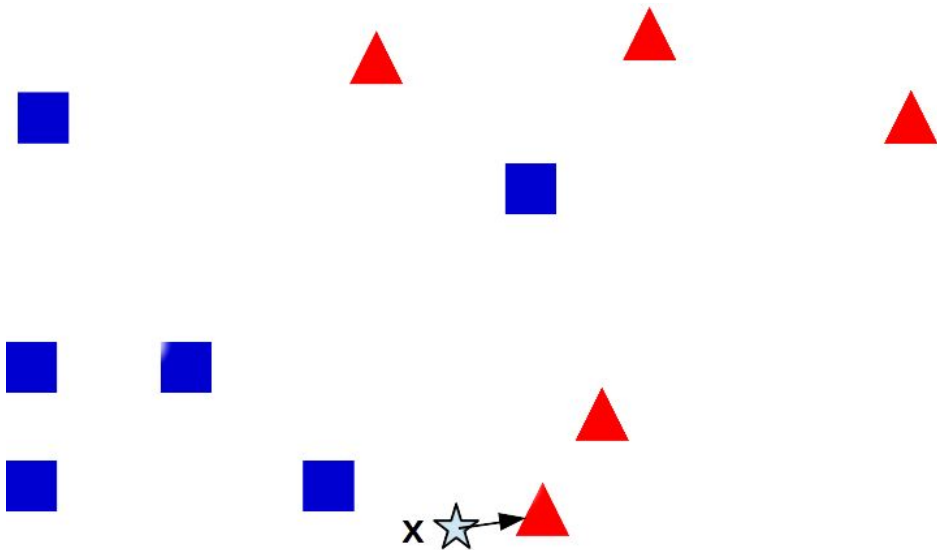


# SUPERVISED LEARNING

**K-Nearest Neighbour, Linear, Logistic...**

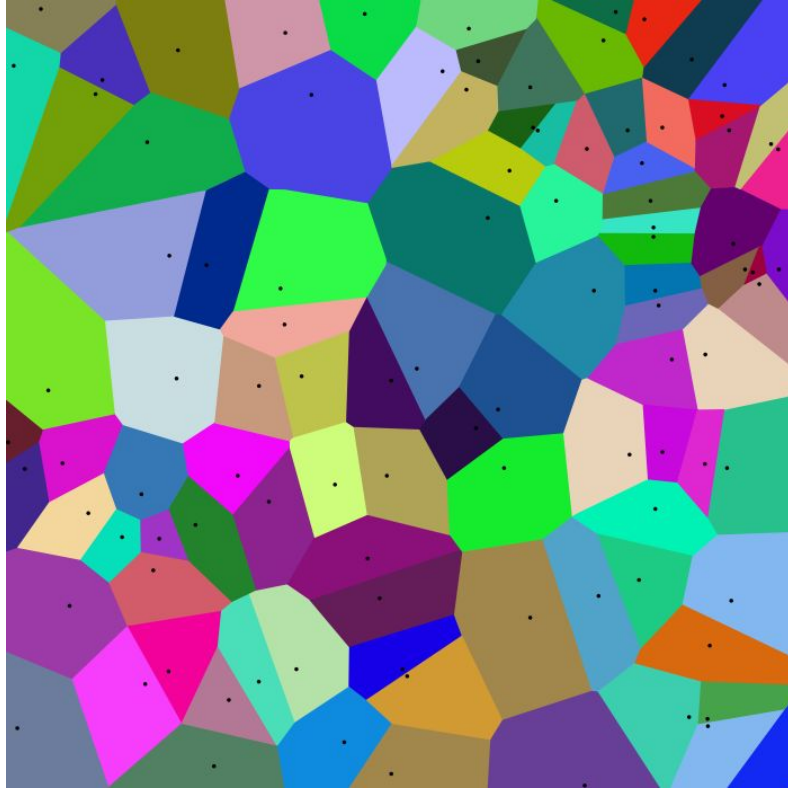
# LEARNING BY COMPUTING DISTANCES - KNN

- The label  $y$  for  $x \in \mathbb{R}^D$  will be label of its **nearest neighbour in training data**.
- Euclidean Distance can be used to find the nearest neighbour.
- NOTE : This method also applies to regression.



