

# DS 4400

## Machine Learning and Data Mining I Spring 2021

Alina Oprea

Associate Professor

Khoury College of Computer Science

Northeastern University

February 16 2021

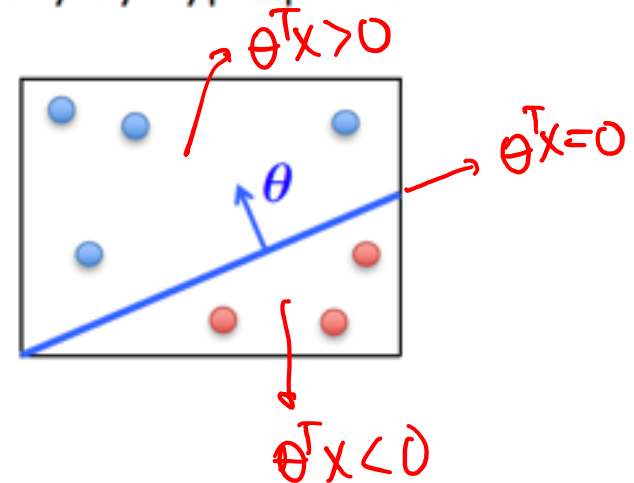
# Outline

- Logistic regression
  - Classification based on probability
- Maximum Likelihood Estimation
  - Application to logistic regression
  - Cross-entropy objective
- Gradient descent for logistic regression
- Logistic regression lab
- Evaluation metrics for classifiers

# Linear Classifiers

- **Linear classifiers:** represent decision boundary by hyperplane

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



$h_\theta(x) = f(\theta^T x)$  linear function

- If  $\theta^T x > 0$  classify "Class 1"
- If  $\theta^T x < 0$  classify "Class 0"

$$\theta_0 + \theta_1 x_1 + \dots + \theta_d x_d = 0$$

All the points  $x$  on the hyperplane satisfy:  $\theta^T x = 0$

# Logistic Regression

- Takes a probabilistic approach to learning discriminative functions (i.e., a classifier)

- $h_{\theta}(x)$  should give  $P(Y = 1|X; \theta)$

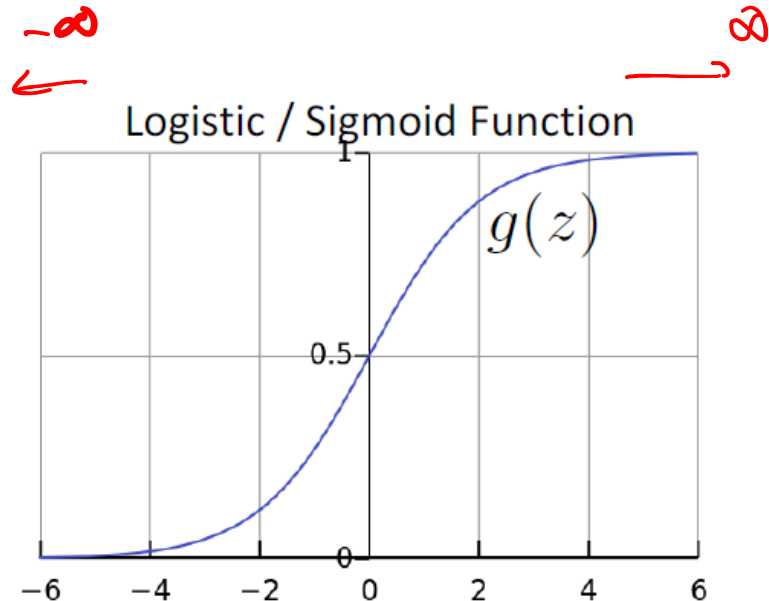
– Want  $0 \leq h_{\theta}(x) \leq 1$

- Logistic regression model:

$$h_{\theta}(x) = \overset{\text{SIGMOID}}{\hat{g}}(\theta^T x)$$

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

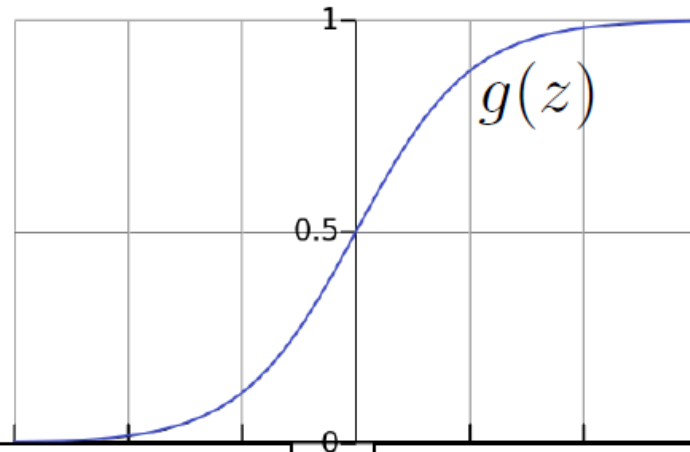
$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$



