# DS 4400

# Machine Learning and Data Mining I
# Spring 2021

Alina Oprea

Associate Professor

Khoury College of Computer Science

Northeastern University

February 11 2021

# Outline

- Classification
  - K Nearest Neighbors (kNN)
- Cross validation
  - K-fold cross validation
  - Leave-one-out cross validation
- Linear classifiers
- Logistic regression

# Gradient Descent vs Closed Form

**Gradient Descent**

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

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update for j = 0 … d

**Closed form**

$$\theta = (X^\top X)^{-1} X^\top y$$

| • Gradient Descent | • Closed Form |
|---|---|
| + Linear increase in d and N | + No parameter tuning |
| + Generally applicable | + Gives the global optimum |
| - Need to choose $\alpha$ and stopping conditions | - Not generally applicable |
| - Might get stuck in local optima | - Slow computation |

# Ridge vs Lasso

- Both methods can be applied to any loss function (regression or classification)
- In both methods, value of regularization parameter $\lambda$ needs to be adjusted
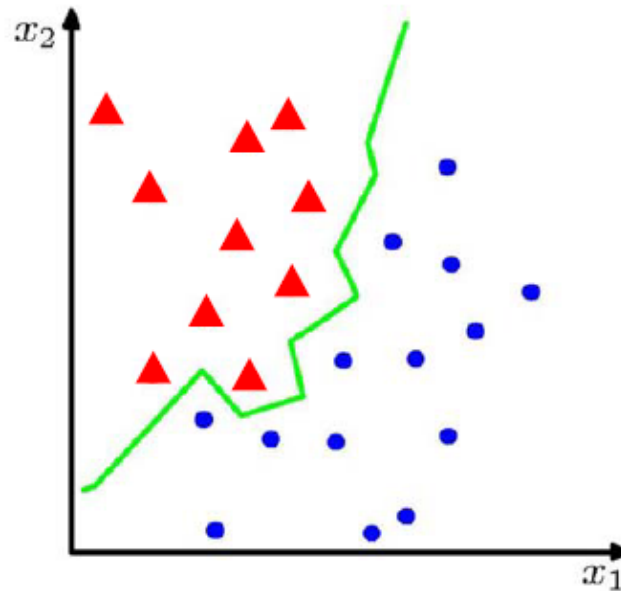- Both reduce model complexity

| - Ridge | - Lasso |
|---|---|
| + Differentiable objective | - Gradient descent needs to be adapted |
| + Gradient descent converges to global optimum | + Results in sparse model |
| - Shrinks all coefficients | + Can be used for feature selection in large dimensions |

# Classification



Binary or discrete
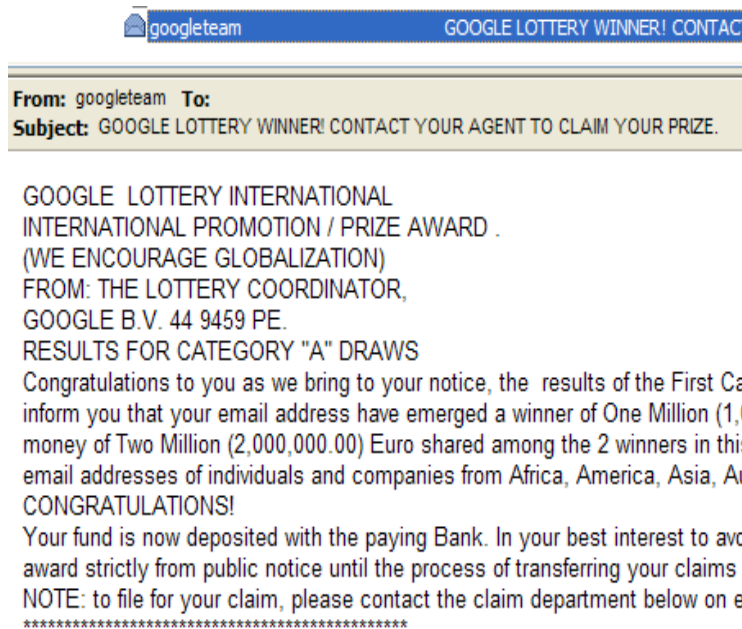
- Suppose we are given a training set of N observations

$$\{x_1, \dots, x_N\} \text{ and } \{y_1, \dots, y_N\}, x_i \in R^d, y_i \in \{0, 1\}$$