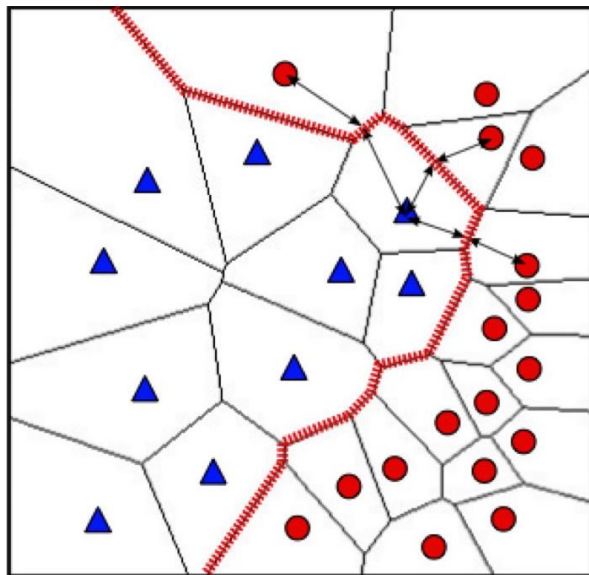
- A simple yet very effective method in practice (if given lots of training data).
- Also called a **memory-based** or **instance-based** or **non-parametric method**.
- No “model” is learned here. Prediction step uses all the training data.
- Requires lots of storage (need to keep all the training data at test time).
- Prediction can be slow at test time
  - For each test point, need to compute its distance from all the training points.
  - Clever data-structures or data-summarization techniques can provide speed-ups.
- Need to be careful in **choosing the distance function** to compute distances (especially when the data dimension  $D$  is very large).
- The 1-NN can suffer if data contains outliers (we will soon see a geometric illustration), or if amount of training data is small. Using more neighbors ( $K > 1$ ) is usually more robust.

# GEOMETRY OF 1-NN



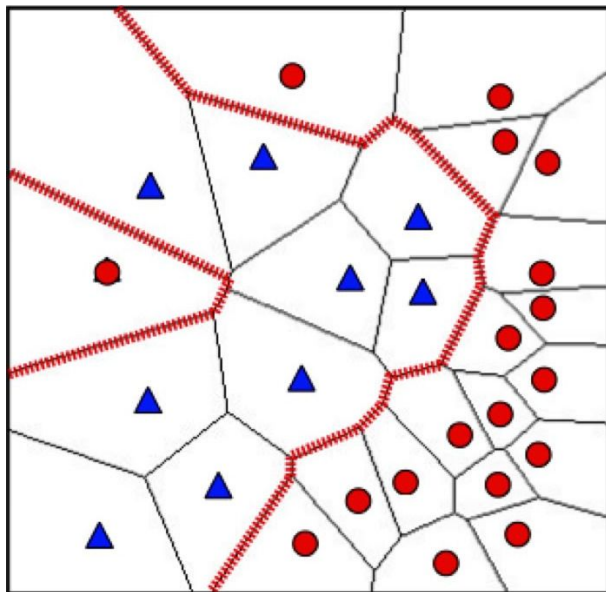
# THE DECISION BOUNDARY OF 1-NN

The decision boundary is composed of hyperplanes that form perpendicular bisectors of pairs of points from different classes