# Logistic Regression

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

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

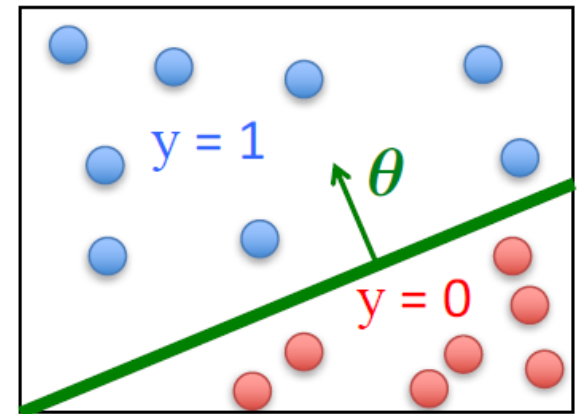


$\theta^T x$  should be large negative values for negative instances

$\theta^T x$  should be large positive values for positive instances

EQUIVALENT TO  $\theta^T x > 0$

- 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 is a linear classifier!

# Maximum Likelihood Estimation (MLE)

Given training data  $X = \{x_1, \dots, x_N\}$  with labels  $Y = \{y_1, \dots, y_N\}$  ,  $y_i \in \{0, 1\}$

What is the likelihood of training data for parameter  $\theta$ ?

Define likelihood function

$$\text{Max}_{\theta} \boxed{L(\theta) = P[Y|X; \theta]} \quad \text{DEF}$$

Assumption: training labels are conditionally independent

$$L(\theta) = \prod_{i=1}^N P[Y = y_i | X = x_i; \theta]$$

General probabilistic method for classifier training

# Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$\log(xy) = \log x + \log y$$

$$\rightarrow L(\theta) = \prod_{i=1}^N P[Y = y_i | X = x_i; \theta]$$

$$\rightarrow \log L(\theta) = \sum_{i=1}^N \log P[Y = y_i | X = x_i; \theta]$$

- They both have the same maximum  $\theta_{MLE}$

# MLE for Logistic Regression

$$P(Y = y_i | X = x_i; \theta) = h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}$$

,  $y_i \in \{0, 1\}$

1)  $y_i = 1$  ,  $P[Y=1 | X=x_i; \theta] = h_{\theta}(x_i)$  FROM DEF OF  $h_{\theta}$  FOR LOGISTIC REGRESSION

2)  $y_i = 0$  ,  $P[Y=0 | X=x_i; \theta] = 1 - h_{\theta}(x_i)$

$$\begin{aligned} \log L(\theta) &= \sum_{i=1}^N \log P[Y=y_i | X=x_i; \theta] = \\ &= \sum_{i=1}^N \log [h_{\theta}(x_i)^{y_i} (1 - h_{\theta}(x_i))^{1-y_i}] \\ &= \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1-y_i) \log (1 - h_{\theta}(x_i))] \end{aligned}$$

Find  $\theta$  to  $\max \log L(\theta)$  GIVEN  $\{x_i, y_i\}; i=1, N$



# How to Train Logistic Regression

1) DEFINE OBJ. / LOSS:  $J(\theta) = -\log L(\theta)$       CROSS-ENTROPY LOSS  
USE GRADIENT DESCENT ON  $J(\theta)$

2) GRADIENT ASCENT  
 $\max J(\theta)$

- 1) Init  $\theta$
- 2) Repeat until convergence
- 3)  $\theta \leftarrow \theta + \alpha \cdot \frac{\partial J(\theta)}{\partial \theta}$

# Cross-Entropy Objective

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

COST | LOSS FOR TRAINING  
EXAMPLE  $i$

# Cross-Entropy Objective

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- Cost of a single instance:

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

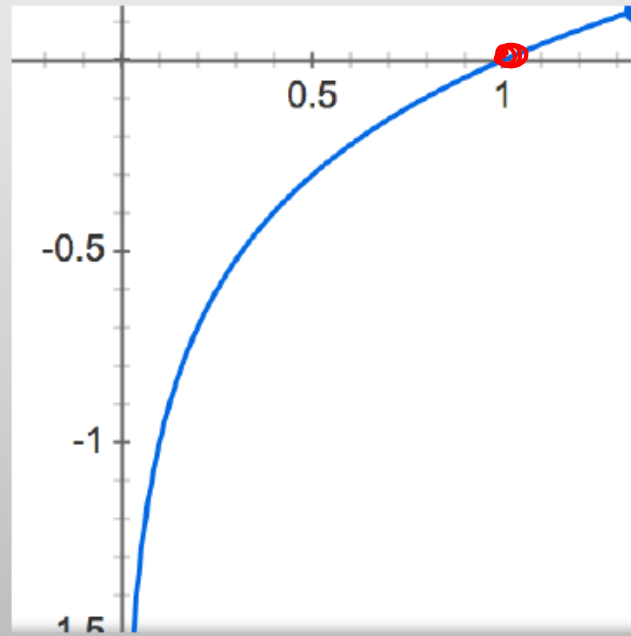
- Can re-write objective function as

$$J(\theta) = \sum_{i=1}^n \underbrace{\text{cost}(h_{\theta}(x_i), y_i)}_{\text{Cross-entropy loss}}$$

# Intuition

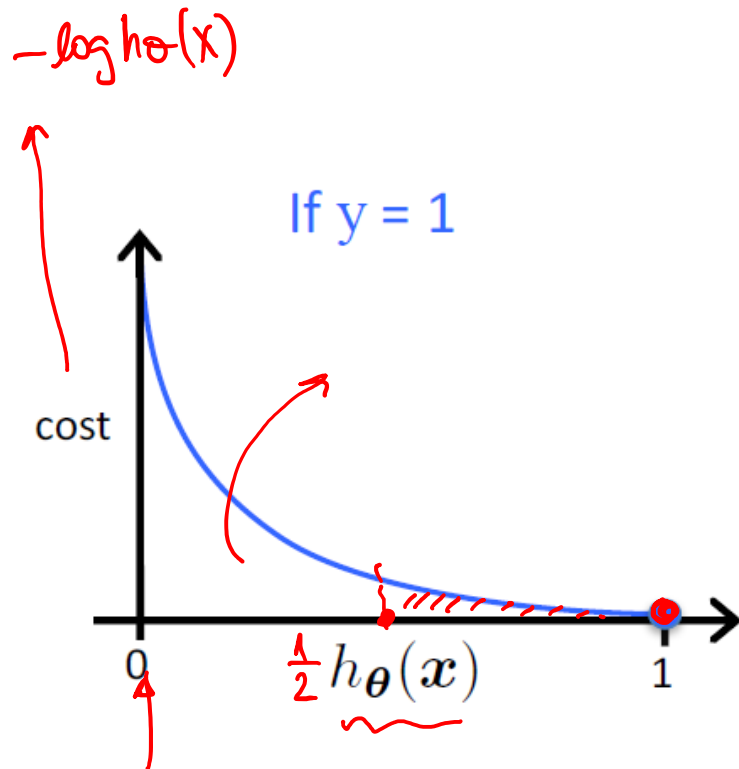
$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