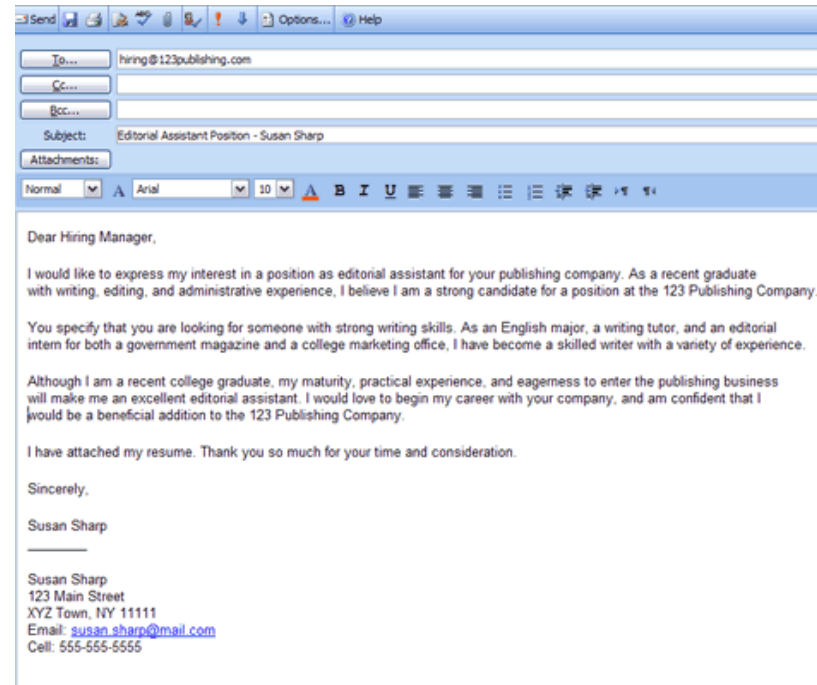- Classification problem is to estimate f(x) from this data such that

$$f(x_i) = y_i$$

# Example 1: Binary classification

## Classifying spam email





**Content-related features**
- Use of certain words
- Word frequencies
- Language
- Sentence

**Structural features**
- Sender IP address
- IP blacklist
- DNS information
- Email server
- URL links (non-matching)

## Binary classification: SPAM or HAM

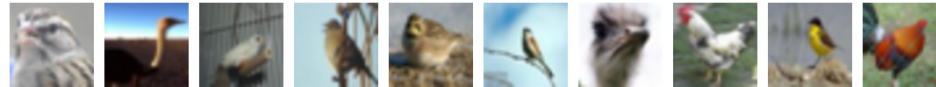# Example 2: Multi-class classification

Image classification



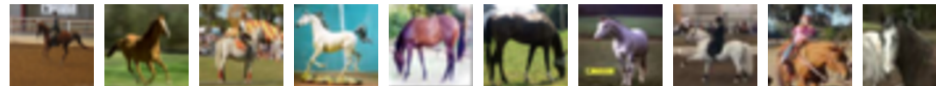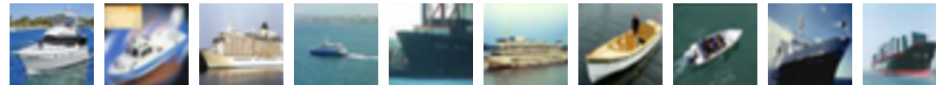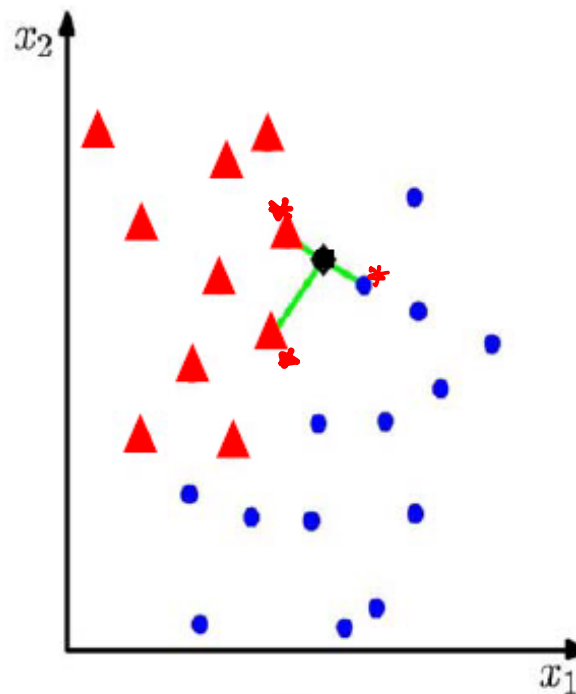| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| airplane | | | | | | | | | | |
| automobile | | | | | | | | | | |
| bird | | | | | | | | | | |
| cat | | | | | | | | | | |
| deer | | | | | | | | | | |
| dog | | | | | | | | | | |
| frog | | | | | | | | | | |
| horse | | | | | | | | | | |
| ship | | | | | | | | | | |
| truck | | | | | | | | | | |