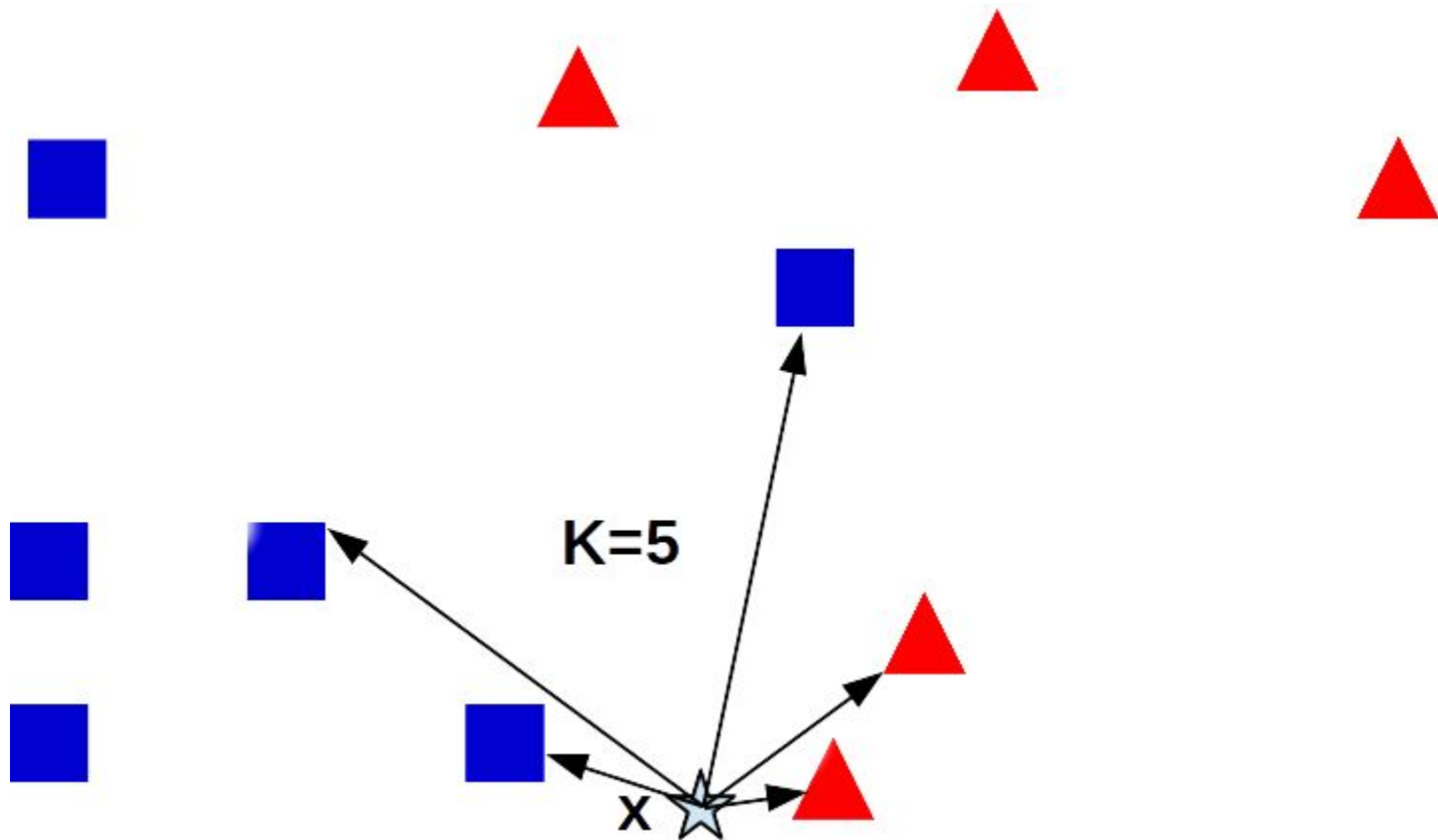
# EFFECT OF OUTLIERS ON 1-NN

An illustration of how the decision boundary can drastically change when the data contains some outliers



# K - NEAREST NEIGHBORS (K - NN)

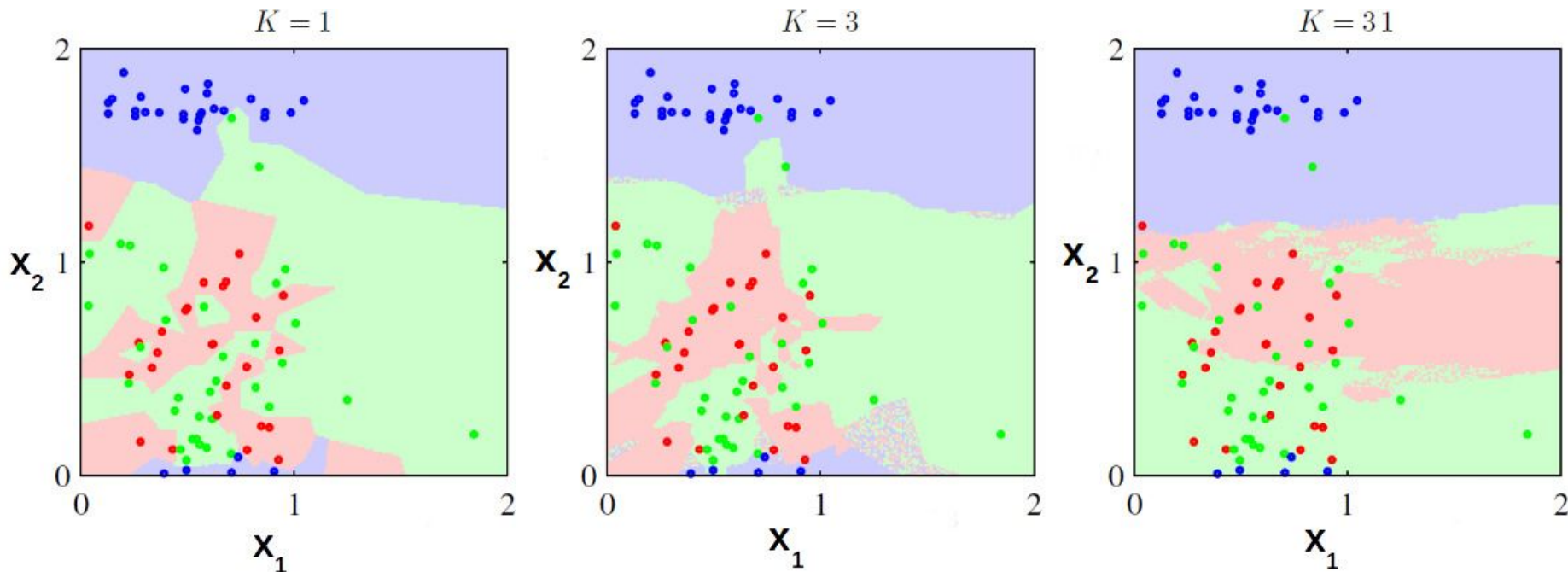
- Makes one-nearest-neighbor more robust by using more than one neighbor.
- The K -NN prediction rule: Take a majority vote (or average) of the labels of  $K > 1$  neighbors in the training data.
- Works for both classification and regression
  - For classification, we usually **take the majority** labels from the K neighbors
  - For regression, we usually **average** the real-valued labels of the K neighbors
- The “right” value of K needs to be selected (e.g., via cross-validation)





