [1.5, -1.4, 10]$$

$\Rightarrow \text{sigmoid} = [0.8176, 0.1978, 1.0000]$   
 $\Rightarrow \text{softmax} \Rightarrow [0.3652, 0.1965, 0.4383]$

```
in[14]: torch.sigmoid(torch.as_tensor([1.5, -1.4, 10.0]).view(-1))
Out[14]: tensor([0.8176, 0.1978, 1.0000])
in[15]: torch.softmax(torch.sigmoid(torch.as_tensor([1.5, -1.4, 10.0]).view(-1)), dim=0)
Out[15]: tensor([0.3652, 0.1965, 0.4383])
```

# Multi-target (Multi-label)

---

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$

Multitarget:  $(x, [1, 0, 1]) \overset{?}{\leftrightarrow} (x, [0.8176, 0.1978, 1.0])$

$(x, \{c1\}) \overset{?}{\leftrightarrow} (x, [0.8176, \blacksquare, \blacksquare])$

$(x, \{c2\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, 0.1978, \blacksquare]) \overset{?}{\leftrightarrow} (x, [\blacksquare, 1 - 0.1978, \blacksquare])$

$(x, \{c3\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, \blacksquare, 1.0])$

$$\begin{aligned} & - [\text{Log } P(c1|x_{c1}) + \text{Log } P(c2|x_{c2}) + \text{Log } P(c3|x_{c3})] \\ & - [\text{Log } (0.81) + \text{Log } (1-0.19) + \text{Log } (1.0)] \end{aligned}$$

# Multiclass

---

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$

Multiclass:  $(x, [0, 1, 0]) \overset{?}{\leftrightarrow} (x, [0.3652, 0.1965, 0.4383])$

$(x, \{1\}) \overset{?}{\leftrightarrow} (x, [0.3652, \blacksquare, \blacksquare]) \overset{?}{\leftrightarrow} (x, [1 - 0.3652, \blacksquare, \blacksquare])$

$(x, \{2\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, 0.1965, \blacksquare])$

$(x, \{3\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, \blacksquare, 0.4383]) \overset{?}{\leftrightarrow} (x, [\blacksquare, \blacksquare, 1 - 0.4383])$

# Multiclass

---

Optimization:  $\text{Min } -\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$

Multiclass:  $(x, [0, 1, 0]) \overset{?}{\leftrightarrow} (x, [0.3652, 0.1965, 0.4383])$

~~$(x, \{1\}) \overset{?}{\leftrightarrow} (x, [0.3652, \blacksquare, \blacksquare]) \overset{?}{\leftrightarrow} (x, [1 - 0.3652, \blacksquare, \blacksquare])$~~

$(x, \{2\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, 0.1965, \blacksquare])$

~~$(x, \{3\}) \overset{?}{\leftrightarrow} (x, [\blacksquare, \blacksquare, 0.4383]) \overset{?}{\leftrightarrow} (x, [\blacksquare, \blacksquare, 1 - 0.4383])$~~

If  $x$  belongs to one class, it does not belong to other classes  
Also, in softmax, if we increase one element, it reduces other (sum=1)

# Multiclass

---

Optimization: Min  $-\sum_{(x,y) \in \mathcal{D}} \text{Log } P(y|x_y)$

Multiclass:  $(x, [0, 1, 0]) \overset{?}{\leftrightarrow} (x, [0.3652, 0.1965, 0.4383])$

$$-\left( [0, 1, 0] \log \begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix} \right) = +\infty$$

$$-\left( [0, 1, 0] \log \begin{bmatrix} 0.3652 \\ 0.1965 \\ 0.4383 \end{bmatrix} \right) = ?$$

$$-\left( [0, 1, 0] \log \begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix} \right) = 0$$

---

# Evaluation

---

---

# Curves

---





# Threshold-based Model

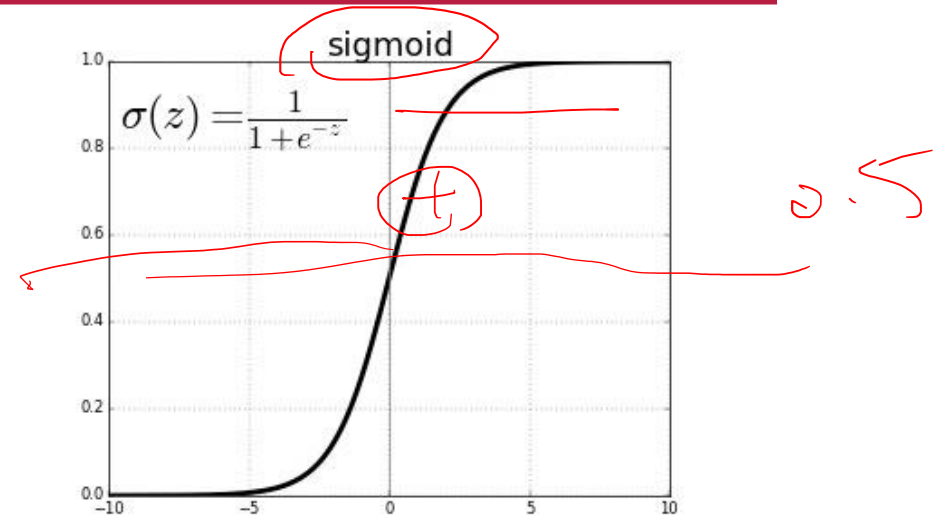
$$P(+ | x) = 1 - P(- | x)$$

$$P(x) = \text{Sigmoid}(f(x))$$

$$P(x) \geq \delta \rightarrow x \text{ is positive}$$

$$P(x) < \delta \rightarrow x \text{ is negative}$$

~~0.6~~  
0.9



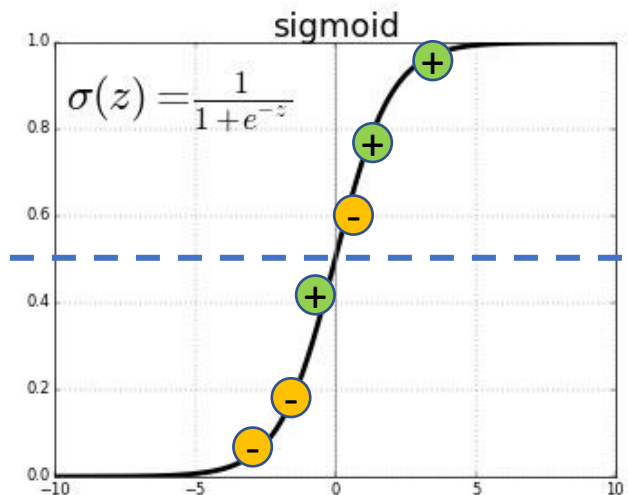
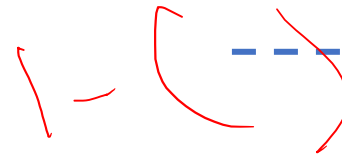
---

# Threshold-based Model

## $\delta=0.5$

---

$P(+|x) = 1 - P(-|x)$   
 $P(x) = \text{Sigmoid}(f(x))$   
 $P(x) \geq 0.5 \rightarrow x \text{ is positive}$   
 $P(x) < 0.5 \rightarrow x \text{ is negative}$