Multi-class classification

# K Nearest Neighbour (K-NN) Classifier

## Algorithm

- For each test point, x, to be classified, find the K nearest samples in the training data

- Classify the point, x, according to the majority vote of their class labels

e.g. K = 3

• applicable to multi-class case

# Distance Metrics

$\text{DIST}(x,y) = \text{NORM}(x-y)$

- Euclidean Distance

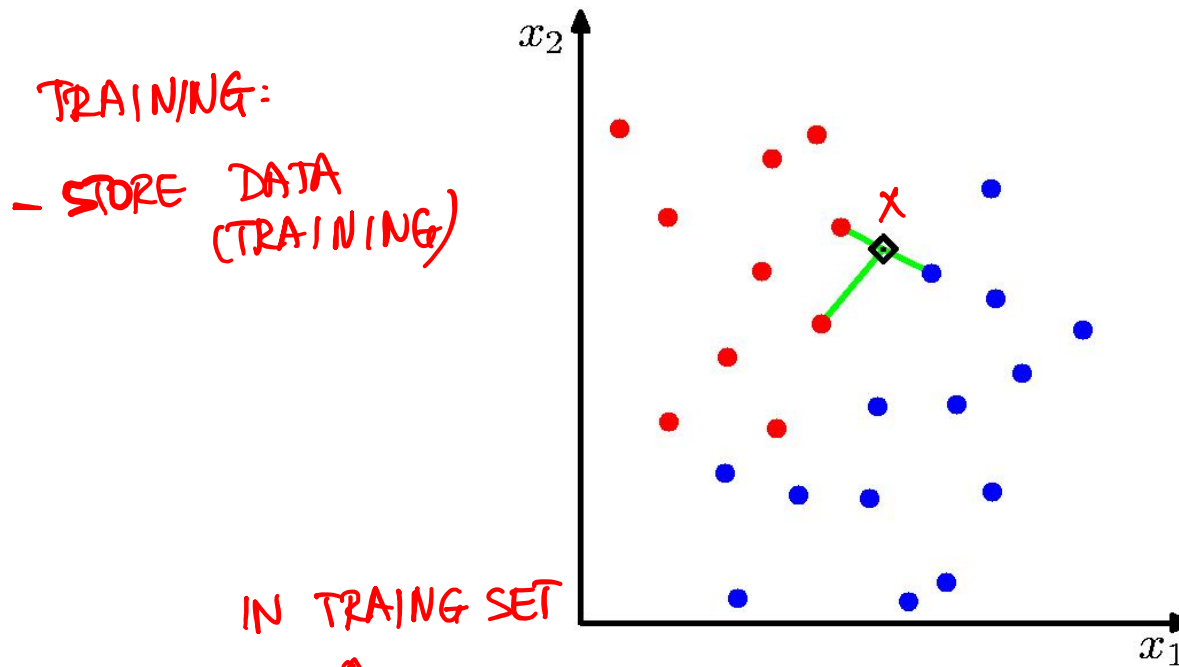$$\sqrt{\left(\sum_{i=1}^{k}(x_i - y_i)^2\right)} \quad L_2$$

- Manhattan Distance

$$\sum_{i=1}^{k}|x_i - y_i| \quad L_1$$

- Minkowski Distance

$$\left(\sum_{i=1}^{k}(|x_i - y_i|)^q\right)^{\frac{1}{q}} \quad L_q$$

# kNN



TRAINING:

- STORE DATA (TRAINING)