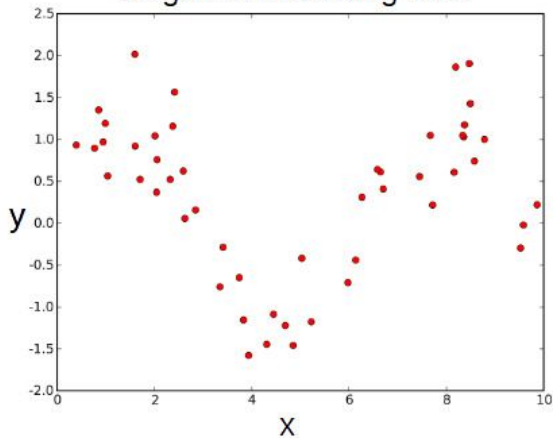
# K-NEAREST NEIGHBORS: DECISION BOUNDARIES

Larger K leads to smoother decision boundaries

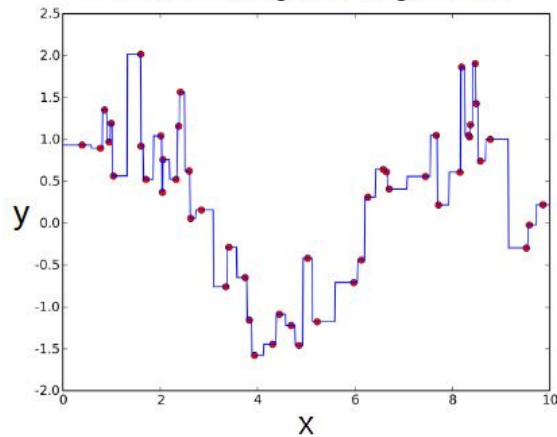


# K -NN BEHAVIOR FOR REGRESSION

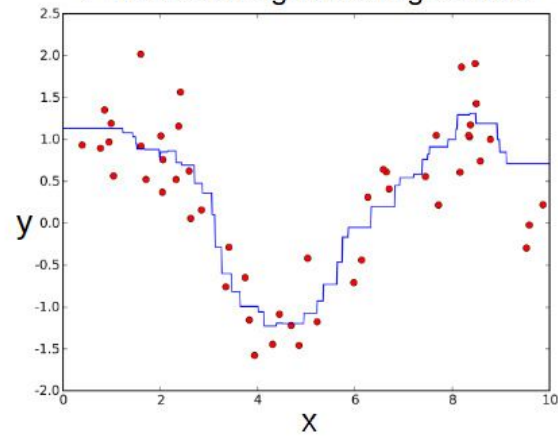
Regression training data



1-nearest neighbor regression



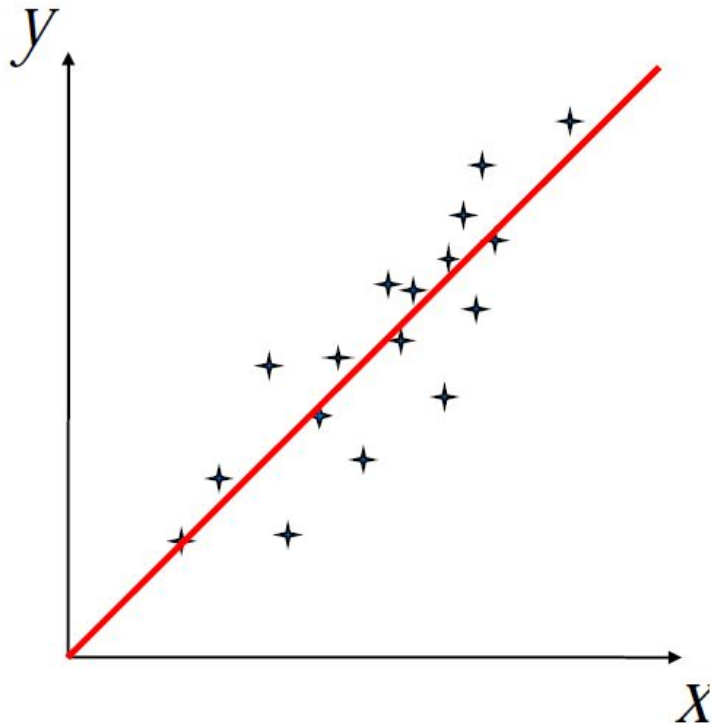
7-nearest neighbors regression



# LINEAR REGRESSION

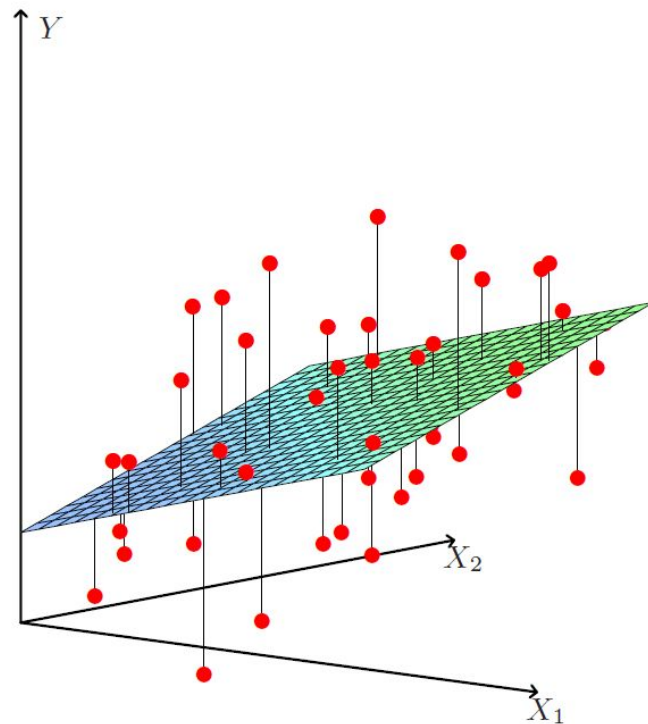
# FITTING A LINE TO THE DATA

- Let's assume the **relationship** between  $x$  and  $y$  to have a linear **model** :  $y = wx$
- Problem boils down to fitting a line to the data
- $w$  is the model parameter (slope of the line here)
- Many  $w$  's (i.e., many lines) can be fit to this data
- Which one is the best?



# FITTING A (HYPER)PLANE TO THE DATA

- For 2-dim. inputs, we can fit a 2-dim. plane to the data.
- In higher dimensions, we can likewise fit a **hyperplane**  $w \cdot x = 0$
- Defined by a D-dim vector  $w$  normal to the plane
- Many planes are possible. Which one is the best?



# LINEAR REGRESSION

**Given:** Training data with  $N$  examples  $\{(\mathbf{x}_n, y_n)\}_{n=1}^N$ ,  $\mathbf{x}_n \in \mathbb{R}^D$ ,  $y_n \in \mathbb{R}$

Assume the following linear model with model parameters  $\mathbf{w} \in \mathbb{R}^D$

$$y_n \approx \mathbf{w}^\top \mathbf{x}_n \quad \Rightarrow \quad y_n \approx \sum_{d=1}^D w_d x_{nd}$$

- The response  $y_n$  is a linear combination of the features of the inputs  $\mathbf{x}_n$
- $\mathbf{w} \in \mathbb{R}^D$  is also called the (regression) **weight vector**
  - Can think of  $w_d$  as weight/importance of  $d$ -th feature in the data