---

# Threshold-based Model

$\delta=0.0 \rightarrow$  Biased Model  $\rightarrow$  All are positives

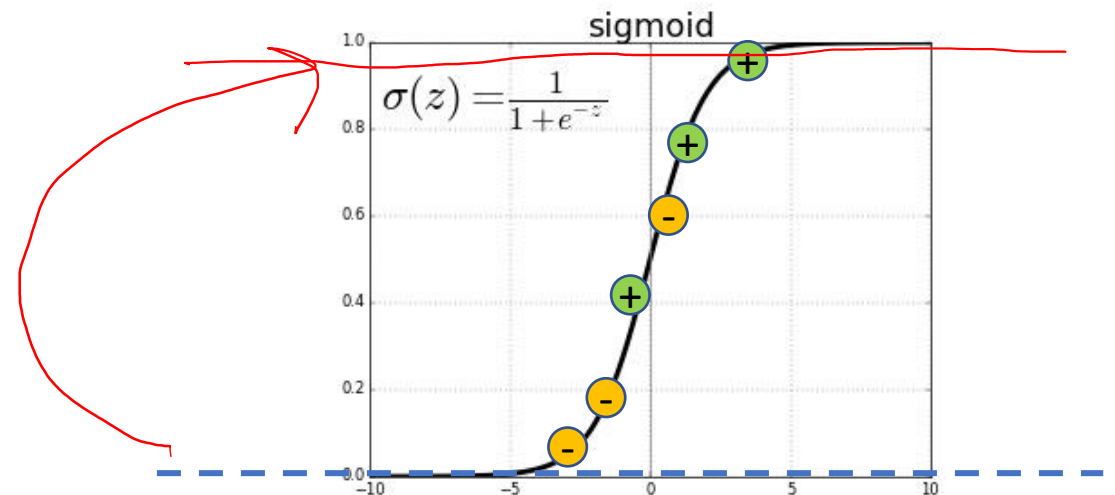
*negative*

$$P(+|x) = 1 - P(-|x)$$

$$P(x) = \text{Sigmoid}(f(x))$$

$$P(x) \geq 0.0 \rightarrow x \text{ is positive}$$

$$P(x) < 0.0 \rightarrow x \text{ is negative}$$



---

# Threshold-based Model

$\delta=1.0 \rightarrow$  Biased Model  $\rightarrow$  All are negatives

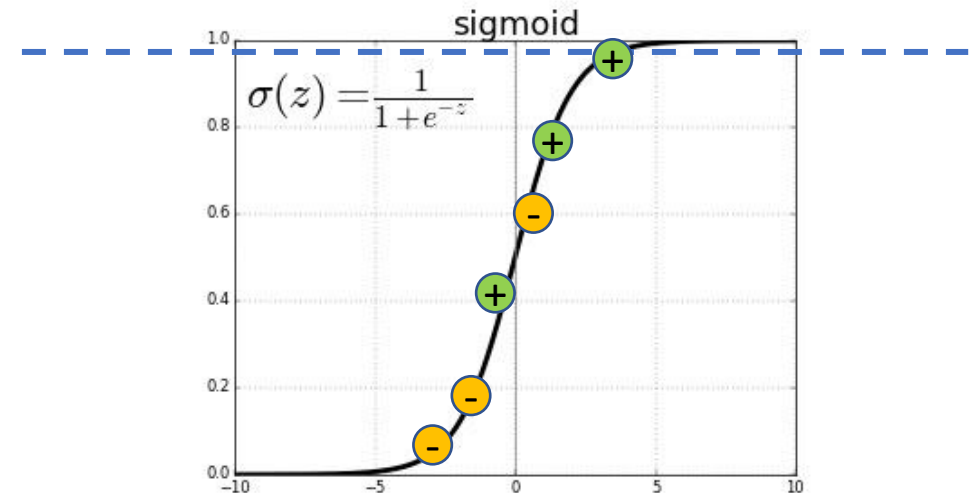
---

$$P(+|x) = 1 - P(-|x)$$

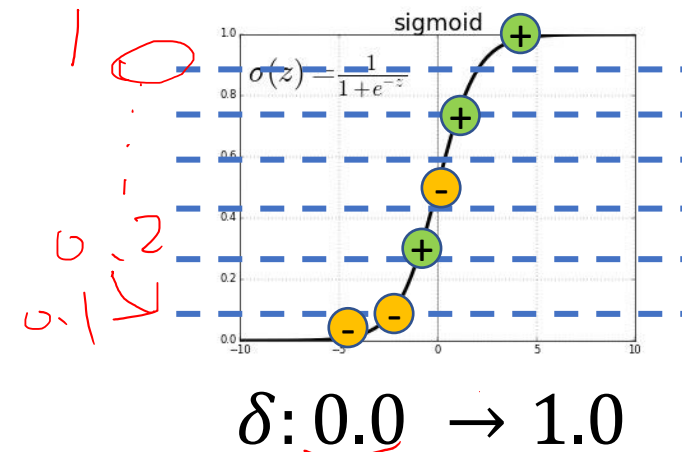
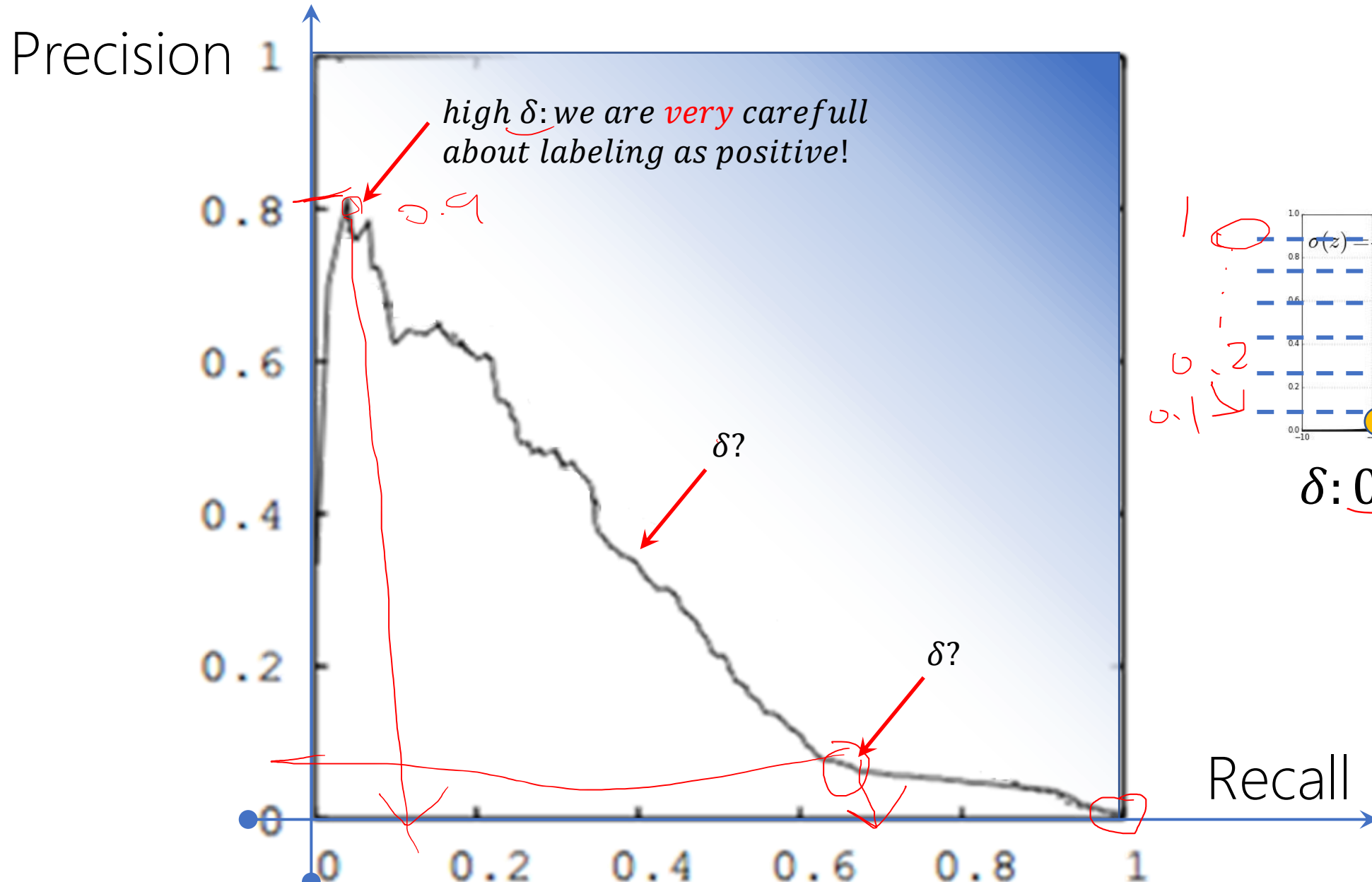
$$P(x) = \text{Sigmoid}(f(x))$$

$$P(x) \geq 1.0 \rightarrow x \text{ is positive}$$

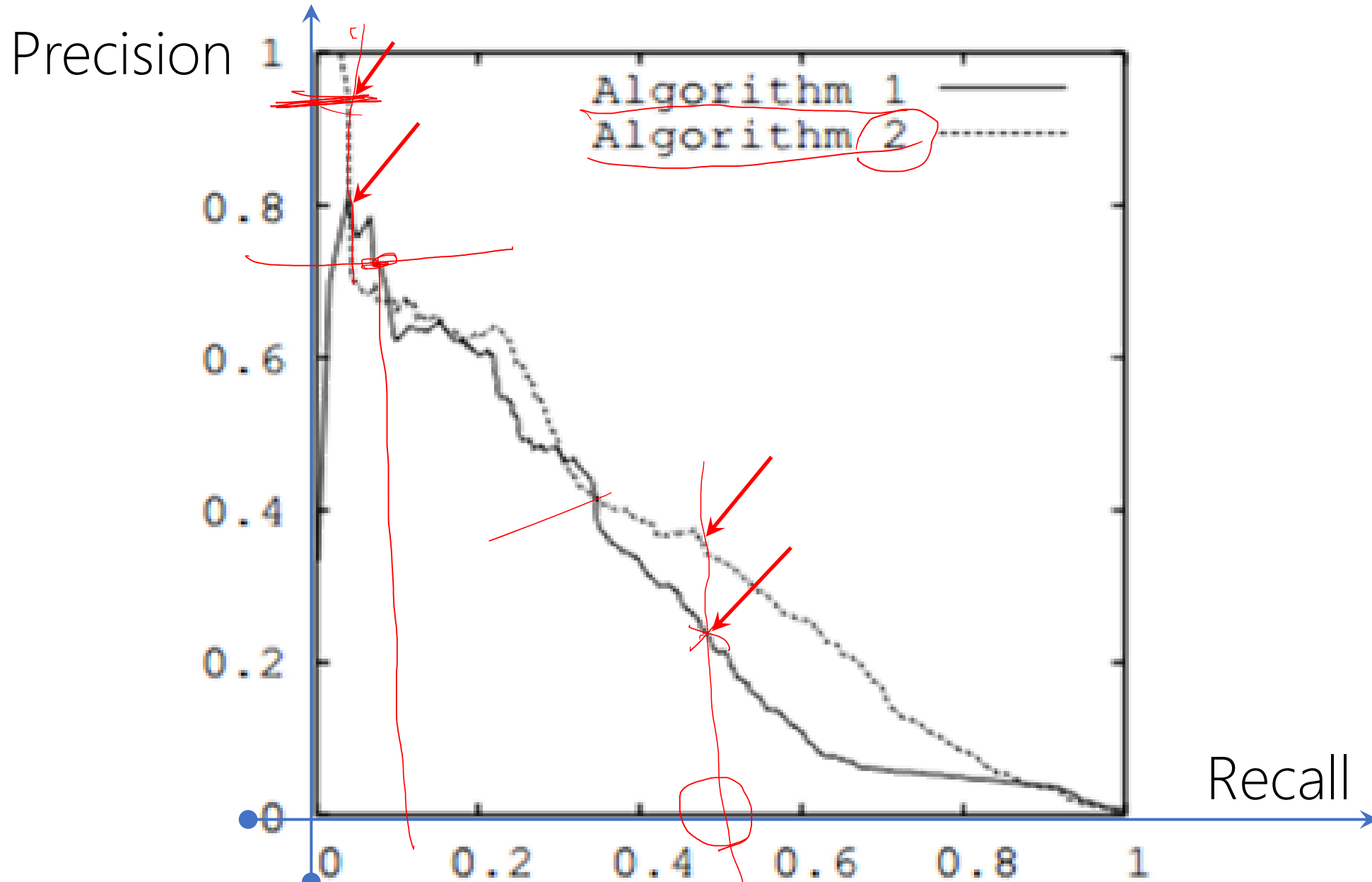
$$P(x) < 1.0 \rightarrow x \text{ is negative}$$



# Precision-Recall Curve: Best $\delta$



# Precision-Recall Curve: Model Comparison



---

# Receiver Operating Characteristic ROC

---

The ROC curve was first developed by electrical engineers and radar engineers during World War II for detecting enemy objects in battlefield!

Missile attack vs. passenger airplane!

# Recall aka True Positive Rate (TPR)

What percentage of positives are captured.