IN TRAING SET

- Algorithm (to classify point $x$)   AT TESTING
  - Find $k$ nearest points to $x$ (according to distance metric)
  - Perform majority voting to predict class of $x$
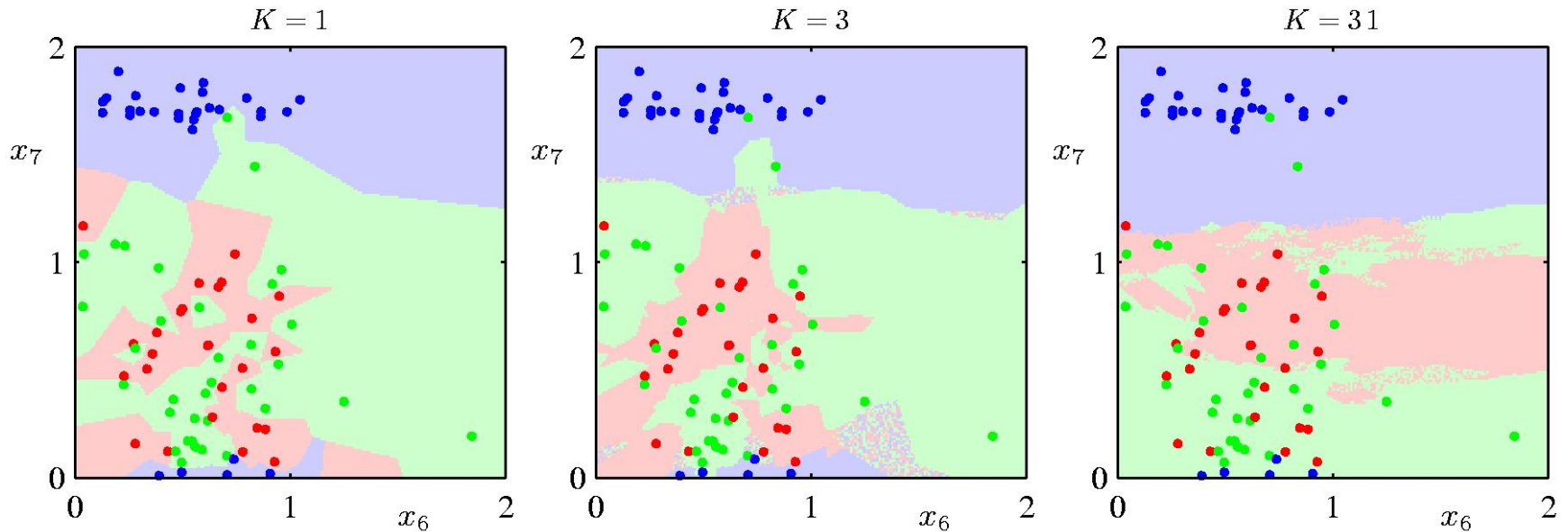
PROS:
  - SIMPLICITY

CONS:
  - COMPUTATIONAL COST
  - FEATURES TREATED UNIFORMLY
  - INSTANCE LEARNER

10

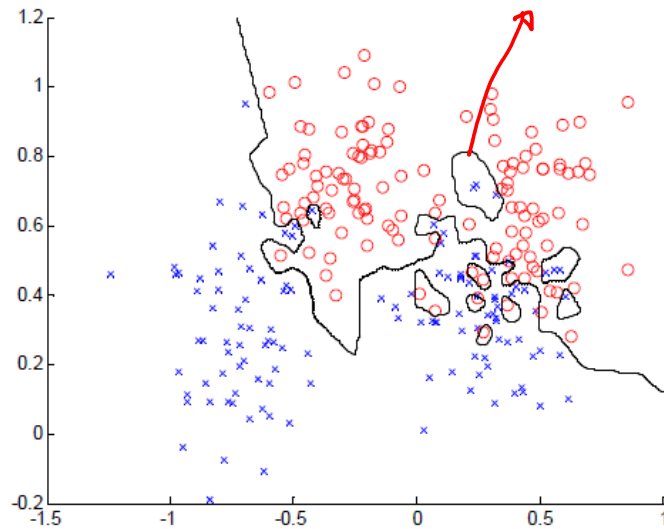# K-Nearest-Neighbours for Multi-class Classification