# LINEAR REGRESSION

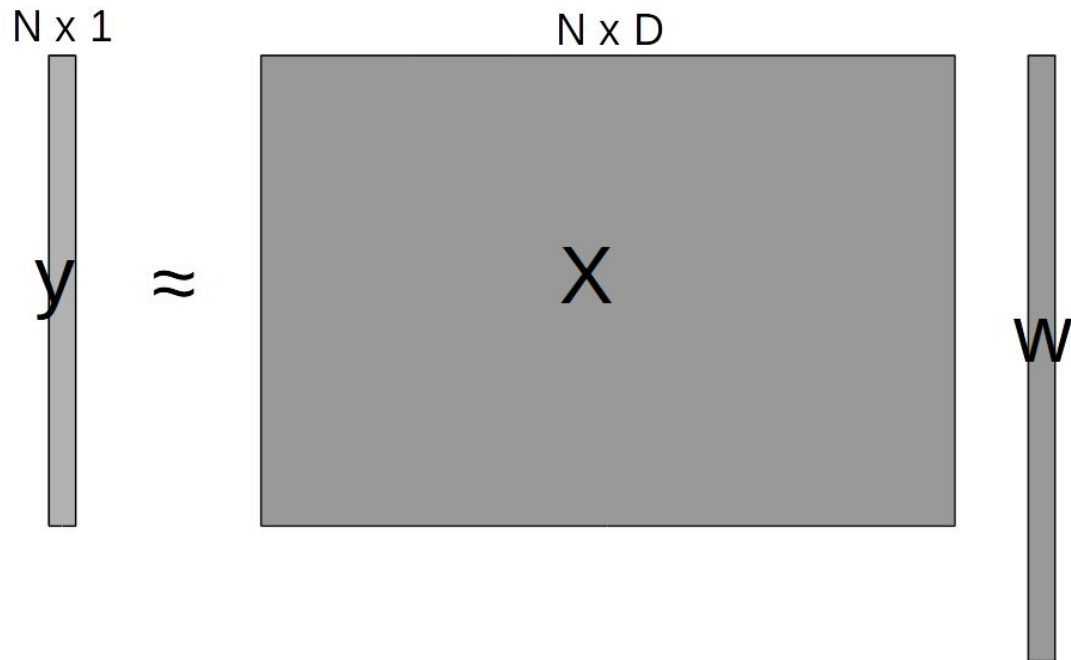
- A simple and interpretable linear model. Can also re-express it compactly for all the  $N$  examples

$$\mathbf{y} \approx \mathbf{X}\mathbf{w} \text{ (akin to a linear system of equations; } \mathbf{w} \text{ being the unknown)}$$

- Notation used here:
  - $\mathbf{w} \in \mathbb{R}^D$  and each  $\mathbf{x}_n \in \mathbb{R}^D$  are  $D \times 1$  column vectors
  - $\mathbf{X} = [\mathbf{x}_1 \ \mathbf{x}_2 \ \dots \ \mathbf{x}_N]^T$  is an  $N \times D$  matrix of features
  - $\mathbf{y} = [y_1 \ y_2 \ \dots \ y_N]^T$  is an  $N \times 1$  column vector of responses

# LINEAR REGRESSION

Linear system of equations with  $w$  being the unknown..





# LINEAR REGRESSION WITH SQUARED LOSS

- Our linear regression model:  $y_n \approx \mathbf{w}^\top \mathbf{x}_n$ . The goal is to learn  $\mathbf{w} \in \mathbb{R}^D$
- Let's use the **squared loss** to define our loss function.

$$\ell(y_n, \mathbf{w}^\top \mathbf{x}_n) = (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- Using the squared loss, the total (empirical) error on the training data

$$L_{\text{emp}}(\mathbf{w}) = \sum_{n=1}^N \ell(y_n, \mathbf{w}^\top \mathbf{x}_n) = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

- We'll estimate  $\mathbf{w}$  by minimizing  $L_{\text{emp}}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  (an optimization problem)

$$\hat{\mathbf{w}} = \arg \min_{\mathbf{w}} \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$$

# LEAST SQUARES LINEAR REGRESSION

- Recall our objective function:  $L_{emp} = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$

- Taking derivative of  $L_{emp}(\mathbf{w})$  w.r.t.  $\mathbf{w}$  and setting to zero

$$\sum_{n=1}^N 2(y_n - \mathbf{w}^\top \mathbf{x}_n) \frac{\partial}{\partial \mathbf{w}} (y_n - \mathbf{w}^\top \mathbf{x}_n) = 0 \quad \Rightarrow \quad \sum_{n=1}^N \mathbf{x}_n (y_n - \mathbf{x}_n^\top \mathbf{w}) = 0$$

- Simplifying further, we get a nice, closed form solution for  $\mathbf{w}$

$$\mathbf{w} = \left( \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^\top \right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

# LEAST SQUARES LINEAR REGRESSION

- Analytic, closed form solution, but has some issues
  - We didn't impose any regularization on  $w$  (thus prone to overfitting)
  - Have to invert a  $D \times D$  matrix; prohibitive especially when  $D$  (and  $N$ ) is large
  - The matrix  $\mathbf{X}^T \mathbf{X}$  may not even be invertible (e.g., when  $D > N$ ). Unique solution not guaranteed

# RIDGE REGRESSION: REGULARIZED LEAST SQUARES

- Least Squares objective:  $L_{emp} = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2$
- No constraints/regularization on  $\mathbf{w}$ . Components  $[w_1, w_2, \dots, w_D]$  of  $\mathbf{w}$  **may become arbitrarily large**. Why is this a bad thing to have?
- Let's add **squared  $l_2$  norm** of  $\mathbf{w}$  as a regularizer:  $R(\mathbf{f}) = R(\mathbf{w}) = \|\mathbf{w}\|^2$
- This results in the so-called “Ridge Regression” model.