		<i>gold standard labels</i>		
		gold positive	gold negative	
<i>system output labels</i>	system positive	<b>true positive</b>	false positive	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	<b>true negative</b>	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

**Figure 4.4** Contingency table



# False Positive Rate (FPR)

What percentage are *incorrectly* captured as positives!

		gold standard labels		
		gold positive	gold negative	
system output labels	system positive	true positive	<b>false positive</b>	precision = $\frac{tp}{tp+fp}$
	system negative	false negative	<b>true negative</b>	
		recall = $\frac{tp}{tp+fn}$		accuracy = $\frac{tp+tn}{tp+fp+tn+fn}$

**Figure 4.4** Contingency table

---

## Perfect Classifier

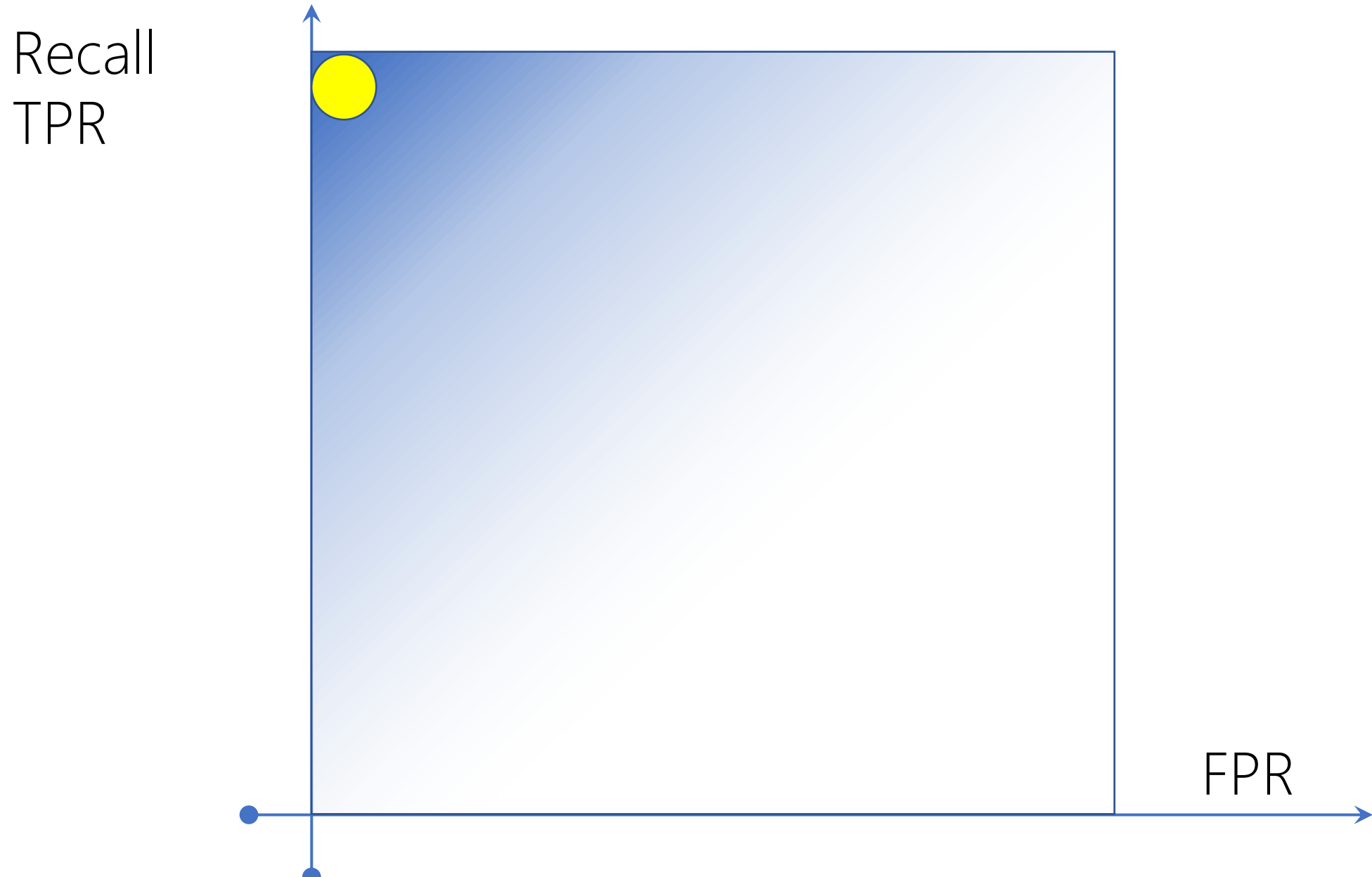
---

	Gold Positive	Gold Negative
Model Positive	N+	0
Model Negative	0	N-

$$\text{TPR} = \frac{N+}{(N+)+0} = 1.0$$

$$\text{FPR} = \frac{0}{0+(N-)} = 0.0$$

# TPR-FPR Curve



---

## Perfect Classifier

---


	Gold Positive	Gold Negative
Model Positive	$N+$ ↓	↑ 0 ↑
Model Negative	0 ↑	↓ $N-$
$TPR = \frac{N+}{(N+)+0} = 1.0$ ↓		$FPR = \frac{0}{0+(N-)} = 0.0$ ↑


---

## Worst Classifier

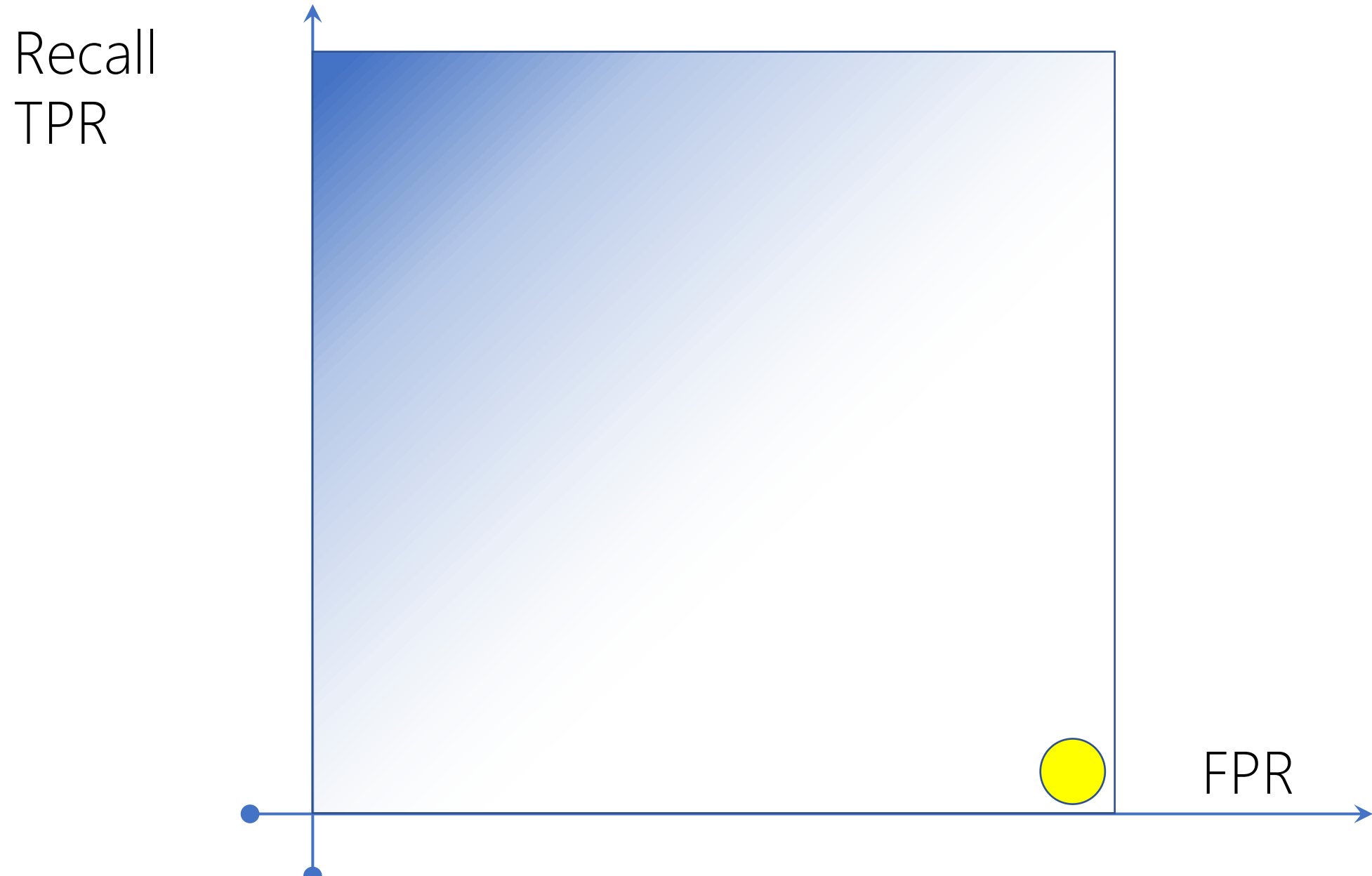
---

	Gold Positive	Gold Negative
Model Positive	0	N-
Model Negative	N+	0

$$\text{TPR} = \frac{0}{(N+) + 0} = 0.0$$


$$\text{FPR} = \frac{N-}{0 + (N-)} = 1.0$$


# TPR-FPR Curve



---

## Uniformly Random Classifier

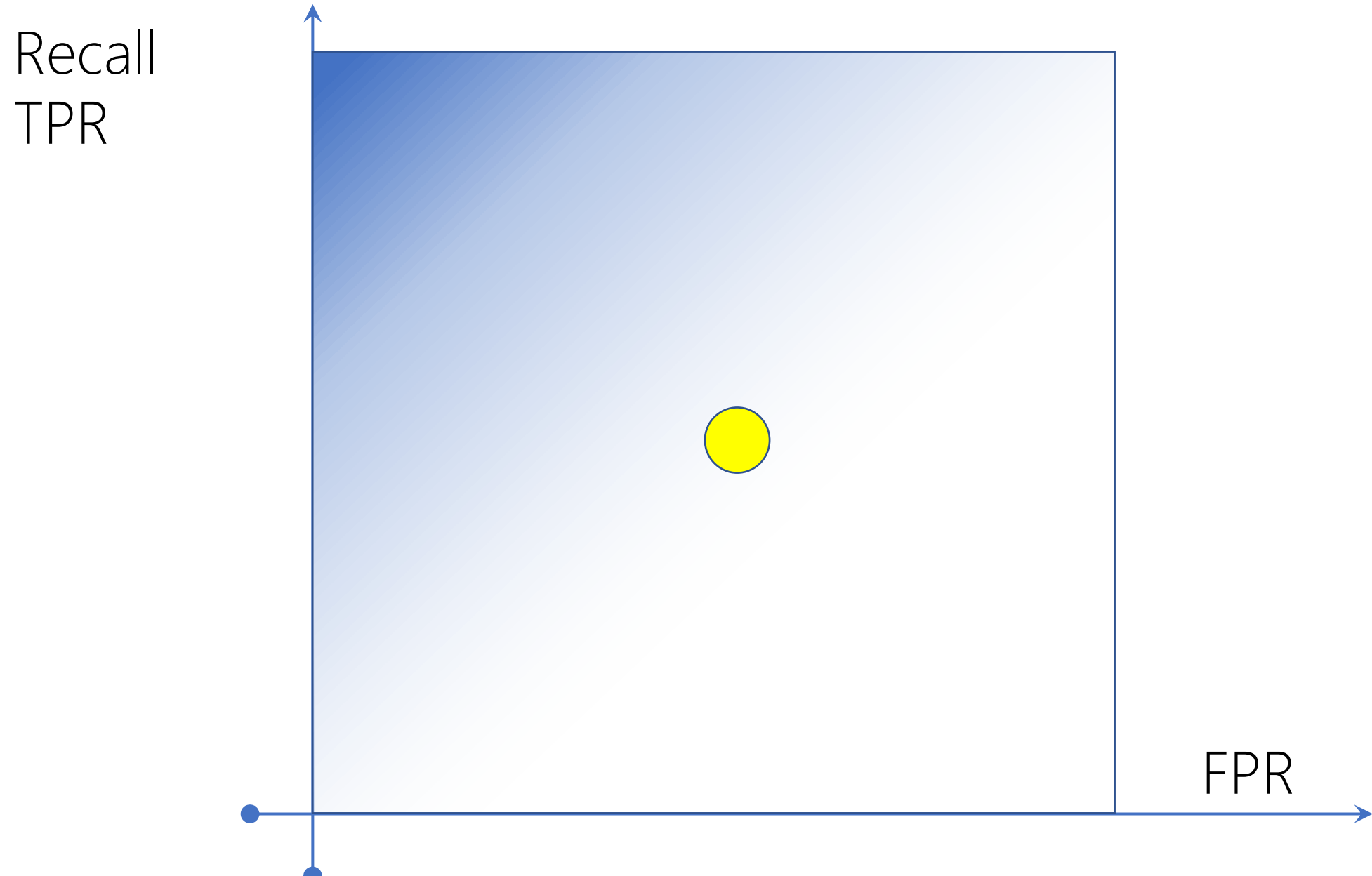
---

	Gold Positive	Gold Negative
Model Positive	?	?
Model Negative	?	?