Vote among multiple classes
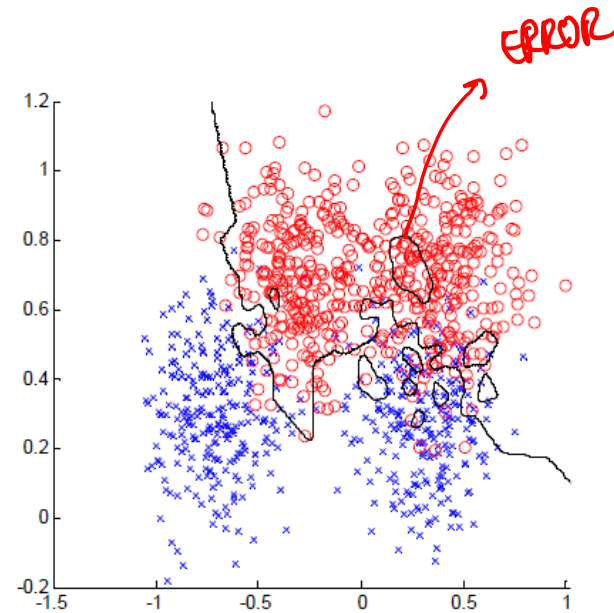
# K = 1    OVERFIT

## Training data



error = 0.0

## Testing data

ERROR


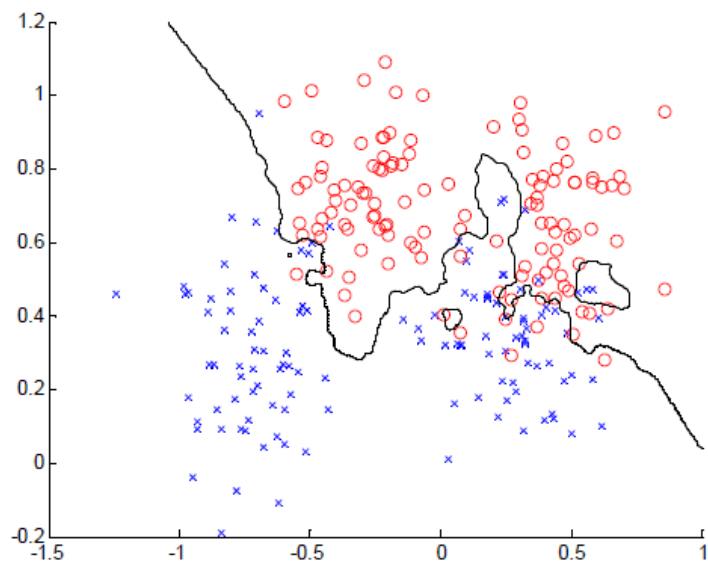
error = 0.15

How to choose k (hyper-parameter)?

# K = 3

Training data



Testing data



error = 0.0760

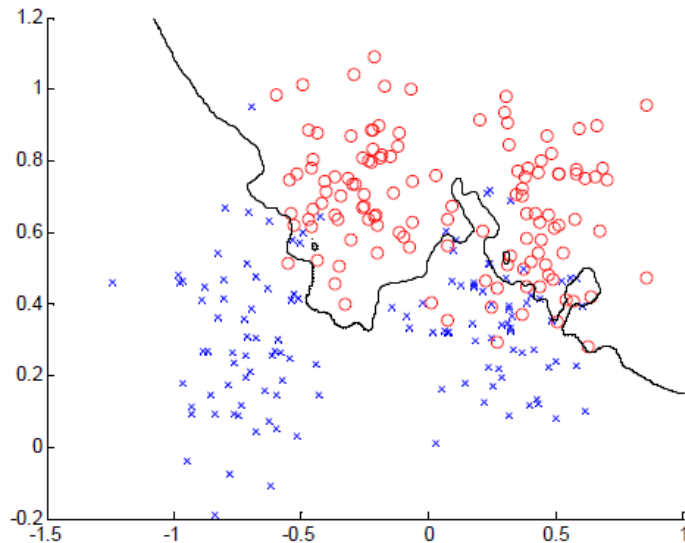error = 0.1340

How to choose k (hyper-parameter)?

**K = 7**

Training data

Testing data



↑ error = 0.1320

error = 0.1110 ↓

How to choose k (hyper-parameter)?

# Bias-Variance Tradeoff for kNN



k= N
UNDERFIT

k=1
OVERFIT

# How Overfitting Affects Prediction



How can we avoid over-fitting without having access to testing data?

