

# Machine Learning

CS 539

Worcester Polytechnic Institute  
Department of Computer Science  
Instructor: Prof. Kyumin Lee

# Project Teams

- Will Buchta, Sarah Semy, Deena Selitsk, Humza Qureshi
- Kai Davidson, Sarah Kogan, Dominic Ranelli
- Bryce Lukens, Harrison Taylor, Gabriel Barbosa
- Axel Luca, Preet Patel, Xiaojun Liu
- Codey Battista, Avinash Bissoondial, Nick Chantre, Paul Crann
- John DeLeo, Jacob Donahue, Joe Molder, Iris Nycz
- Srisaranya Pujari, Kangjian Wu, Matthew Yeo, Yang Yi

# Previous Class...

Regularization

# Previous Class...

Regularization

Perceptron

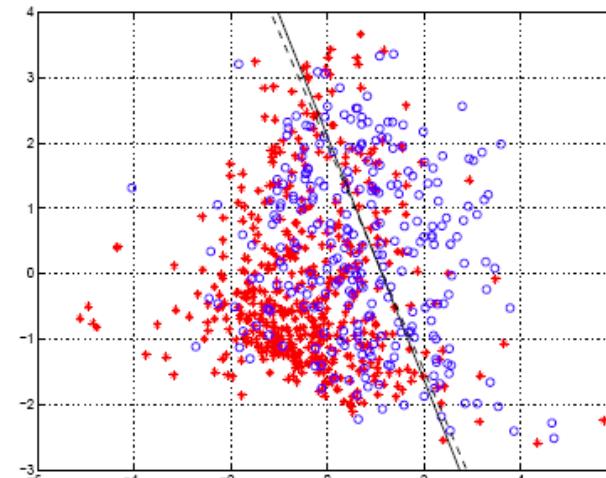
# Logistic Regression

(models probability of output in terms of input)

# Classification based on Probability

- Instead of just predicting the class, give the probability of the instance being that class
  - i.e., learn  $p(y | x)$
- Comparison to perceptron:
  - Perceptron doesn't produce probability estimate
  - Perceptron (and other discriminative classifiers) are only interested in producing a discriminative model

- Recall that:
$$0 \leq p(\text{event}) \leq 1$$
$$p(\text{event}) + p(\neg\text{event}) = 1$$



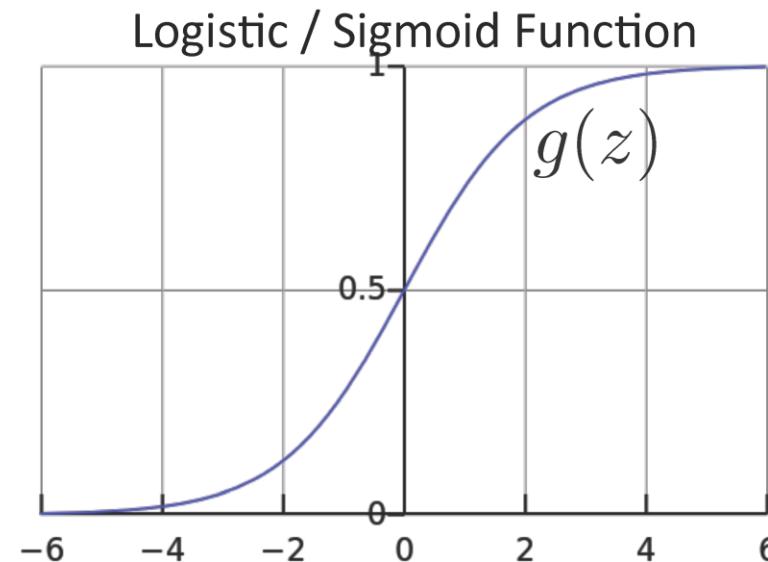
# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)
- $h_{\theta}(x)$  should give  $p(y = 1 | x; \theta)$ 
  - Want  $0 \leq h_{\theta}(x) \leq 1$
- Logistic regression model:

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



# Interpretation of Hypothesis Output

$$h_{\theta}(x) = \text{estimated } p(y = 1 | x; \theta)$$

Example: Cancer diagnosis from tumor size

$$\mathbf{x} = \begin{bmatrix} x_0 \\ x_1 \end{bmatrix} = \begin{bmatrix} 1 \\ \text{tumorSize} \end{bmatrix}$$

$$h_{\theta}(\mathbf{x}) = 0.7$$

→ Tell patient that 70% chance of tumor being malignant

Note that:  $p(y = 0 | x; \theta) + p(y = 1 | x; \theta) = 1$

Therefore,  $p(y = 0 | x; \theta) = 1 - p(y = 1 | x; \theta)$

# Another Interpretation

- Equivalently, logistic regression assumes that

$$\log \frac{p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta})}{p(y = 0 \mid \mathbf{x}; \boldsymbol{\theta})} = \theta_0 + \theta_1 x_1 + \dots + \theta_d x_d$$

odds of  $y = 1$  (Ratio of  $y=1$  and  $y=0$ )

- In other words, logistic regression assumes that the log odds is a linear function of  $\mathbf{x}$