TPR=0.5

FPR=0.5

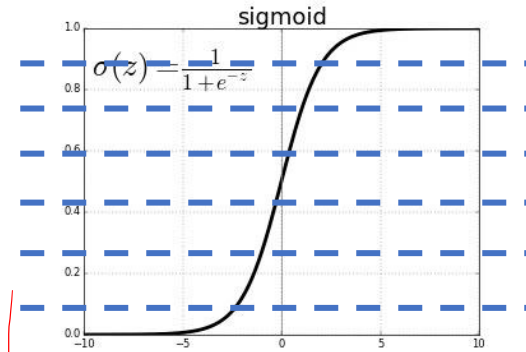
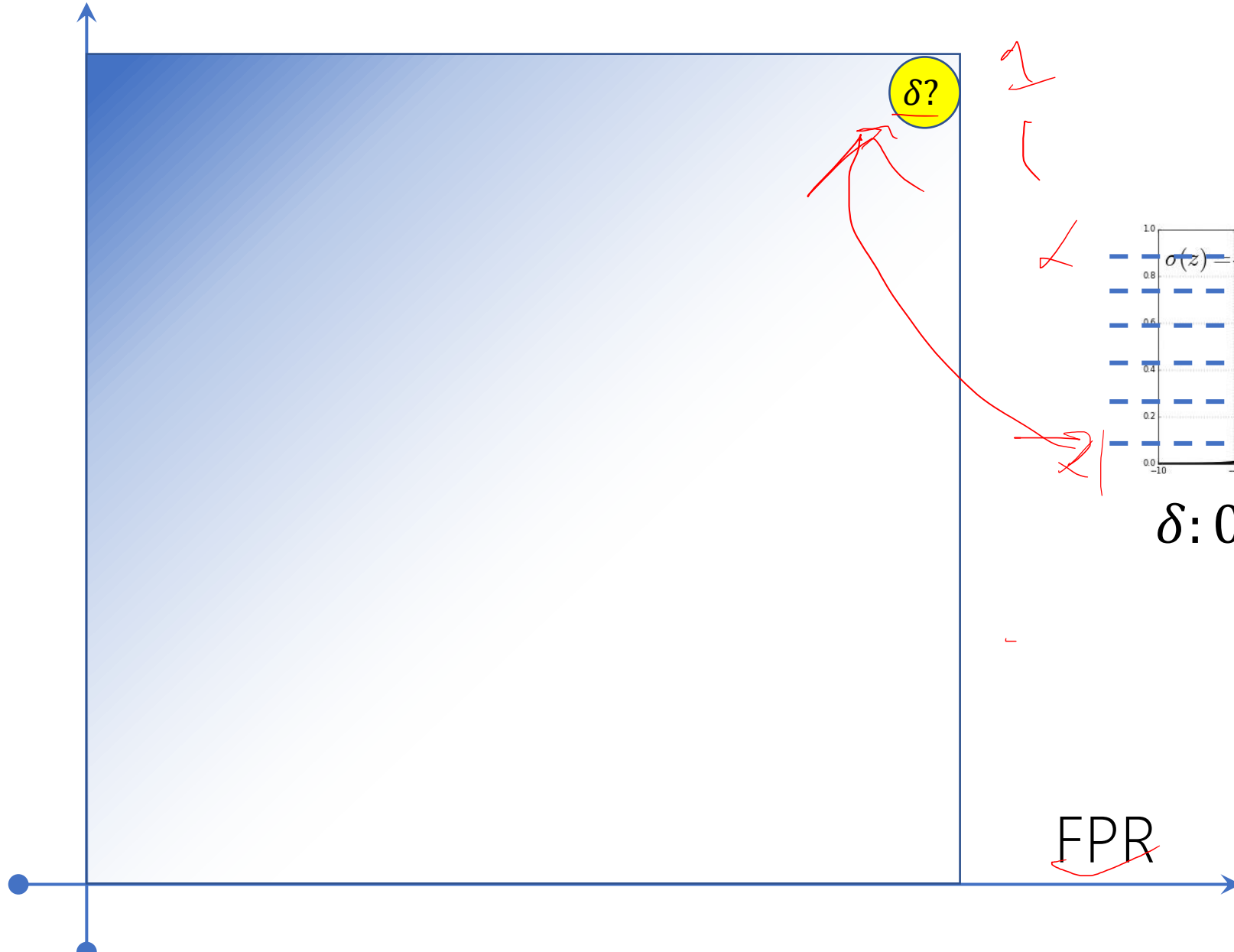
# TPR-FPR Curve





# TPR-FPR Curve

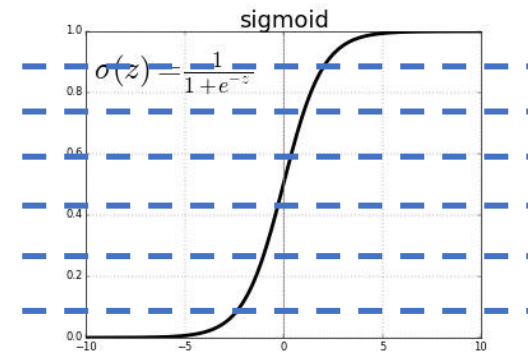
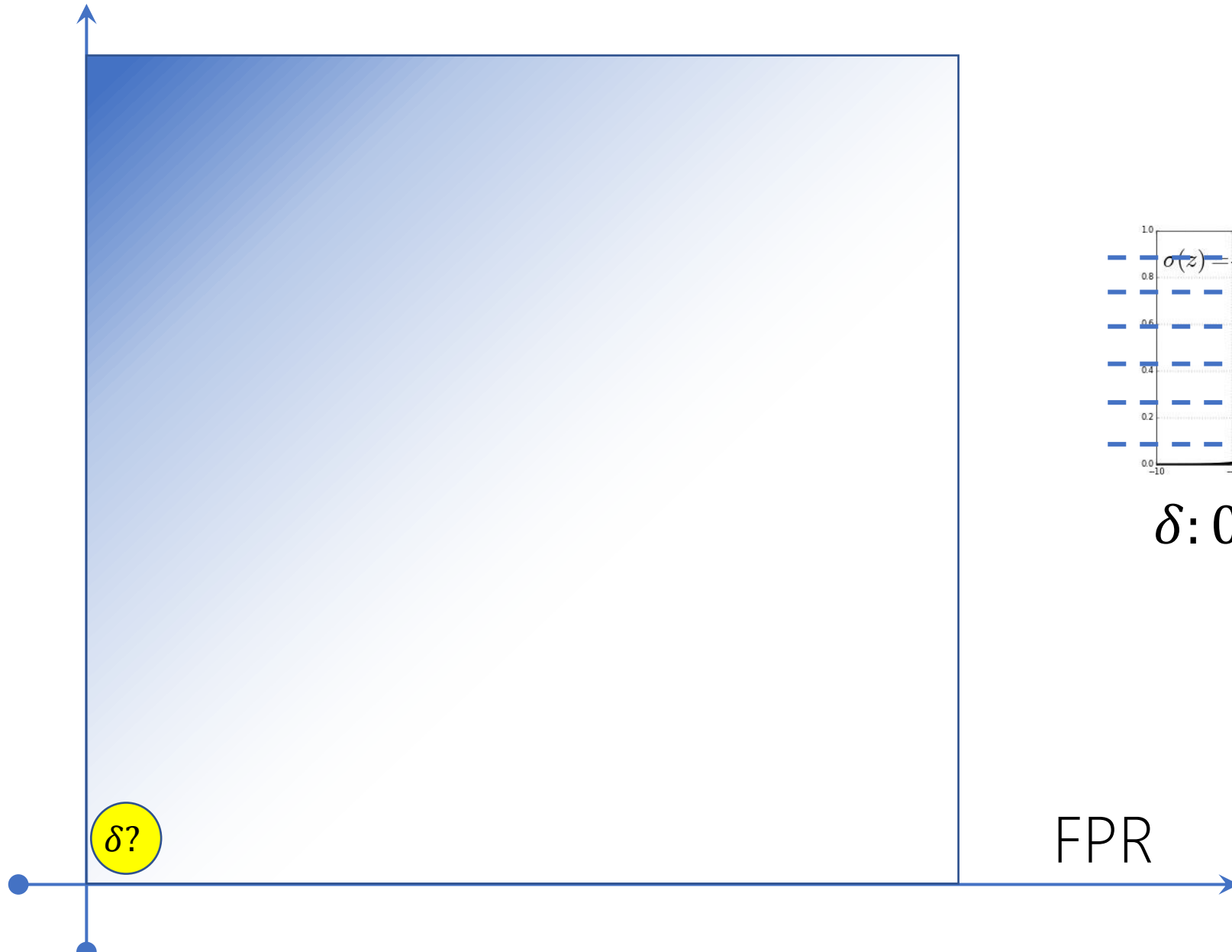
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# TPR-FPR Curve