# Cross Validation

## As K increases:

- Classification boundary becomes smoother
- Training error can increase

## Choose (learn) K by cross-validation

- Split training data into training and validation
- Hold out validation data and measure error on this

MEASURES ERROR

ML MODEL

| TRAIN | VAL. |

TR. DATA

CV METHOD

- SELECT FRACTION OF TRAINING DATA AT RANDOM R TRAIN ON THIS
- USE THE REST FOR VALIDATION

$K=1$

EXP1: | TRAIN 1 | VAL 1 |

EXP2:

EXP 10: | TRAIN 10 | VAL 10 |

$\downarrow$

AVG. ERROR
$K = 1$

$k=3$

$k=7$

PICK K THAT MIN AVG. VALIDATION ERROR

17

# Cross Validation



**Training Data**

**1st Partition**
Validation Set
Training Data

**2nd Partition**
Training Data
Validation Set
Training Data

**k$^{th}$ Partition**
Training Data
Validation Set

Test Data

Compute error metrics in each fold
Average error across folds

1. k-fold CV

– Split training data into k partitions (folds) of equal size

– Pick the optimal value of hyper-parameter according to error metric averaged over all folds

# Cross Validation

123 | n

123 | n
123 | n
123 | n
123 | n

*n EXP*

*TRAIN n-1 POINT*
*VAL 1 POINT*

## 2. Leave-one-out CV (LOOCV)

– k=n (validation set only one point)

- Pros: Less bias

- Cons: More expensive to implement, higher variance

- Recommendation: perform k-fold CV with k=5 or k=10

# Cross-Validation Takeaways

- General method to estimate performance of ML model at testing and select hyper-parameters
  - Improves model generalization
  - Avoids overfitting to training data
- Techniques for CV: k-fold CV and LOOCV
- Compare to regularization

– CV FOR TUNING HYPER-PARAMS

– LASSO & RIDGE ARE APPLICABLE WHEN OPT OBY

       – LINEAR REG

       – SVM

       – LOGISTIC REG.

# Cross Validation Applications

- PICK $\lambda$ IN REGULARIZATION ( LASSO & RIDGE)

      — COMBINING CV + REG.

- FIND DEGREE OF POLY IN POLY REGRESSION
- FIND NUMBER OF KNOTS IN SPLINE REG.

# Linear classifiers

- A **hyperplane** partitions space into 2 half-spaces
  - Defined by the normal vector $\theta \in \mathbb{R}^{d+1}$

MODEL PARAM

  - $\theta$ is orthogonal to any vector lying on the hyperplane

  - Assumed to pass through the origin
    - This is because we incorporated bias term $\theta_0$ into it by $x_0 = 1$

$\theta$ orthogonal $\Rightarrow$ $\theta^T x = 0$, $\forall x$ ON HYPERPLANE

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

- Consider classification with +1, -1 labels …

$f = \text{sign.}$     If $\theta^T x > 0 \Rightarrow$ CLASSIFY +1
$< 0$          $-1$

# Linear Classifiers

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

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_d \end{bmatrix} \qquad \boldsymbol{x}^{\mathsf{T}} = \begin{bmatrix} 1 & x_1 & \cdots & x_d \end{bmatrix}$$



$h_\theta(x) = f(\theta^T x)$ linear function
- If $\theta^T x > 0$ classify "Class 1"
- If $\theta^T x < 0$ classify "Class 0"

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

# Linear vs Non-Linear Classifiers



LOG. REGRESSION
LINEAR SVM
LDA
PERCEPTRON

class1
class2

KNN
DECISION TREES
KERNEL SVM
BOOSTING
RANDOM FOREST
NEURAL NETWORK

# Classification Based on Probability

- Instead of just predicting the class, give the *probability of the instance being in that class*

LEARN $\quad P[Y = 1 \mid X = x]$

- Consider binary classifier with classes 0 and 1

LEARN $\quad \boxed{P[Y=1 \mid X = x]} + P[Y=0 \mid X = x] = 1$

$$P[Y=0 \mid X = x] = 1 - P[Y=1 \mid X = x]$$

- Advantages: interpretability and confidence of output

# Logistic Regression

- Setup
  $$x_i \in \mathbb{R}^d$$
  - Training data: $\{x_i, y_i\}$, for $i = 1, \ldots, N$
  - Labels: $y_i \in \{0,1\}$
- Goals
  - Learn $P(Y = 1 | X = x) = h_\theta(x)$
- Highlights
  - Probabilistic output
  - At the basis of more complex models (e.g., neural networks)
  - Supports regularization (Ridge, Lasso)
  - Can be trained with Gradient Descent