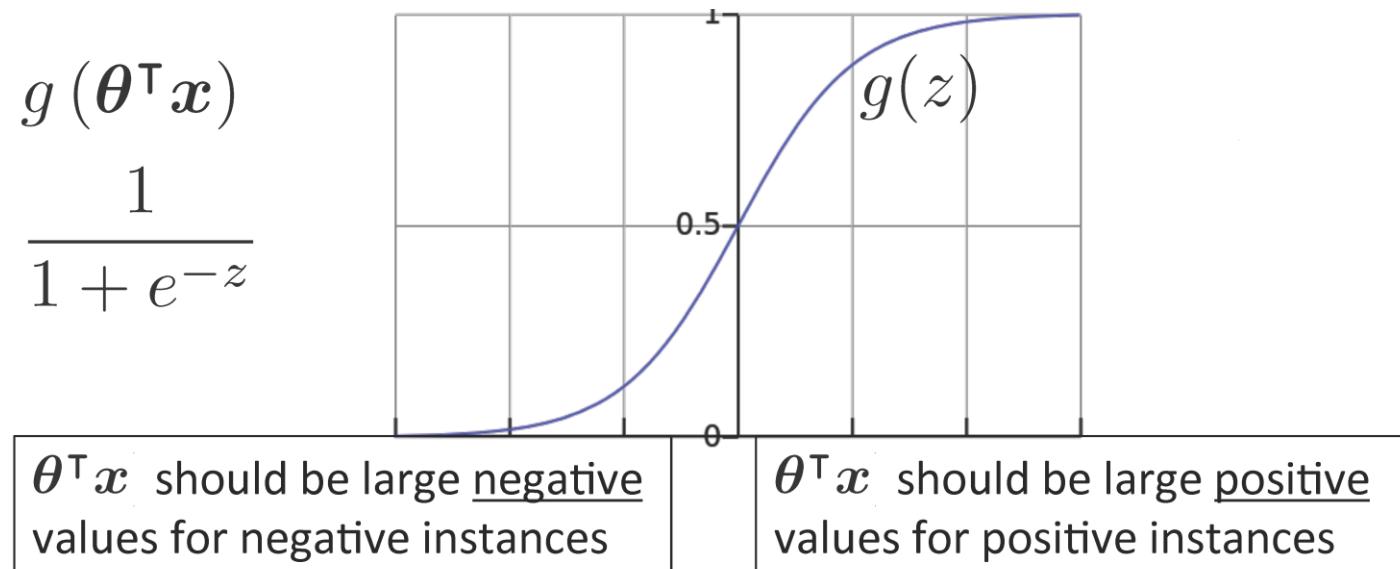
**Side Note:** the odds in favor of an event is the quantity  $p / (1 - p)$ , where  $p$  is the probability of the event

E.g., If I toss a fair dice, what are the odds that I will have a 6?

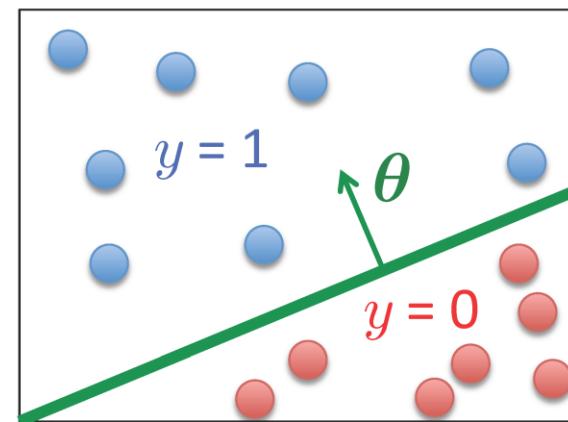
# Logistic Regression

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1 + e^{-z}}$$



- Assume a threshold and...
  - Predict  $y = 1$  if  $h_{\theta}(x) \geq 0.5$
  - Predict  $y = 0$  if  $h_{\theta}(x) < 0.5$



# Logistic Regression

- Given  $\left\{ \left( \mathbf{x}^{(1)}, y^{(1)} \right), \left( \mathbf{x}^{(2)}, y^{(2)} \right), \dots, \left( \mathbf{x}^{(n)}, y^{(n)} \right) \right\}$   
where  $\mathbf{x}^{(i)} \in \mathbb{R}^d$ ,  $y^{(i)} \in \{0, 1\}$
- Model:  $h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^\top \mathbf{x})$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \quad \mathbf{x}^\top = [1 \quad x_1 \quad \dots \quad x_d]$$

# Logistic Regression Objective Function

- Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

- Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^\top \mathbf{x}}}$$

results in a non-convex optimization

# Intuition Behind the Objective (log loss)

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- Cost of a single instance:

$$\text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

- Can re-write objective function as

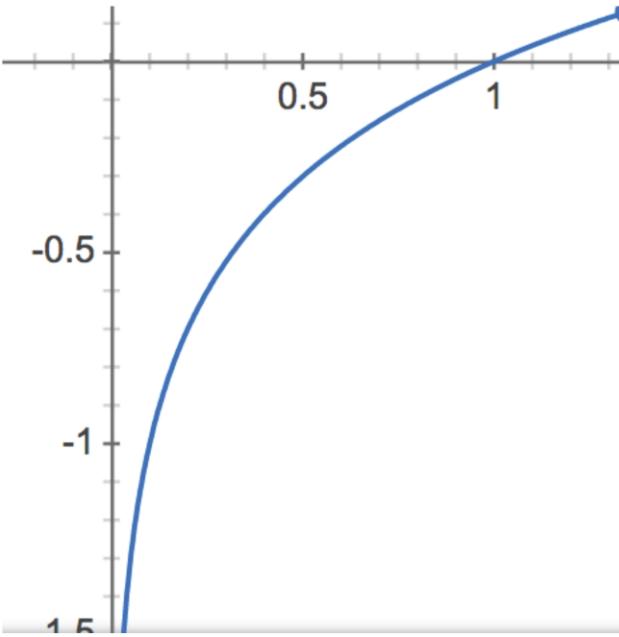
$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \text{cost}\left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)}\right)$$

Compare to linear regression:  $J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$

# Intuition Behind the Objective

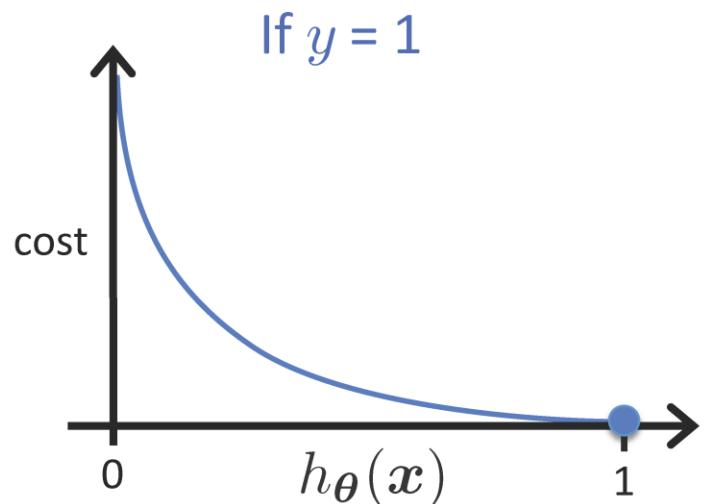
$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