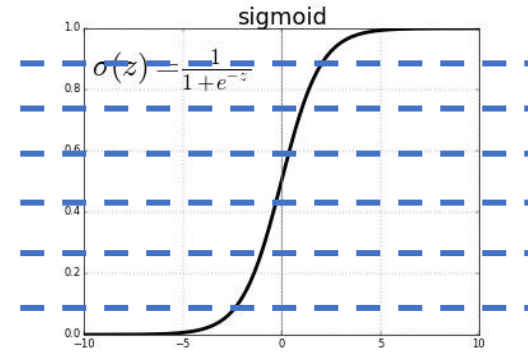
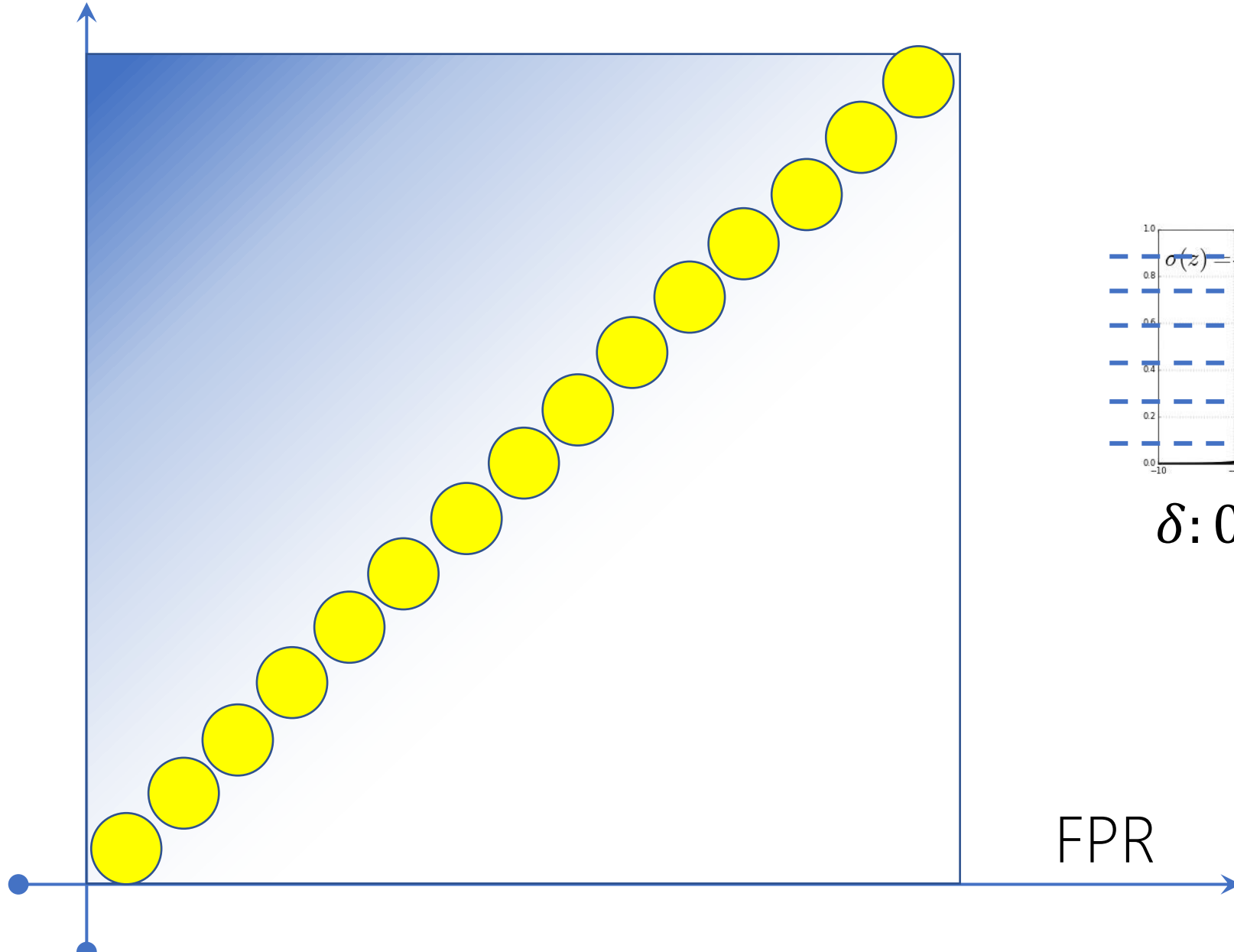
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# What is this model?

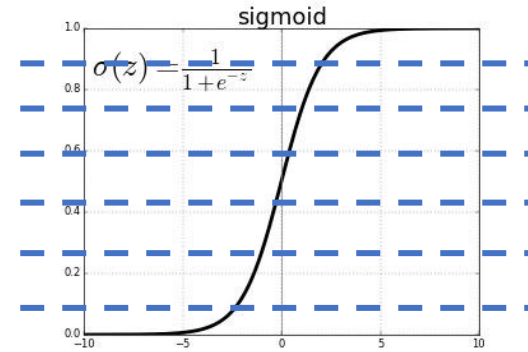
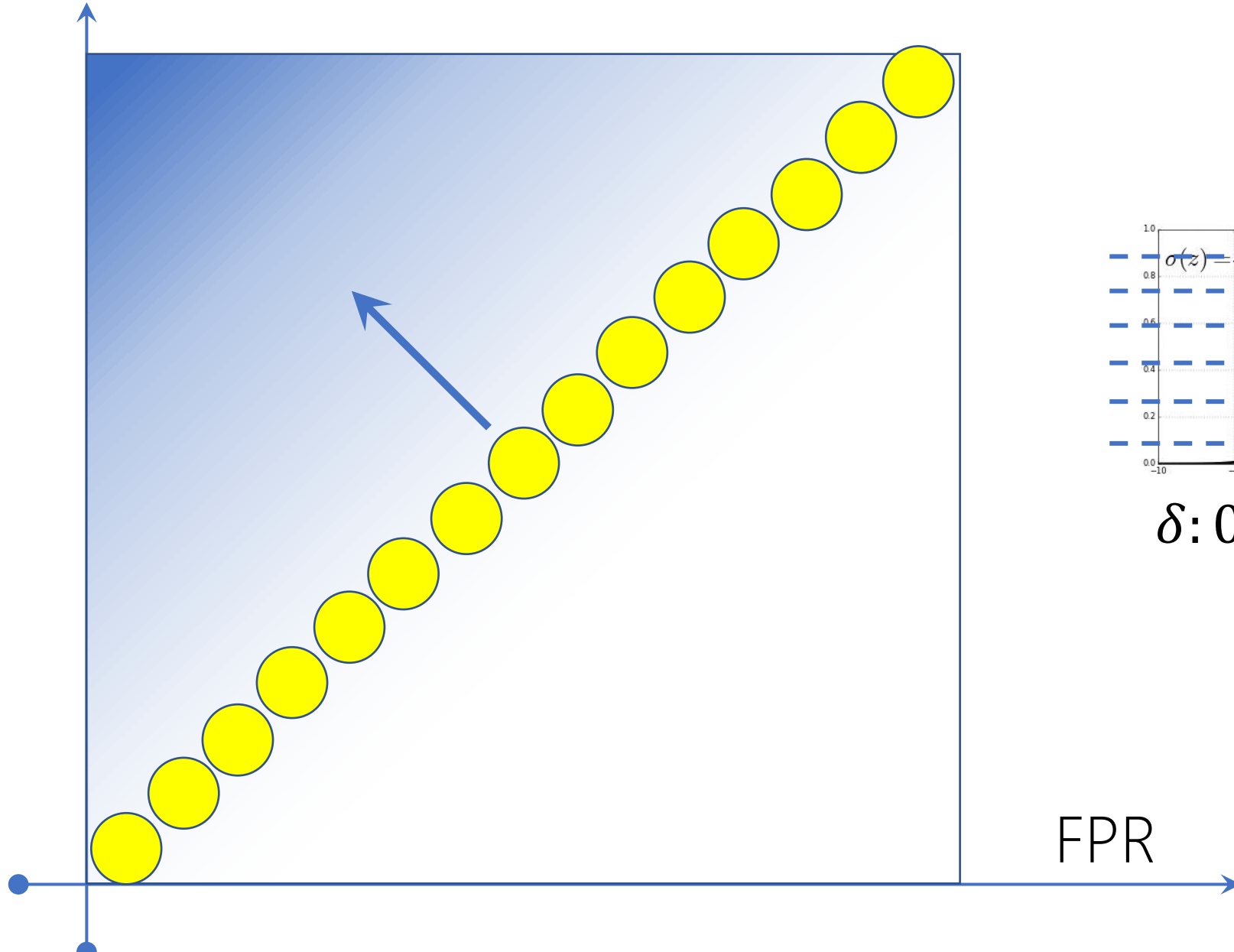
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# ROC

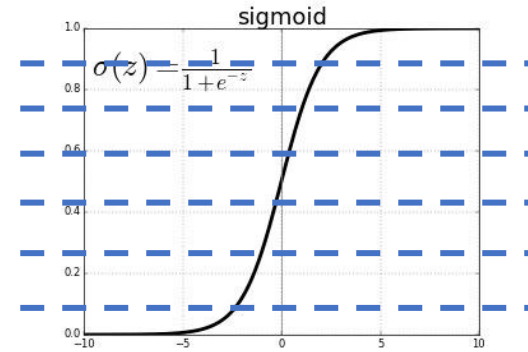
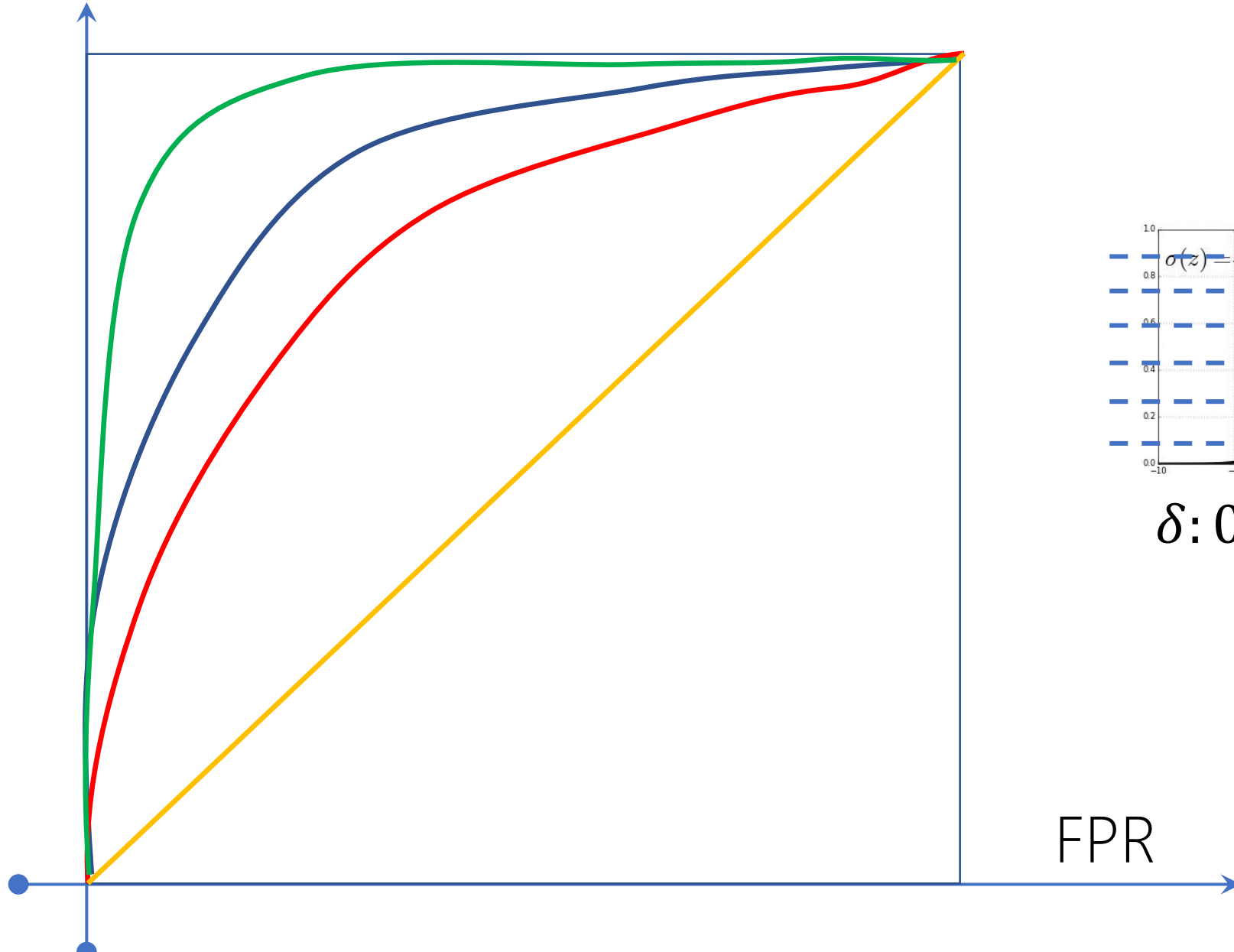
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# ROC