# Interpretation of Model Output

$h_{\boldsymbol{\theta}}(\boldsymbol{x})$ = estimated $\qquad P(Y = 1 | X; \theta)$

Example: Cancer diagnosis from tumor size

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

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

→ Tell patient that 70% chance of tumor being malignant

Note that: $\qquad P(Y = 0 | X; \theta) + P(Y = 1 | X; \theta) = 1$

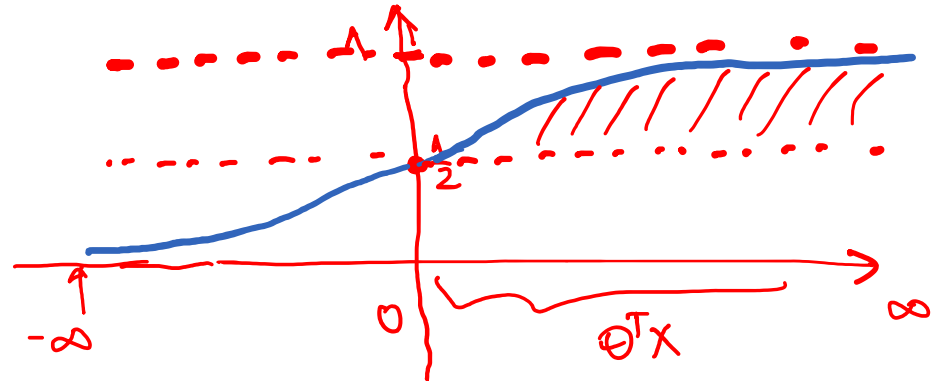Therefore, $\qquad P(Y = 0 | X; \theta) = 1 - P(Y = 1 | X; \theta)$

# 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 $0 \leq h_\theta(x) \leq 1$

$$h_\theta(x) = g(\theta^T x) \in [0, 1]$$

SIGMOID: $\quad g(z) = \dfrac{1}{1 + e^{-z}}$

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



28

# LR Predictions

$$h_\theta(x) = P[y=1 \mid X=x] = \frac{1}{1+e^{-\theta^T x}}$$

SIGMOID FUNCTION
IS NOT LINEAR

- Predict $Y = 1$ if:

$$P[y=1 \mid X=x] > \frac{1}{2}$$

$$\frac{1}{1+e^{-\theta^T x}} > \frac{1}{2}$$
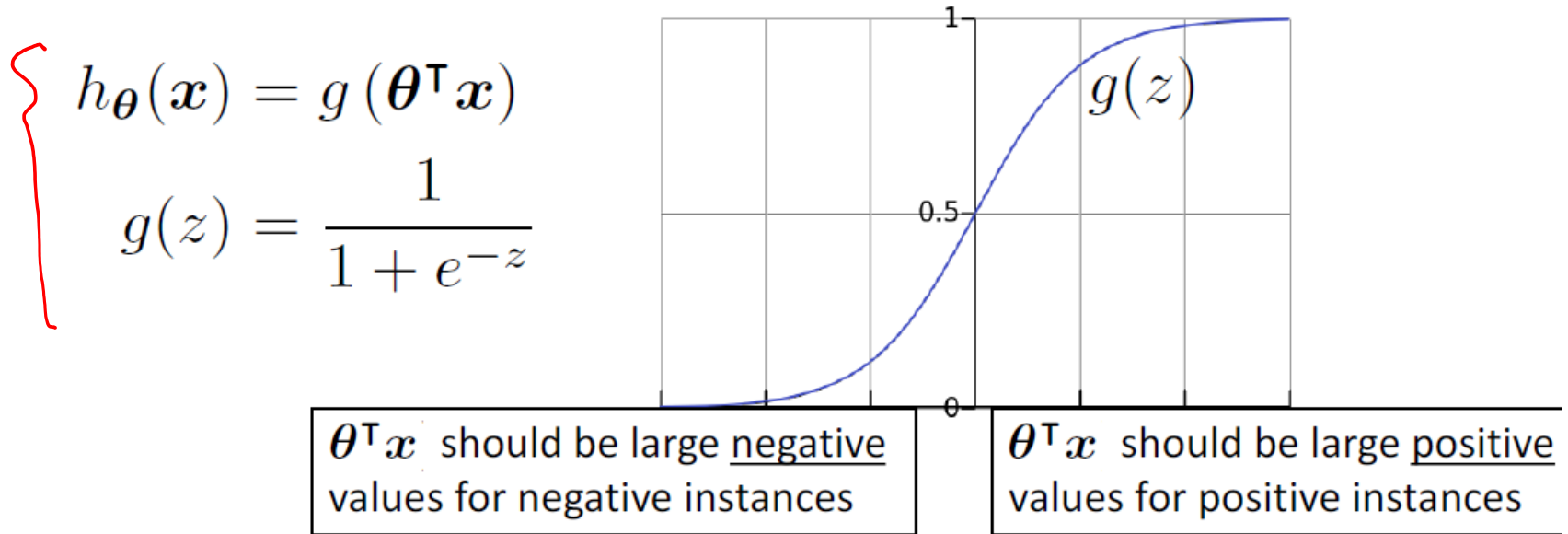
$$1 + e^{-\theta^T x} < 2$$