$$L_{reg} = \sum_{n=1}^N (y_n - \mathbf{w}^\top \mathbf{x}_n)^2 + \lambda \|\mathbf{w}\|^2$$

- Note that  $\|\mathbf{w}\|^2 = \mathbf{w}^\top \mathbf{w} = \sum_{d=1}^D w_d^2$

# RIDGE REGRESSION: REGULARIZED LEAST SQUARES

- Minimizing  $L_{\text{reg}}$  will prevent components of  $\mathbf{w}$  from becoming very large.  
Why is this nice?
- Taking derivative of  $L_{\text{reg}}$  w.r.t.  $\mathbf{w}$  and setting to zero gives (verify yourself)

$$\mathbf{w} = \left( \sum_{n=1}^N \mathbf{x}_n \mathbf{x}_n^{\top} + \lambda \mathbf{I}_D \right)^{-1} \sum_{n=1}^N y_n \mathbf{x}_n = (\mathbf{X}^{\top} \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^{\top} \mathbf{y}$$

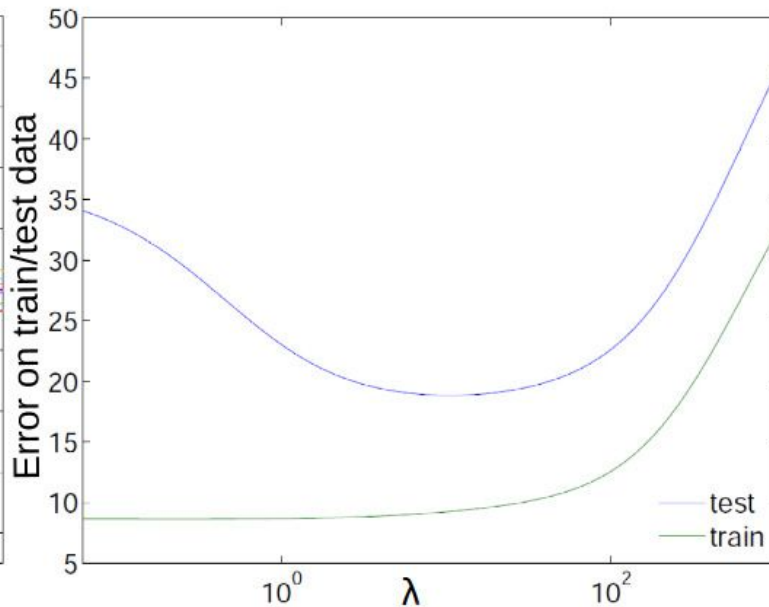
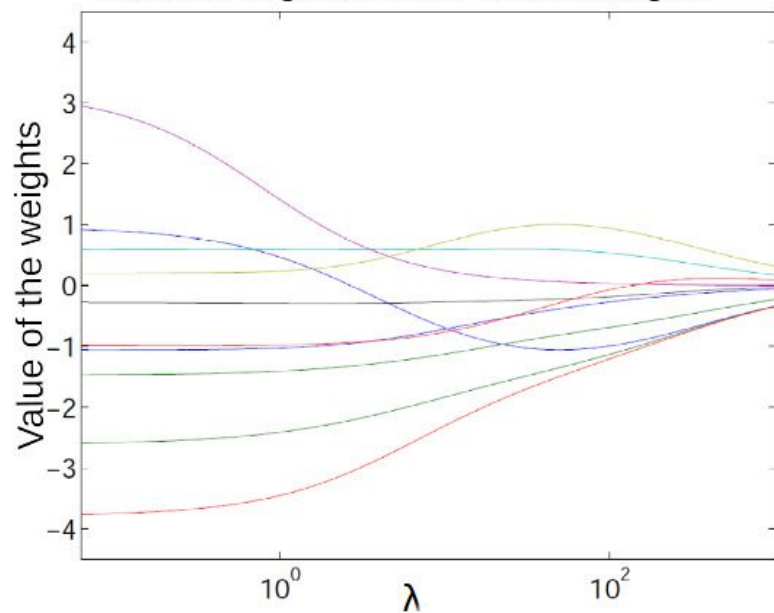
# INTUITIVELY, WHY SMALL WEIGHTS ARE GOOD?