Aside: Recall the plot of  $\log(z)$



# Intuition Behind the Objective

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

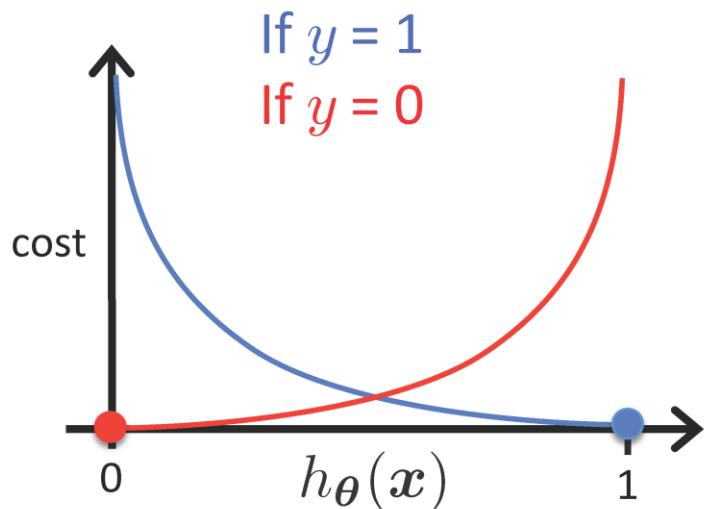


If  $y = 1$

- Cost = 0 if prediction is correct
- As  $h_{\theta}(\mathbf{x}) \rightarrow 0$ , cost  $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties
  - e.g., predict  $h_{\theta}(\mathbf{x}) = 0$ , but  $y = 1$

# Intuition Behind the Objective

$$\text{cost}(h_{\theta}(x), y) = \begin{cases} -\log(h_{\theta}(x)) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(x)) & \text{if } y = 0 \end{cases}$$



If  $y = 0$

- Cost = 0 if prediction is correct
- As  $(1 - h_{\theta}(x)) \rightarrow 0$ , cost  $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties

# Regularized Logistic Regression

$$J(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

- We can regularize logistic regression exactly as before:

$$\begin{aligned} J_{\text{regularized}}(\boldsymbol{\theta}) &= J(\boldsymbol{\theta}) + \frac{\lambda}{2} \sum_{j=1}^d \theta_j^2 \\ &= J(\boldsymbol{\theta}) + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2 \end{aligned}$$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize  $\boldsymbol{\theta}$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta}) \quad \begin{matrix} \text{simultaneous update} \\ \text{for } j = 0 \dots d \end{matrix}$$

Use the natural logarithm ( $\ln = \log_e$ ) to cancel with the  $\exp()$  in  $h_{\boldsymbol{\theta}}(\mathbf{x})$

# Gradient Descent for Logistic Regression

$$J_{\text{reg}}(\boldsymbol{\theta}) = -\frac{1}{n} \sum_{i=1}^n \left[ y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2} \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

Want  $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize  $\boldsymbol{\theta}$
- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \alpha \lambda \theta_j$$

<https://stats.stackexchange.com/questions/278771/how-is-the-cost-function-from-logistic-regression-derived>

# Gradient Descent for Logistic Regression

- Initialize  $\theta$

- Repeat until convergence (simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \frac{1}{n} \sum_{i=1}^n \left( h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)} \right) x_j^{(i)} - \alpha \lambda \theta_j$$

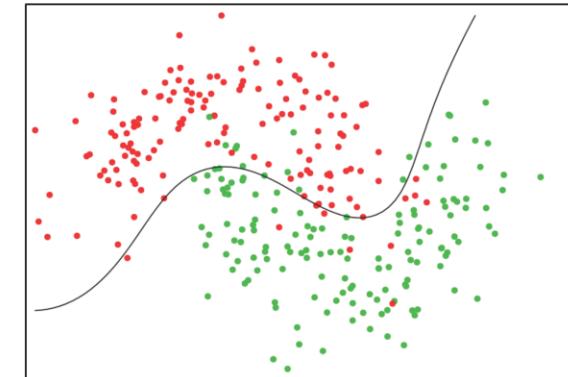
This looks IDENTICAL to linear regression!!!

- However, the form of the model is very different:

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}$$

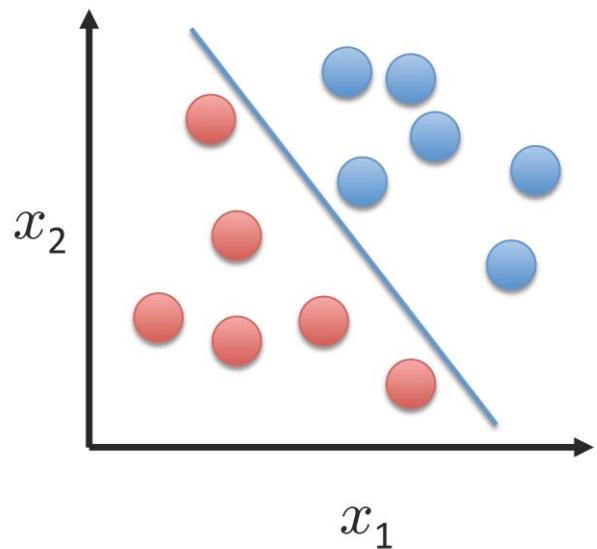
# Non-Linear Decision Boundary

- Can apply basis function expansion to features, same as with linear regression

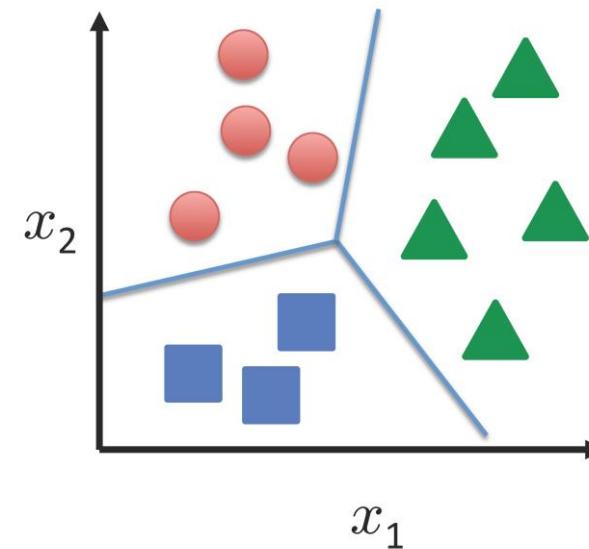
$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ x_2 \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ \vdots \end{bmatrix}$$