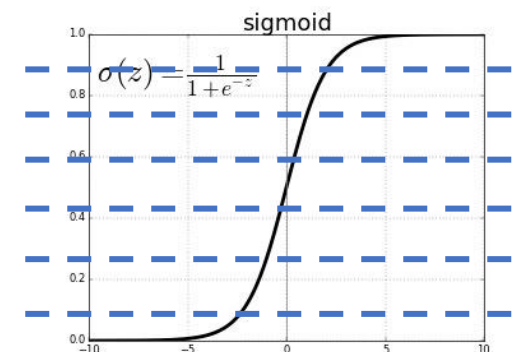
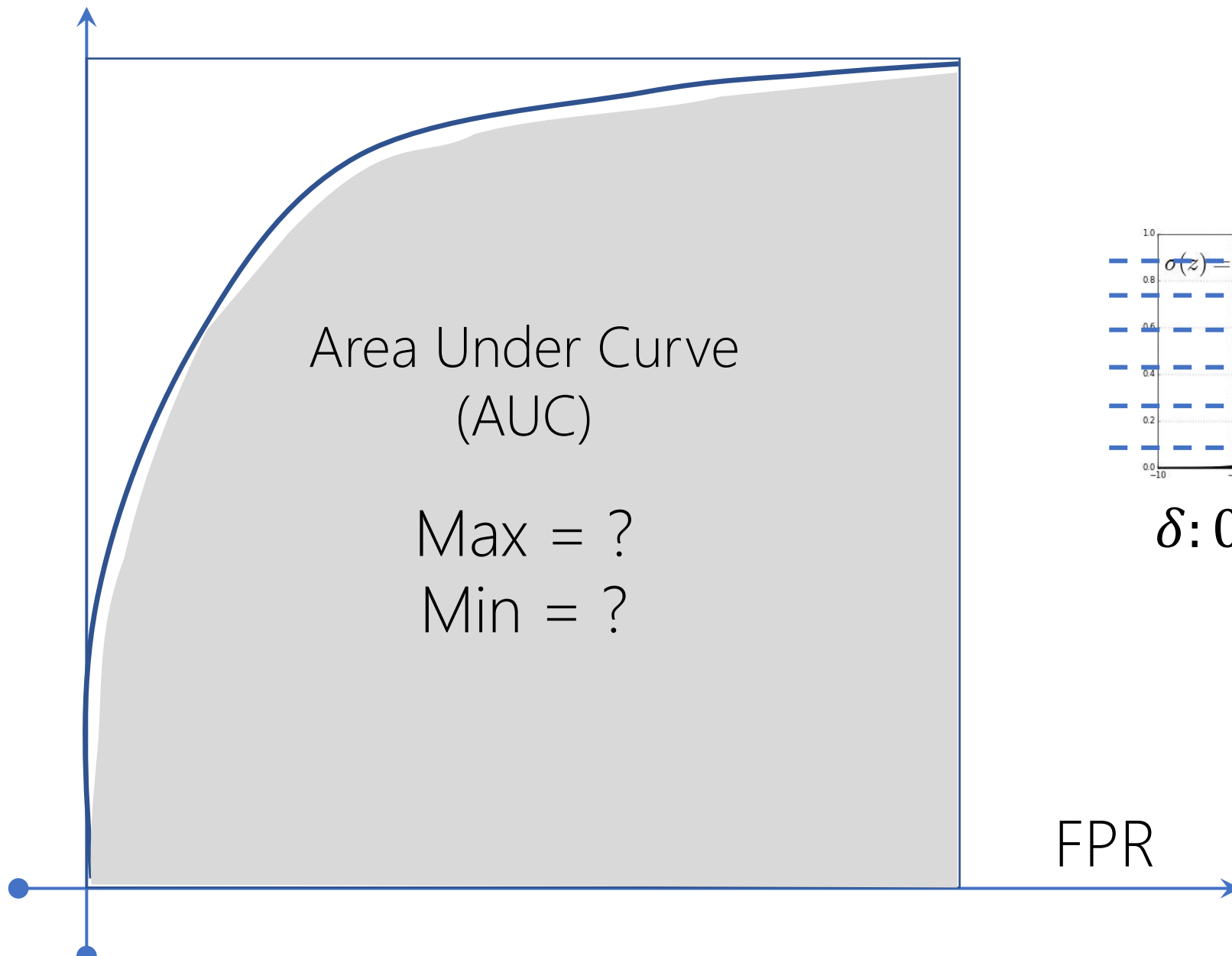
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# Area Under Curve (AUC): Single Real Point

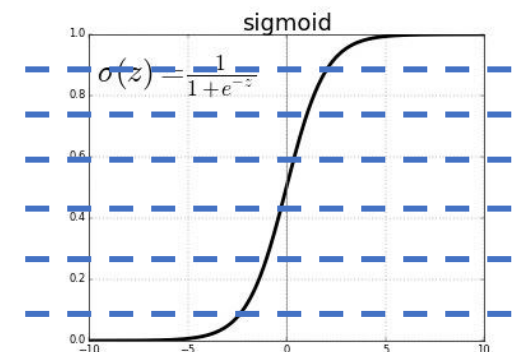
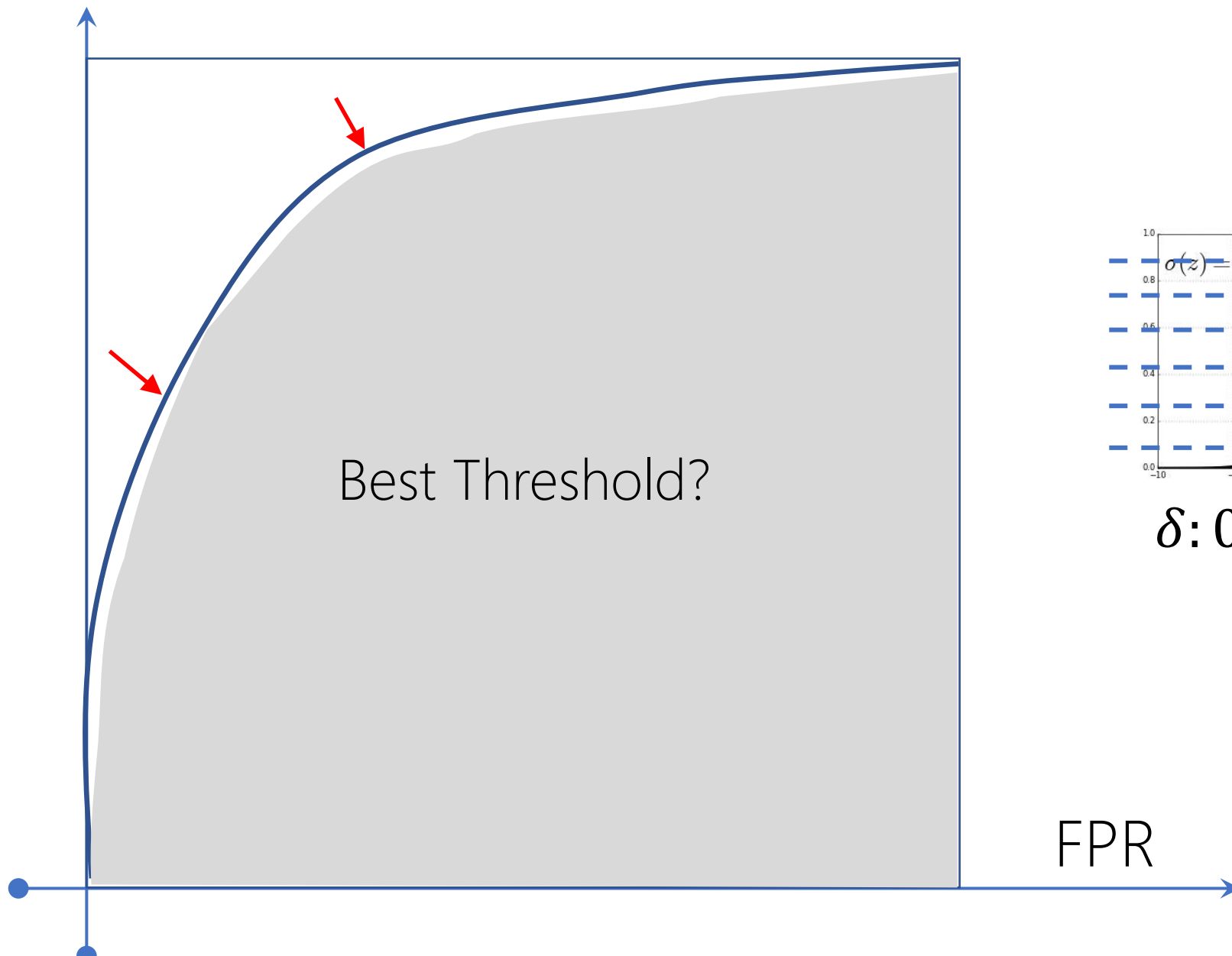
Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