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



# Intuition

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

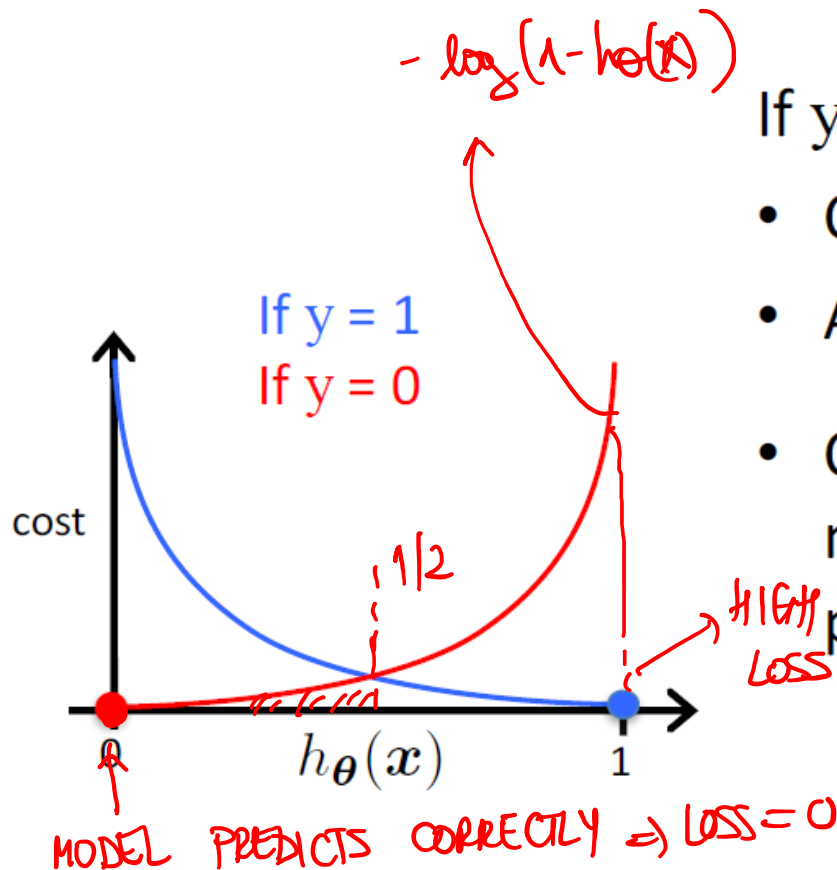


If  $y = 1$

- 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

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



If  $y = 0$

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

# Gradient Descent for Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))] \quad c_i(\theta)$$

$$J(\theta) = - \sum_{i=1}^N c_i(\theta)$$

FOR ALL  $j=0, \dots, d$

$$\frac{\partial J(\theta)}{\partial \theta_j} = - \sum_{i=1}^N \frac{\partial c_i(\theta)}{\partial \theta_j}$$

$$\begin{aligned} \frac{\partial h_{\theta}(x_i)}{\partial \theta_j} &= \frac{\partial g(\theta^T x_i)}{\partial \theta_j} = \\ &= g(\theta^T x_i) (1 - g(\theta^T x_i)) \cdot \boxed{\frac{\partial (\theta^T x_i)}{\partial \theta_j}} \\ &\quad \downarrow \\ &\quad x_{ij} \end{aligned}$$

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

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

$$g'(z) = + \frac{1}{(1 + e^{-z})^2} \cdot e^{-z}$$

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

CHAIN  
RULE

CALCULUS:

$$\left(\frac{1}{x}\right)' = -\frac{1}{x^2}$$

$$(\log x)' = \frac{1}{x}$$

$$(e^x)' = e^x$$

$$(e^{-x})' = -e^{-x}$$

# Gradient Computation

$$C_i(\theta) = y_i \log h_\theta(x_i) + (1 - y_i) \log (1 - h_\theta(x_i))$$

$$\frac{\partial h_\theta(x_i)}{\partial \theta_j} = g(\theta^T x_i) (1 - g(\theta^T x_i)) x_{ij}$$

$$\frac{\partial C_i(\theta)}{\partial \theta_j} = y_i \cdot \frac{1}{h_\theta(x_i)} \cdot \frac{\partial h_\theta(x_i)}{\partial \theta_j} - (1 - y_i) \frac{1}{1 - h_\theta(x_i)} \cdot \frac{\partial h_\theta(x_i)}{\partial \theta_j}$$

$$= y_i \cdot \frac{1}{h_\theta(x_i)} \cdot g(\theta^T x_i) (1 - g(\theta^T x_i)) x_{ij} - (1 - y_i) \frac{1}{1 - h_\theta(x_i)} g(\theta^T x_i) (1 - g(\theta^T x_i)) x_{ij}$$

$$\begin{aligned} &= [y_i (1 - g(\theta^T x_i)) - (1 - y_i) g(\theta^T x_i)] x_{ij} \\ &= [y_i - y_i g(\theta^T x_i) - g(\theta^T x_i) + y_i g(\theta^T x_i)] x_{ij} = \\ &= [y_i - h_\theta(x_i)] x_{ij} \end{aligned}$$



$$J(\theta) = - \sum_{i=1}^N c_i(\theta)$$

$$\frac{\partial J(\theta)}{\partial \theta_j} = - \sum_{i=1}^N \frac{\partial c_i(\theta)}{\partial \theta_j} = \sum_{i=1}^N [h_{\theta}(x_i) - y_i] x_{ij}$$

# Gradient Descent for Logistic Regression

CROSS-ENTROPY LOSS

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

$C_i(\theta)$

1) LINEAR REG,  $h_{\theta}(x_i) = \theta^T x_i$   
 2) LOGISTIC REG  
 $h_{\theta}(x_i) = \frac{1}{1 + e^{-\theta^T x_i}}$