# Multi-Class Classification

Binary classification:



Multi-class classification:

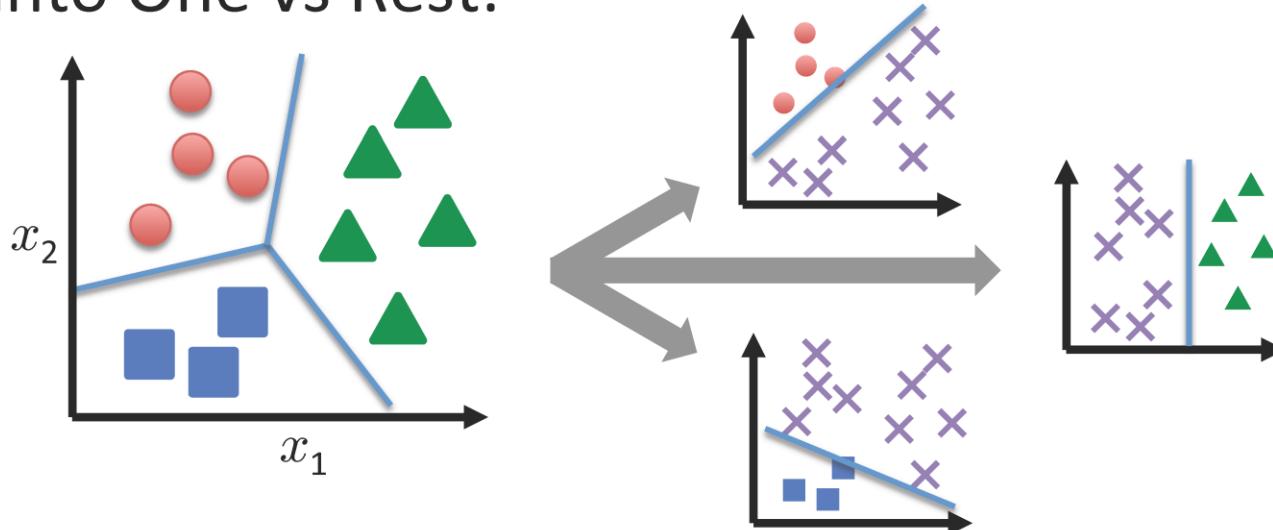


Disease diagnosis: healthy / cold / flu / pneumonia

Object classification: desk / chair / monitor / bookcase

# Multi-Class Logistic Regression

Split into One vs Rest:



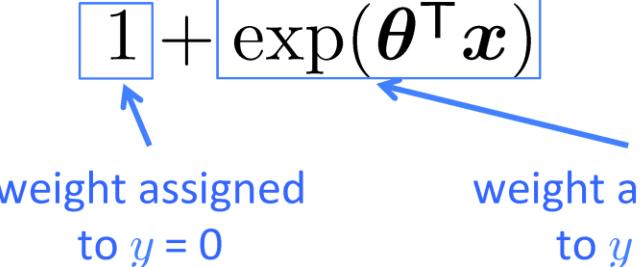
- Train a logistic regression classifier for each class  $i$  to predict the probability that  $y = i$  with

$$h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^\top \mathbf{x})}$$

# Multi-Class Logistic Regression

- For 2 classes:

$$h_{\theta}(x) = \frac{1}{1 + \exp(-\theta^T x)} = \frac{\exp(\theta^T x)}{1 + \exp(\theta^T x)}$$

  
weight assigned to  $y = 0$                                     weight assigned to  $y = 1$

- For  $C$  classes  $\{1, \dots, C\}$ :

$$p(y = c \mid x; \theta_1, \dots, \theta_C) = \frac{\exp(\theta_c^T x)}{\sum_{c=1}^C \exp(\theta_c^T x)}$$

– Called the **softmax** function (normalized exponential)

# Implementing Multi-Class Logistic Regression

- Use  $h_c(\mathbf{x}) = \frac{\exp(\boldsymbol{\theta}_c^\top \mathbf{x})}{\sum_{c=1}^C \exp(\boldsymbol{\theta}_c^\top \mathbf{x})}$  as the model for class  $c$
- Gradient descent simultaneously updates all parameters for all models
  - Same derivative as before, just with the above  $h_c(\mathbf{x})$
- Predict class label as the most probable label

$$\max_c h_c(\mathbf{x})$$

# Advanced Evaluation Metrics

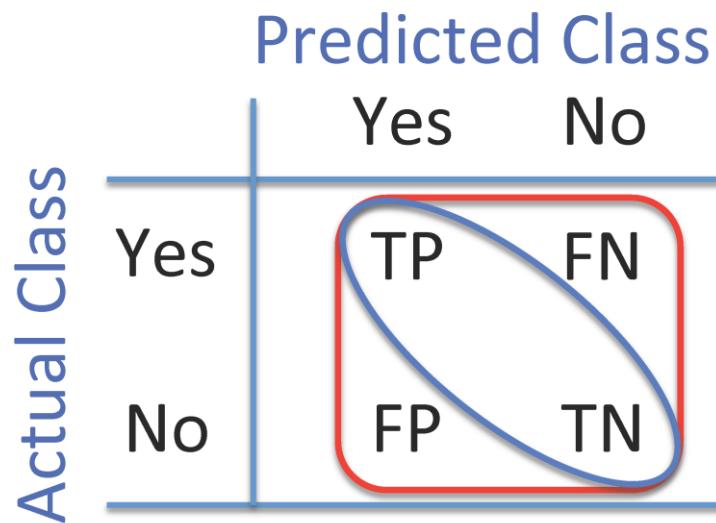
# Confusion Matrix

Given a dataset of  $P$  positive instances and  $N$  negative instances:

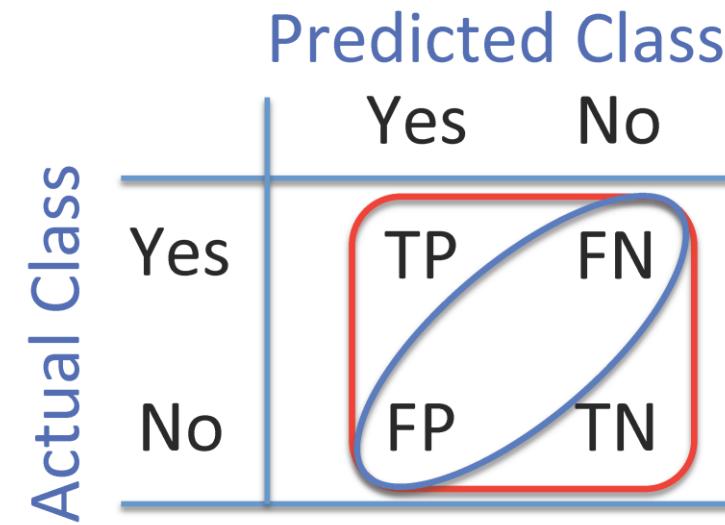
		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

# Accuracy & Error

Given a dataset of  $P$  positive instances and  $N$  negative instances:



$$\text{accuracy} = \frac{TP + TN}{P + N}$$



$$\begin{aligned}\text{error} &= 1 - \frac{TP + TN}{P + N} \\ &= \frac{FP + FN}{P + N}\end{aligned}$$

# Precision & Recall

## Precision

- the fraction of positive predictions that are correct
- $P(\text{is pos} | \text{predicted pos})$

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

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

## Recall

- fraction of positive instances that are identified
- $P(\text{predicted pos} | \text{is pos})$

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

		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN

# Precision & Recall

## Precision

- the fraction of positive predictions that are correct
- $P(\text{is pos} \mid \text{predicted pos})$

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

## Recall

- fraction of positive instances that are identified
- $P(\text{predicted pos} \mid \text{is pos})$

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

- 
- You can get high recall (but low precision) by only predicting positive
  - Recall is a non-decreasing function of the # positive predictions
  - Typically, precision decreases as either the number of positive predictions or recall increases
  - Precision & recall are widely used in information retrieval

# F-Measure

- Combined measure of precision/recall tradeoff

$$F_1 = 2 \times \frac{\text{precision} \times \text{recall}}{\text{precision} + \text{recall}}$$

- This is the harmonic mean of precision and recall
- In the  $F_1$  measure, precision and recall are weighted evenly
- Can also have biased weightings that emphasize either precision or recall more ( $F_2 = 2 \times \text{recall}$ ;  $F_{0.5} = 2 \times \text{precision}$ )

# A Word of Caution

- Consider binary classifiers A, B, C:

		A	.	B	.	C	.
		1	0	1	0	1	0
Predictions	1	0.9	0.1	0.8	0	0.78	0
	0	0	0	0.1	0.1	0.12	0.1

- Clearly A is useless, since it always predicts 1
- B is slightly better than C
- But, here are the performance metrics:

Metric	A	B	C
Accuracy	0.9	0.9	0.88
Precision	0.9	1.0	1.0
Recall	1.0	0.888	0.8667
F-score	0.947	0.941	0.9286

# Receiver Operating Characteristic (ROC)

ROC curves assess predictive behavior independent of error costs or class distributions

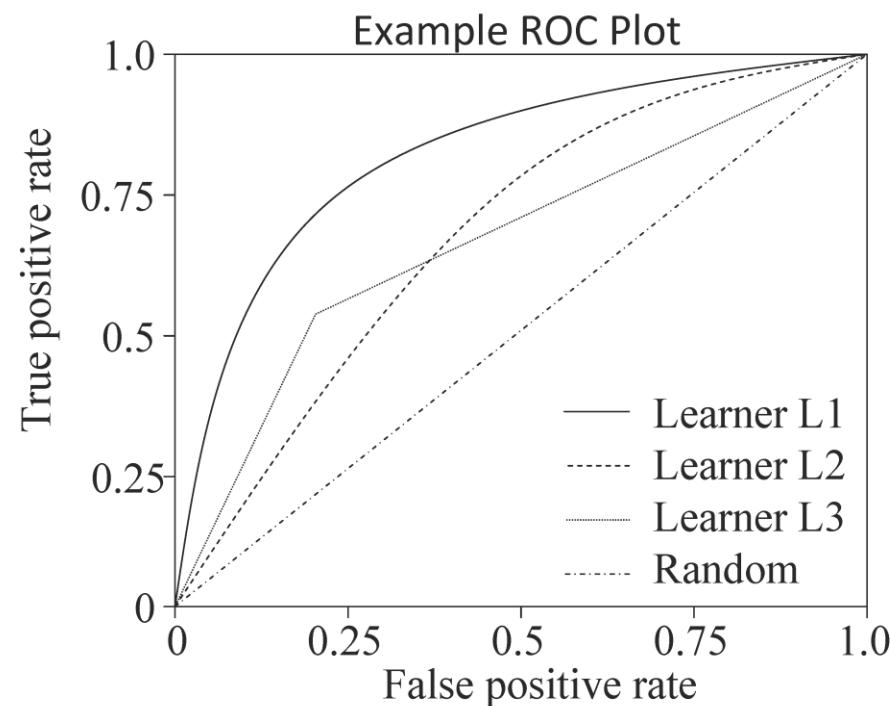
- Originated from signal detection theory
- Common in medical diagnosis, now used for ML