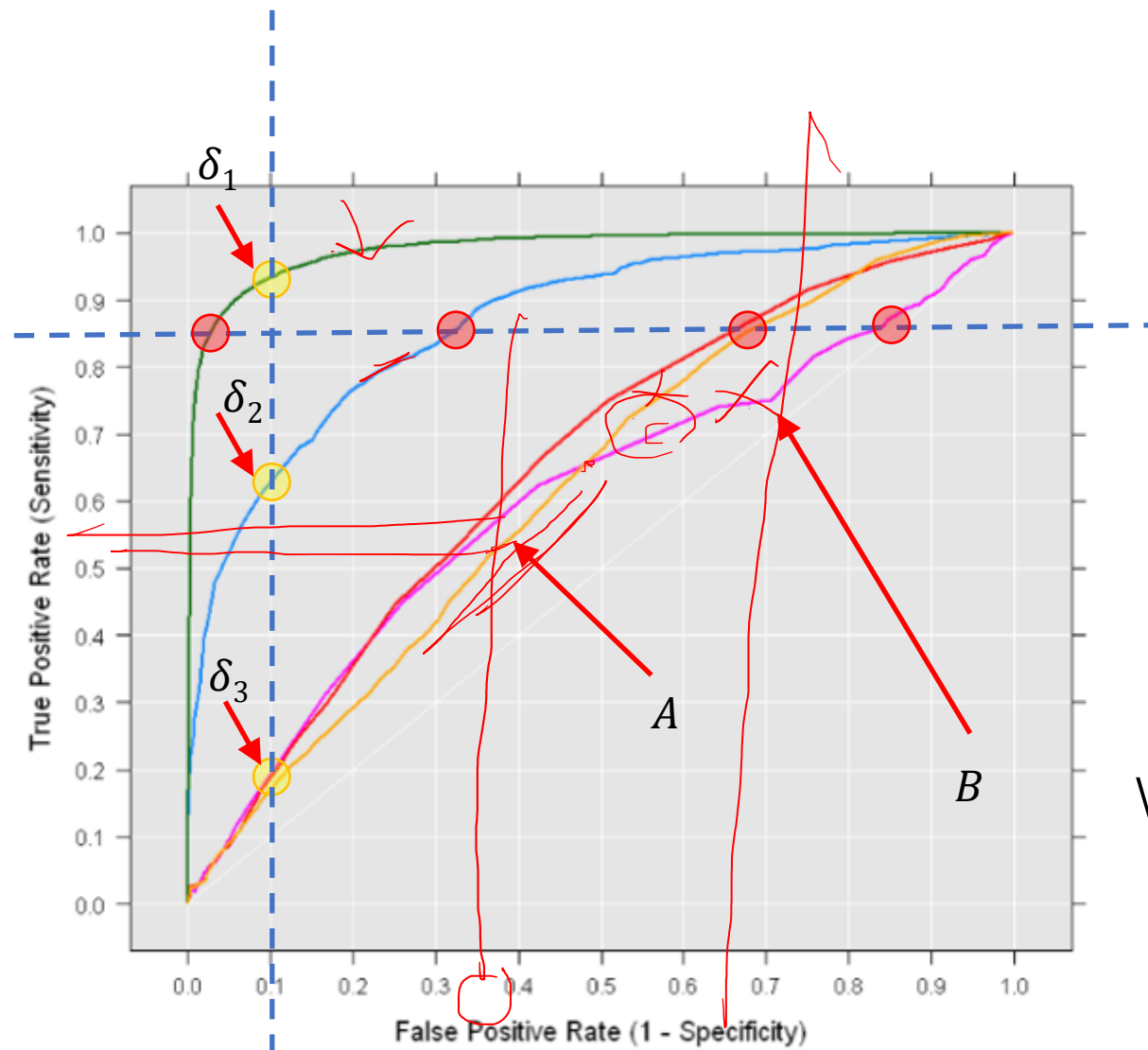
# Area Under Curve (AUC): Single Real Point

Recall  
TPR



$\delta: 0.0 \rightarrow 1.0$

# ROC: Model Comparison



Which one? A or B?

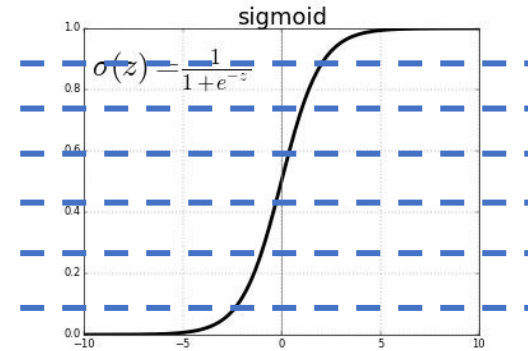
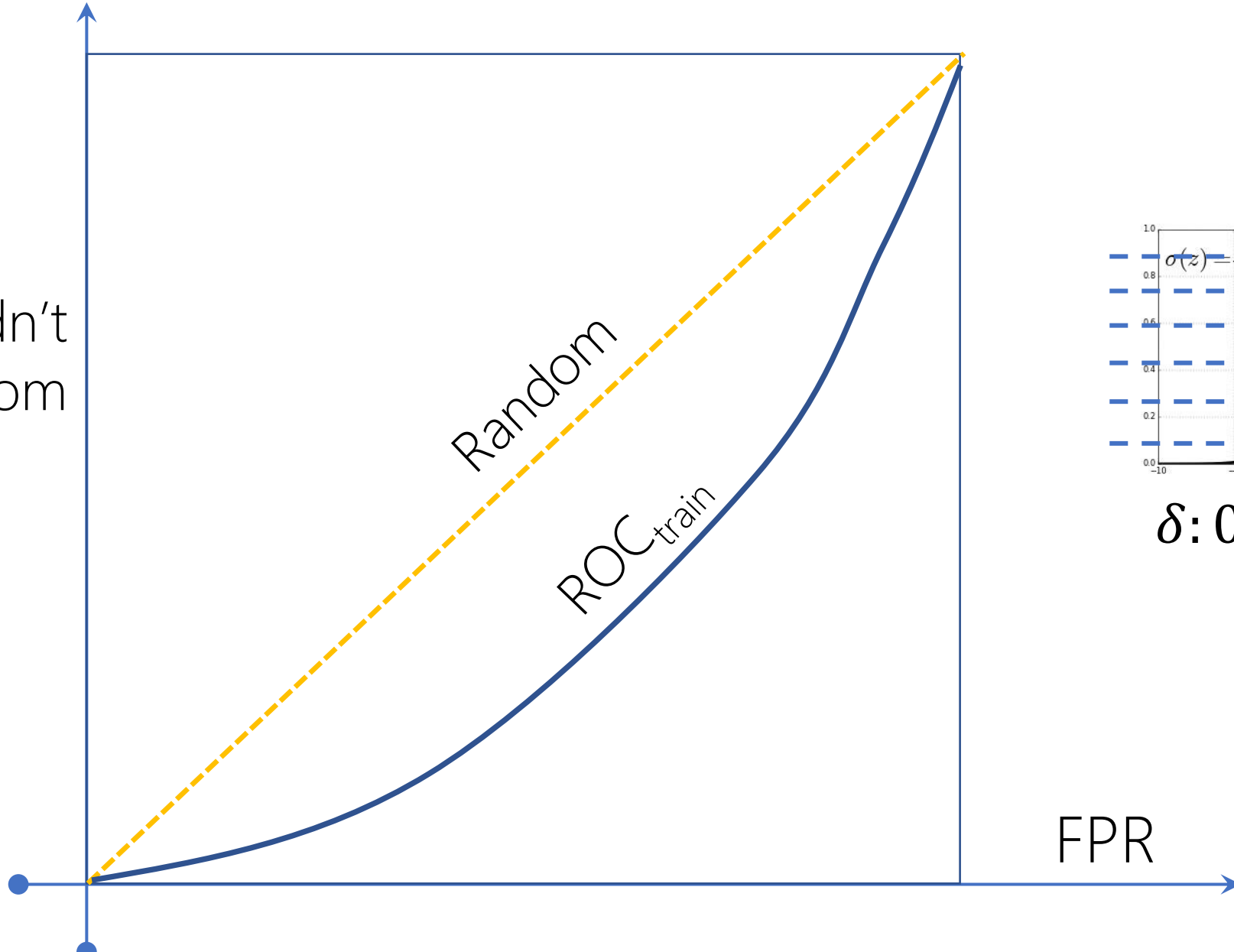


# ROC<sub>train</sub>

Recall  
TPR

The model couldn't  
learn anything from  
training set!

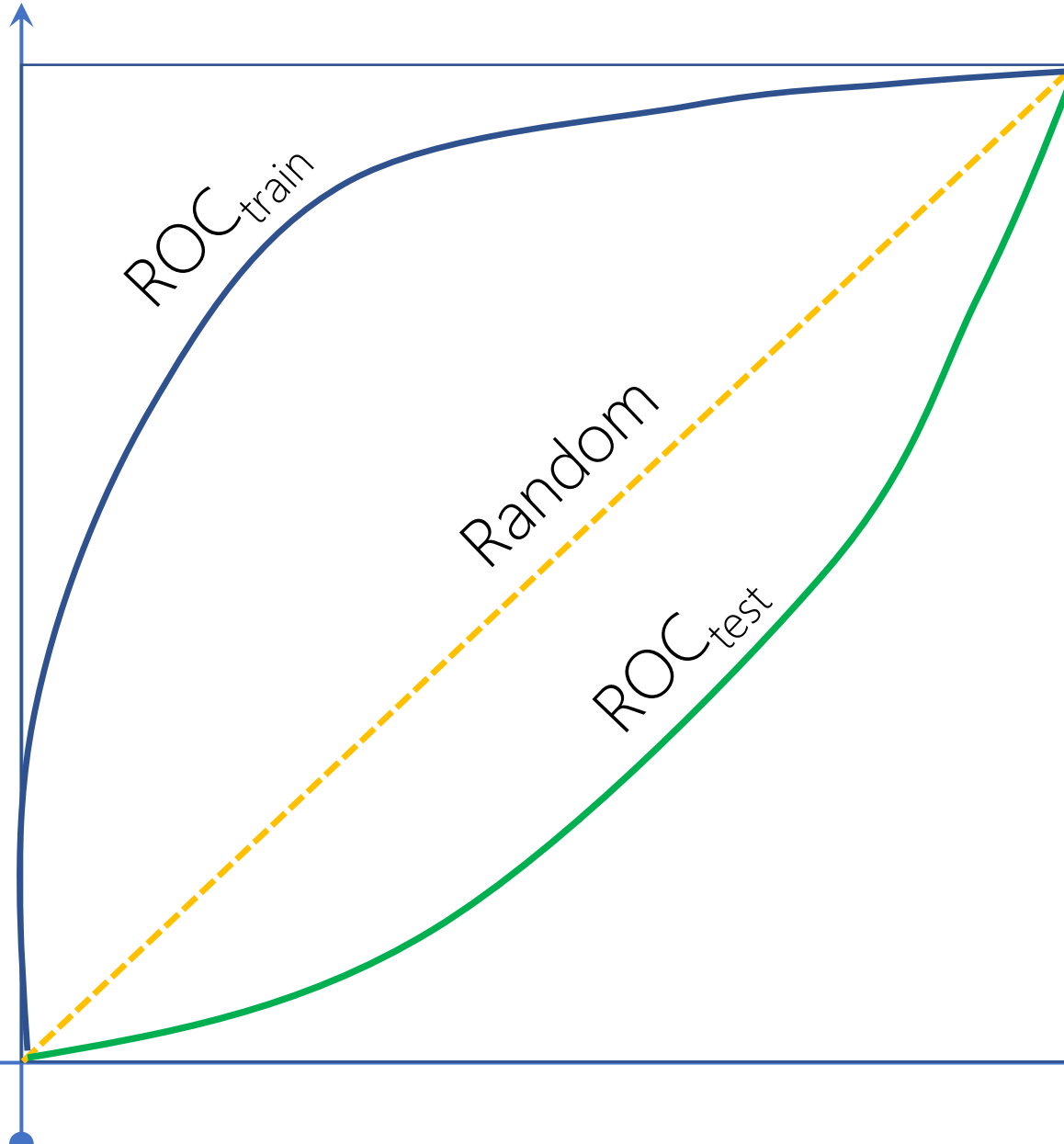
ROC<sub>Test</sub> = ?