$$e^{-\theta^T x} < 1$$

$$e^{\theta^T x} > 1$$

$$\boxed{\theta^T x > 0}$$

LINEAR CLASSIFIER

# Logistic Regression

$$h_{\boldsymbol{\theta}}(\boldsymbol{x}) = g\left(\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}\right)$$

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

$g(z)$

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>negative</u> values for negative instances

$\boldsymbol{\theta}^{\mathsf{T}}\boldsymbol{x}$ should be large <u>positive</u> values for positive instances

- Assume a threshold and...   $\Rightarrow \theta^{\mathsf{T}}x > 0$

  - Predict $Y = 1$ if $\boxed{h_{\boldsymbol{\theta}}(\boldsymbol{x}) \geq 0.5}$

  - Predict $Y = 0$ if $h_{\boldsymbol{\theta}}(\boldsymbol{x}) < 0.5$

$\theta^{\mathsf{T}}x > 0$

y = 1

$\theta$

y = 0

$\theta^{\mathsf{T}}x < 0$

**Logistic Regression is a linear classifier!**

# How to Pick Loss Function?

1) MSE

$$J(\theta) = \frac{1}{N} \sum_{i=1}^{M} \left[ h_\theta(x_i) - y_i \right]^2$$

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

NOT CONVEX.

2) 0-1 LOSS:

$$J(\theta) = \frac{1}{M} \sum_{i=1}^{N} I\left[ h_\theta(x_i) \neq y_i \right]$$

$$= \begin{cases} 1, & h_\theta(x_i) \neq y_i \\ 0, & \text{OTHERWISE} \end{cases}$$

ERROR RATE

1) NOT DIFFERENTIABLE

2) DOES NOT MEASURE RATE OF ERROR

# Maximum Likelihood Estimation (MLE)

Given training data $X = \{x_1, \ldots, x_N\}$ with labels
$Y = \{y_1, \ldots, y_N\}$

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

Define likelihood function

$$L(\theta) = P\left[Y \mid X; \theta\right] = P\left[y_1, \ldots, y_N \mid x_1, \ldots, x_N; \theta\right]$$

Find $\theta$ to max $L(\theta)$

$$L(\theta) = \prod_{i=1}^{N} P\left(Y = y_i \mid X = x_i; \theta\right)$$

HOW LIKELY IS TRAINING DATA UNDER MODEL PARAM $\theta$

# Log Likelihood

- Max likelihood is equivalent to maximizing log of likelihood

$$\max \quad \log L(\theta) = \sum_{i=1}^{N} \log P\left[Y = y_i \mid X = x_i ; \theta\right]$$

$x_i \in \mathbb{R}^D$

TRAINING
EXAMPLE
(d- FEATURE)

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

LABELS

$y_i = 1$

$P[Y = 1 \mid X = x_i ; \theta]$

$\parallel$

$h_\theta(x_i)$

$y_i = 0$

$1 - h_\theta(x_i)$

$\log L(\theta) \rightarrow \text{DEPENDS ON} \quad \theta, x_i, y_i$

# Review

- K nearest neighbors is the first example of classifier
  - Instance learner
- Cross-validation should be performed to
  - Improve generalization and avoid over-fitting
  - Choose hyper parameters (k in kNN)
- Logistic regression is a linear classifier that predicts class probability