Plots TP rate vs FP Rate

$$\text{TP rate} = \frac{\text{TP}}{\text{P}}$$

$$\text{FP rate} = \frac{\text{FP}}{\text{N}}$$

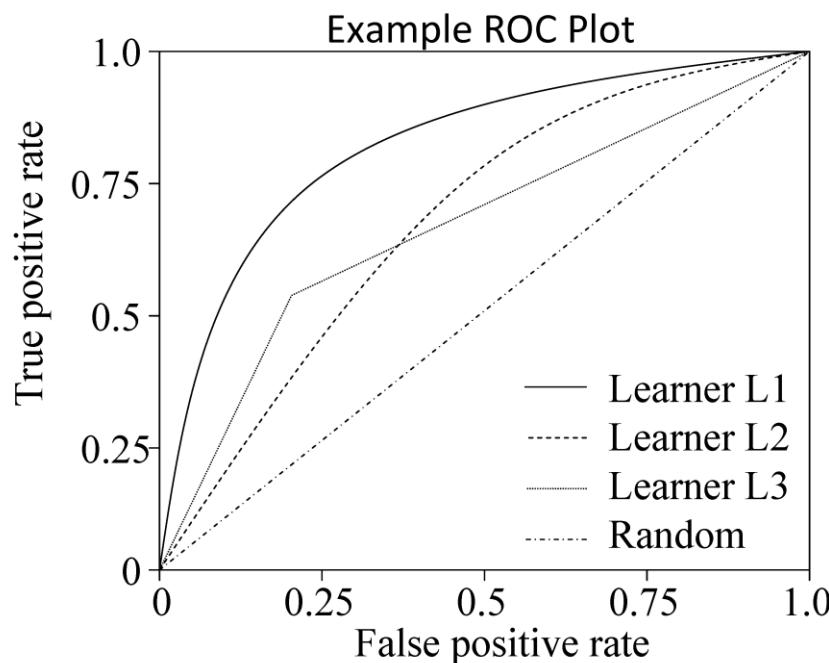
		Predicted Class	
		Yes	No
Actual Class	Yes	TP	FN
	No	FP	TN



# Performance Depends on Threshold

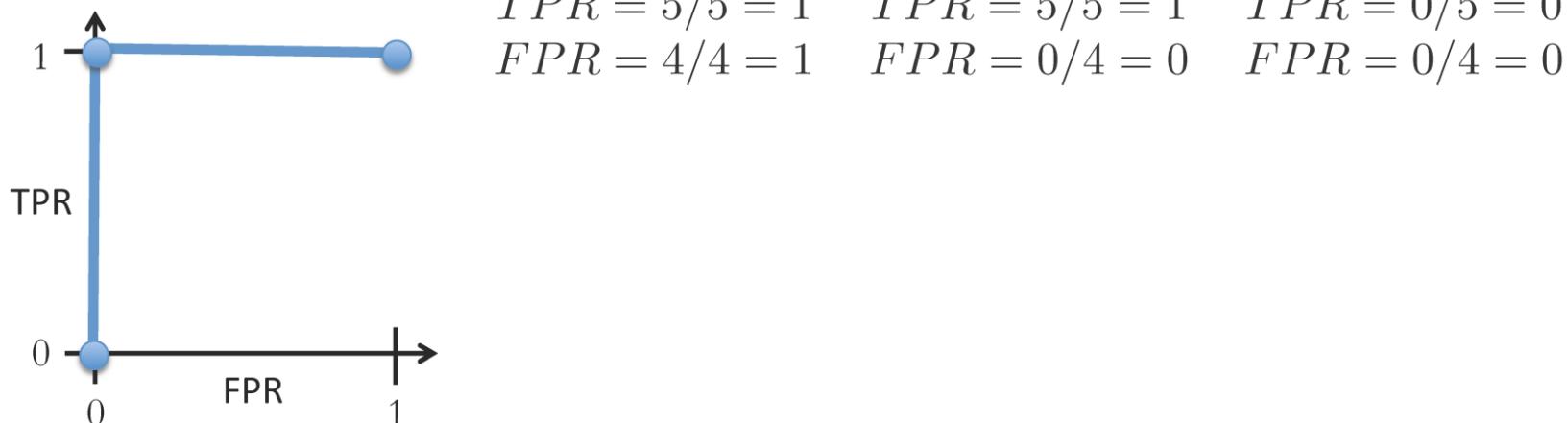
Predict positive if  $P(y = 1 \mid x) \geq \theta$ , otherwise negative

- Number of TPs and FPs depend on threshold  $\theta$
- As we vary  $\theta$ , we get different (TPR, FPR) points



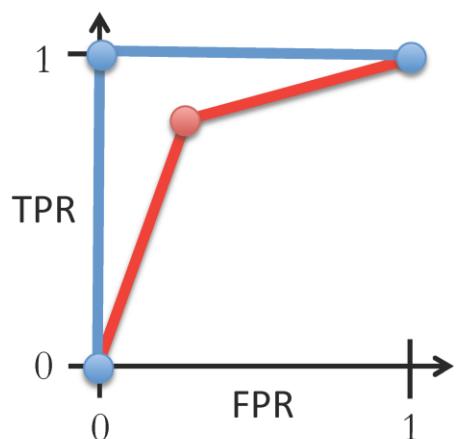
# ROC Example

$i$	$y_i$	$p(y_i = 1 \mid \mathbf{x}_i)$	$h(\mathbf{x}_i \mid \theta = 0)$	$h(\mathbf{x}_i \mid \theta = 0.5)$	$h(\mathbf{x}_i \mid \theta = 1)$
1	1	0.9	1	1	0
2	1	0.8	1	1	0
3	1	0.7	1	1	0
4	1	0.6	1	1	0
5	1	0.5	1	1	0
6	0	0.4	1	0	0
7	0	0.3	1	0	0
8	0	0.2	1	0	0
9	0	0.1	1	0	0