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

- Initialize  $\theta$
- Repeat until convergence

(simultaneous update for  $j = 0 \dots d$ )

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i) x_{ij}$$

$\frac{\partial C_i(\theta)}{\partial \theta_j}$

$\frac{\partial J}{\partial \theta_j}$

,  $j = 1, \dots, d$

$x_{i0} = 1$

# Gradient Descent for Logistic Regression

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

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

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^N (h_{\theta}(x_i) - y_i)$$

$$\theta_j \leftarrow \theta_j - \alpha \sum_{i=1}^N \underbrace{(h_{\theta}(x_i) - y_i)}_{g(\theta^T x_i)} x_{ij}$$

**This looks IDENTICAL to Linear Regression!**

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

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

# Regularized Logistic Regression

$$J(\theta) = - \sum_{i=1}^N [y_i \log h_{\theta}(x_i) + (1 - y_i) \log (1 - h_{\theta}(x_i))]$$

- We can regularize logistic regression exactly as before:

$$J_{\text{regularized}}(\theta) = J(\theta) + \underbrace{\lambda \sum_{j=1}^d \theta_j^2}_{\text{REGULARIZATION TERM}}$$
$$= J(\theta) + \lambda \|\theta_{[1:d]}\|_2^2$$

LASSO REG: L2 regularization

$$J_{\text{LASSO}}(\theta) = J(\theta) + \lambda \underbrace{\sum_{j=1}^d |\theta_j|}_{\|\theta\|_1}$$

# Logistic Regression

## Lab Example

# Classifier Evaluation

- Classification is a supervised learning problem
  - Prediction is binary or multi-class
- Classification techniques
  - Linear classifiers
    - ~~Perceptron (online or batch mode)~~
    - Logistic regression (probabilistic interpretation)
  - Instance learners
    - kNN: need to store entire training data
- Cross-validation should be used for parameter selection and estimation of model error

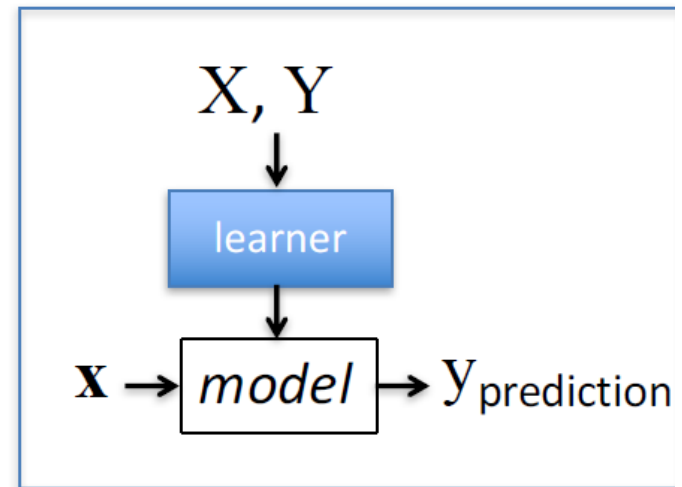
# Evaluation of classifiers

**Given:** labeled training data  $X, Y = \{\langle \mathbf{x}_i, y_i \rangle\}_{i=1}^n$

- Assumes each  $\mathbf{x}_i \sim \mathcal{D}(\mathcal{X})$

**Train the model:**

$model \leftarrow classifier.train(X, Y)$



**Apply the model to new data:**

- Given: new unlabeled instance  $x \sim \mathcal{D}(\mathcal{X})$

$y_{\text{prediction}} \leftarrow model.predict(\mathbf{x})$

VALIDATION  
TESTING

# Classification Metrics

$$\text{accuracy} = \frac{\# \text{ correct predictions}}{\# \text{ test instances}}$$

$$\text{error} = 1 - \text{accuracy} = \frac{\# \text{ incorrect predictions}}{\# \text{ test instances}}$$

- Training set accuracy and error
- Testing set accuracy and error



# 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

$$P = TP + FN$$

$$N = TN + FP$$

$$All = P + N$$

$$Accuracy = \frac{TP + TN}{P + N}$$

$$Error = \frac{FN + FP}{P + N}$$

# Review

- Maximum Likelihood Estimation (MLE) is a general statistical method for parameter estimation
- Logistic regression is a linear classifier that predicts class probability
  - Cross-entropy objective derived with MLE
- Logistic regression can be trained with Gradient Descent

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
- Thanks!