$\delta: 0.0 \rightarrow 1.0$

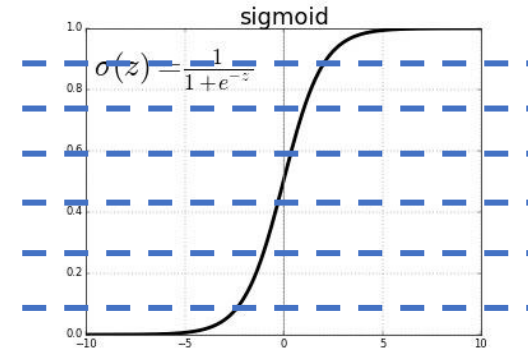
# ROC<sub>test</sub>

Recall  
TPR



The model performs well on train set. It means it learnt!

But performs poor in test set. It only know train set.  
Cannot generalize!

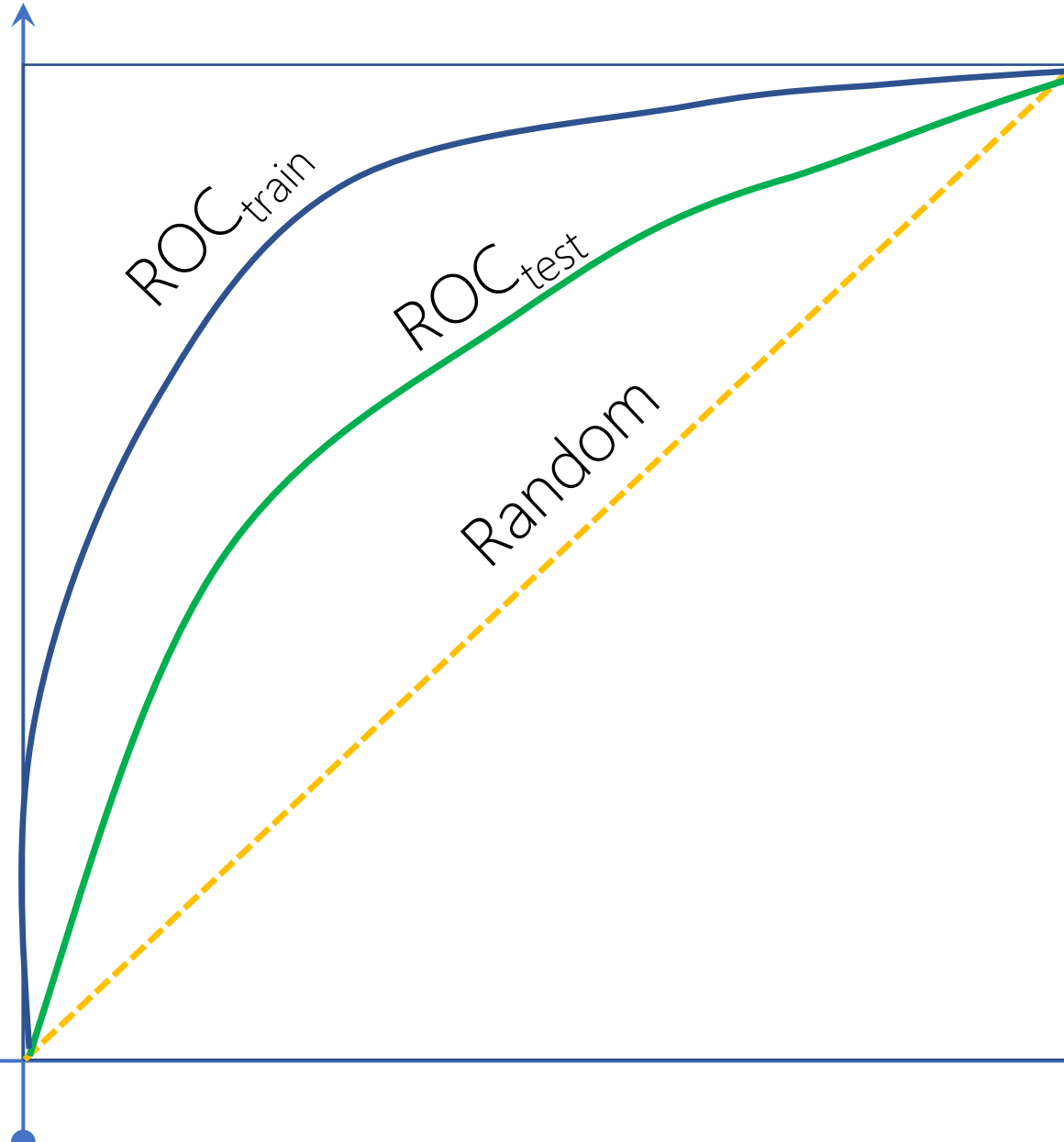


$\delta: 0.0 \rightarrow 1.0$

FPR

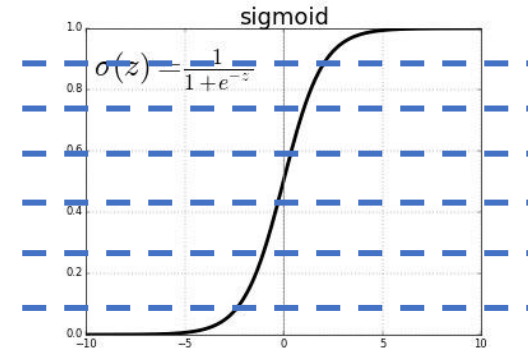
# ROC<sub>test</sub>

Recall  
TPR



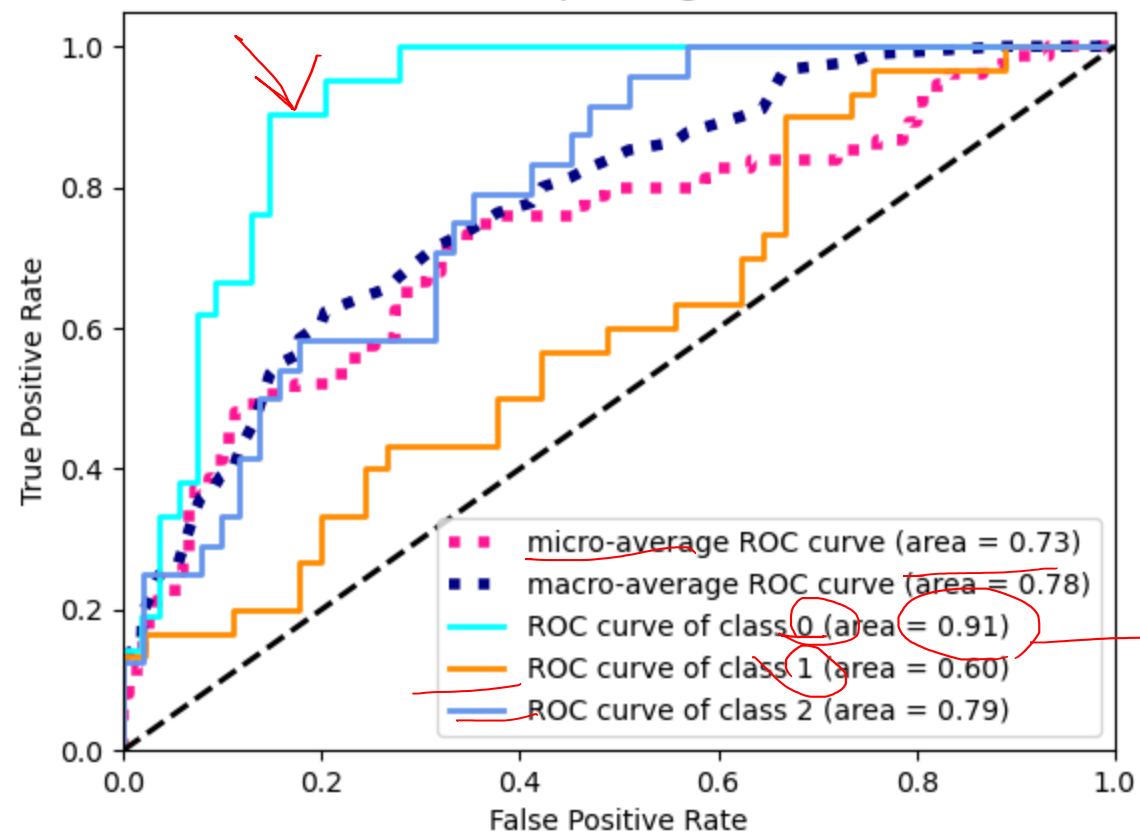
The model performs well on train set. It means it learnt!

But performs poor in test set. It only know train set.  
Cannot generalize!



$\delta: 0.0 \rightarrow 1.0$

Some extension of Receiver operating characteristic to multi-class



---

So Far, the LR model and its output

---

---

How to input text as  $X$  to LR?

---

$X$  is a vector!

How to map text into vector?



Text:  $\{w_1, w_5, w_6\}$

$[10001100-]$   
(14)

*Ceci n'est pas une pipe.*

# Word vector space model

***The Treachery of Images***



Artist	René Magritte
Year	1929
Medium	Oil on canvas
Movement	Surrealism
Dimensions	60.33 cm × 81.12 cm (23.75 in × 31.94 in)
Location	Los Angeles County Museum of Art <sup>[1]</sup>



- Phonetics and Phonology  
knowledge about linguistic sounds
- Morphology  
knowledge of the formation and internal structure of words
- Syntax  
knowledge of the structural relationships between words
- Semantics  
knowledge of meaning
- Pragmatics  
knowledge of the relationship of meaning to the goals & intentions of the speaker
- Discourse  
knowledge about linguistic units larger than a single utterance

Task: Engaging in Natural Language Communication



---

# Semiotics: The Science of Symbols

---

Semantics: Relation between signs and things to which they refer: meaning; sense

Syntactics: Relations among signs in formal structures

Pragmatics: Relation between signs and sign-using agents