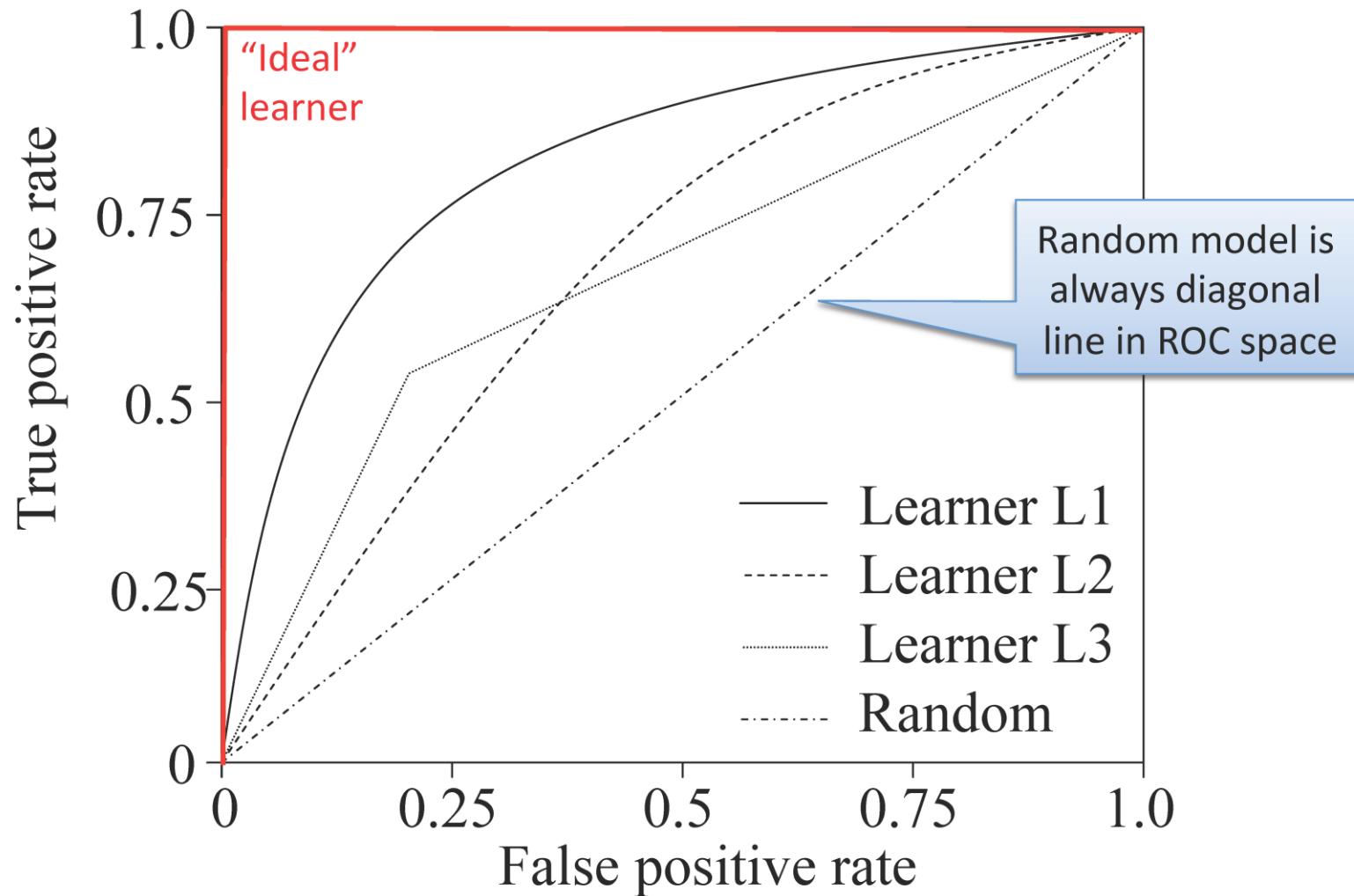
# ROC Example

$i$	$y_i$	$p(y_i = 1 \mid \mathbf{x}_i)$	$h(\mathbf{x}_i \mid \theta = 0)$	$h(\mathbf{x}_i \mid \theta = 0.5)$	$h(\mathbf{x}_i \mid \theta = 1)$
1	1	0.9	1	1	0
2	1	0.8	1	1	0
3	1	0.7	1	1	0
4	1	0.6	1	1	0
5	1	<b>0.2</b>	1	<b>0</b>	0
6	0	<b>0.6</b>	1	<b>1</b>	0
7	0	0.3	1	0	0
8	0	0.2	1	0	0
9	0	0.1	1	0	0

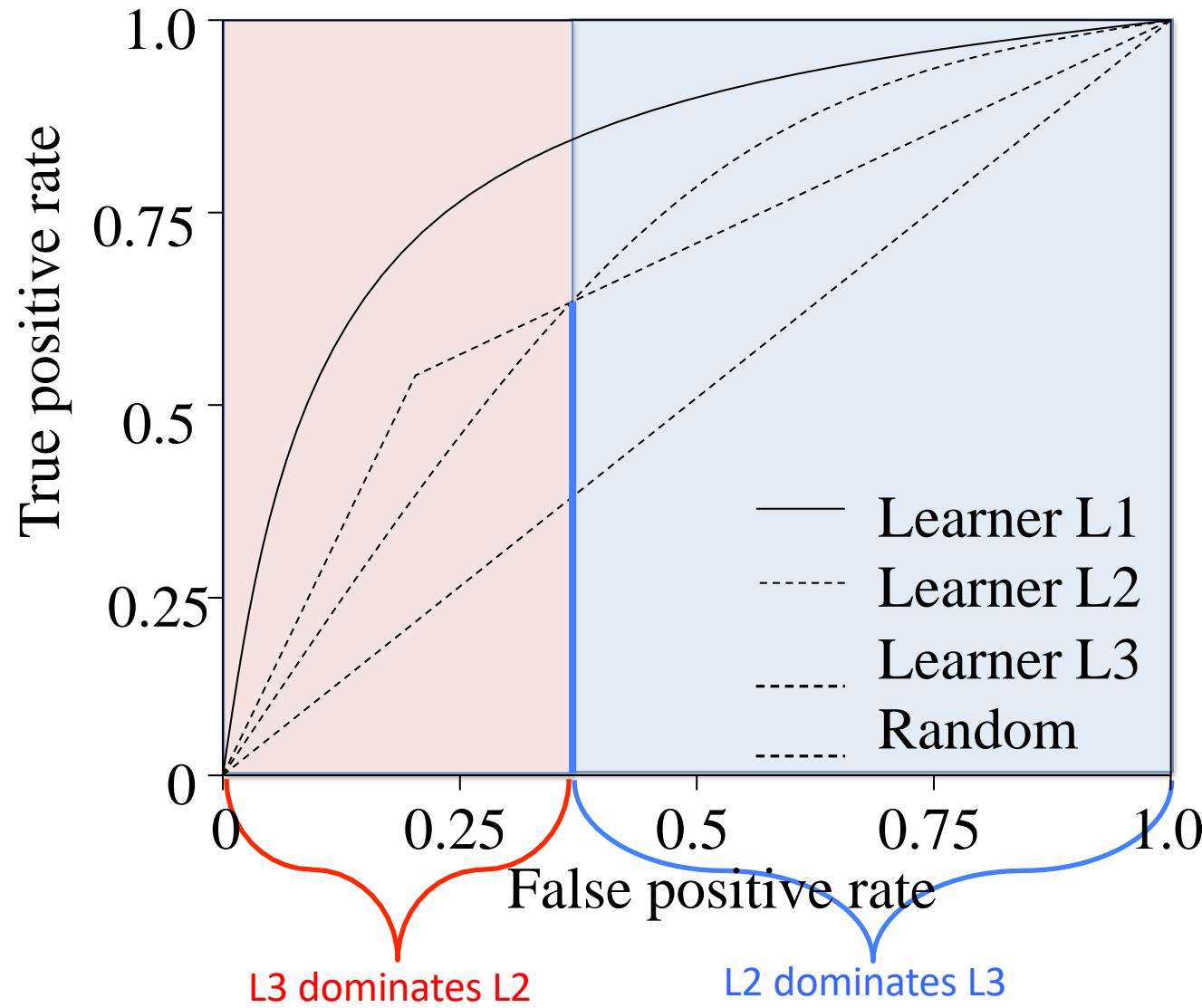


$$\begin{array}{lll} TPR = 5/5 = 1 & TPR = 4/5 = 0.8 & TPR = 0/5 = 0 \\ FPR = 4/4 = 1 & FPR = 1/4 = 0.25 & FPR = 0/4 = 0 \end{array}$$

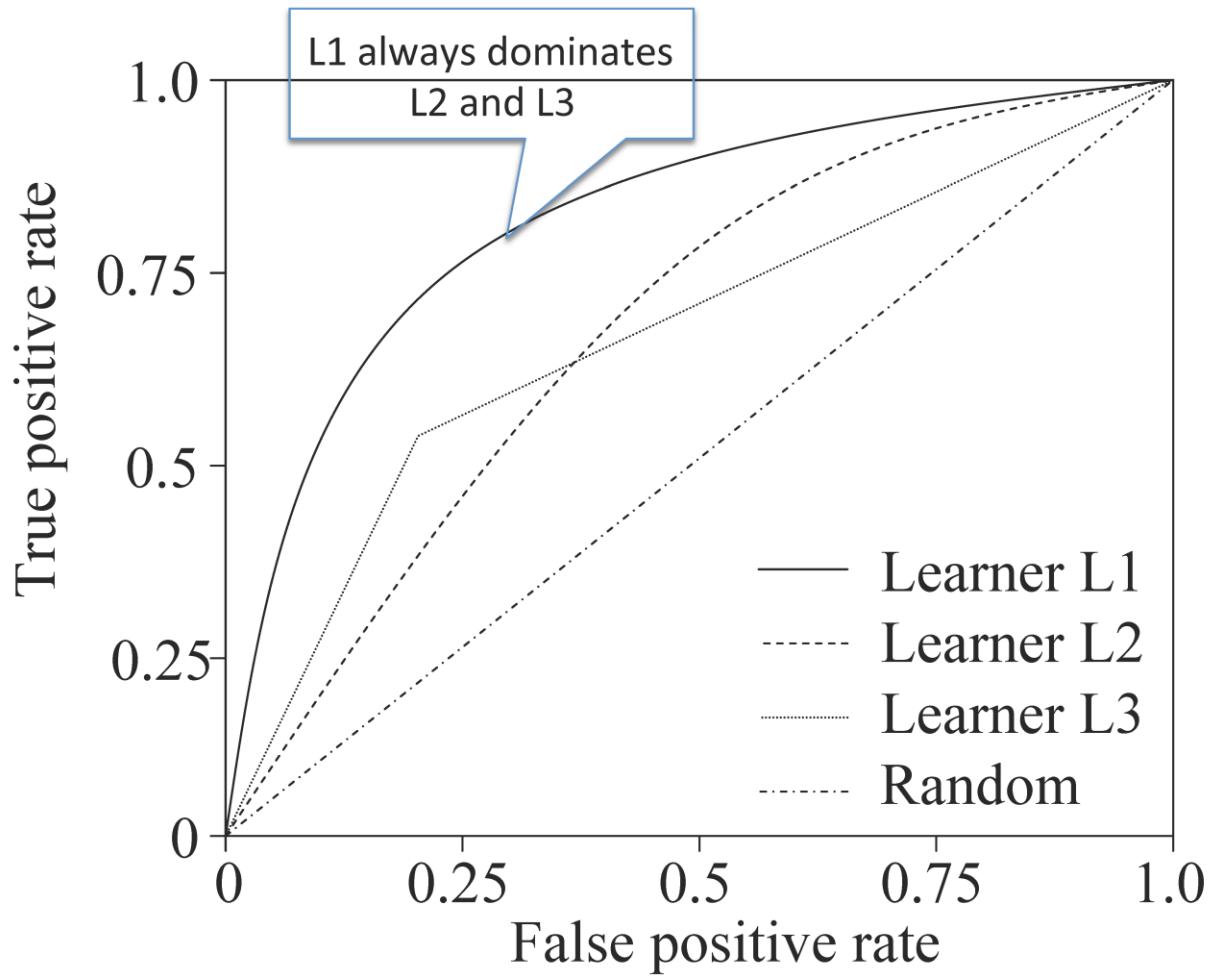
# Receiver Operating Characteristic (ROC)



# Receiver Operating Characteristic (ROC)



# Receiver Operating Characteristic (ROC)



# Area Under the ROC Curve (AUC)

- Can take area under the ROC curve to summarize performance as a single number