- Small weights ensure that the function  $y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$  is **smooth** (i.e., we expect similar  $\mathbf{x}$ 's to have similar  $y$ 's). Below is an informal justification:
- Consider two points  $\mathbf{x}_n \in \mathbb{R}^D$  and  $\mathbf{x}_m \in \mathbb{R}^D$  that are exactly similar in all features **except the  $d$ -th feature** where they differ by a small value, say  $\epsilon$
- Assuming a simple/smooth function  $f(\mathbf{x})$ ,  $y_n$  and  $y_m$  should also be close
- However, as per the model  $y = f(\mathbf{x}) = \mathbf{w}^T \mathbf{x}$ ,  $y_n$  and  $y_m$  will differ by  $\epsilon w_d$
- Unless we constrain  $w_d$  to have a small value, the difference  $w_d$  would also be very large (which isn't what we want).
- That's why regularizing (via  $l_2$  regularization) and making the individual components of the weight vector small helps

# RIDGE REGRESSION: EFFECT OF REGULARIZATION

- Consider ridge regression on some data with 10 features (thus the weight vector  $w$  has 10 components)

Effect of regularization on the weights



# SOLUTION VIA GRADIENT-BASED METHODS

- Both least squares and ridge regression require matrix inversion.

$$\text{Least Squares } \mathbf{w} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y}$$

$$\text{Ridge } \mathbf{w} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I}_D)^{-1} \mathbf{X}^\top \mathbf{y}$$

- This can be **computationally very expensive** when  $D$  is very large.
- We can instead solve for  $\mathbf{w}$  more efficiently using generic/specialized optimization methods on the respective loss functions ( $L_{\text{emp}}$  or  $L_{\text{reg}}$ ).



# SOLUTION VIA GRADIENT-BASED METHODS

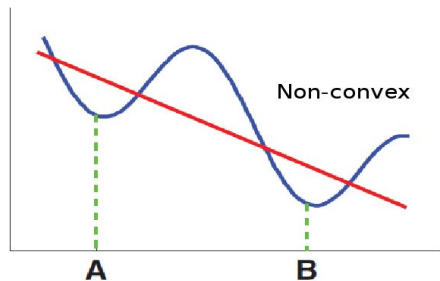
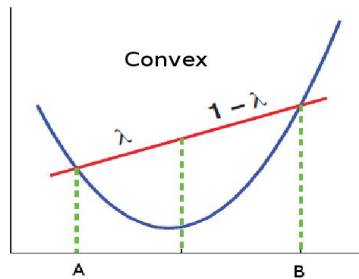
- A simple scheme can be the following iterative gradient-descent procedure
  - Start with an initial value of  $\mathbf{w} = \mathbf{w}^{(0)}$
  - Update  $\mathbf{w}$  by moving along the gradient of the loss function  $L$  ( $L_{\text{emp}}$  or  $L_{\text{reg}}$ ), where  $\eta$  is the learning rate
  - Repeat until converge
- For unreg. least squares, the gradient is

$$\mathbf{w}^{(t)} = \mathbf{w}^{(t-1)} - \eta \frac{\partial L}{\partial \mathbf{w}} \Big|_{\mathbf{w}=\mathbf{w}^{(t-1)}}$$

$$\frac{\partial L}{\partial \mathbf{w}} = - \sum_{n=1}^N \mathbf{x}_n (y_n - \mathbf{x}_n^{\top} \mathbf{w})$$

# GRADIENT-BASED METHODS: SOME NOTES

- Guaranteed to converge to a local minima
- Converge to global minima if the function is **convex**



- Note: The squared loss function in linear regression is convex
  - With  $l_2$  regularizer, it becomes strictly convex (single global minima).
- Learning rate is important (should not be too large or too small)
- Can also use stochastic/online gradient descent for more speed-ups. Require computing the gradients using only one or a small number of examples

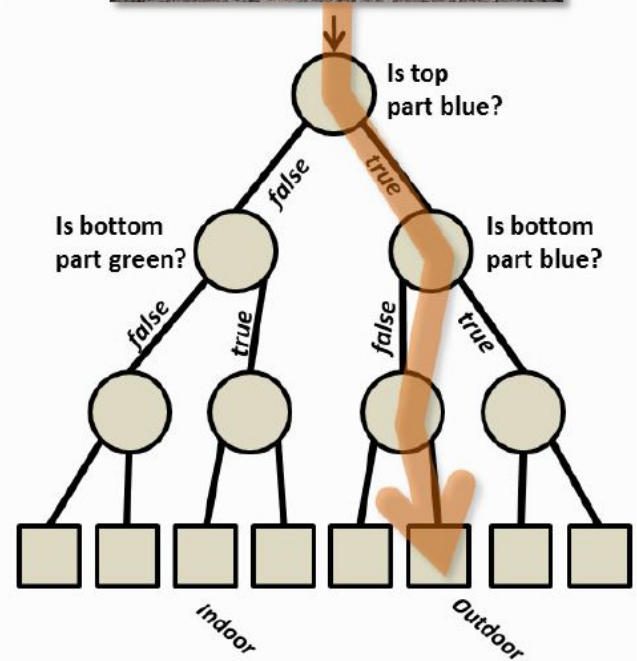
# DECISION TREES

# A CLASSIFICATION PROBLEM

**Indoor or Outdoor ?**



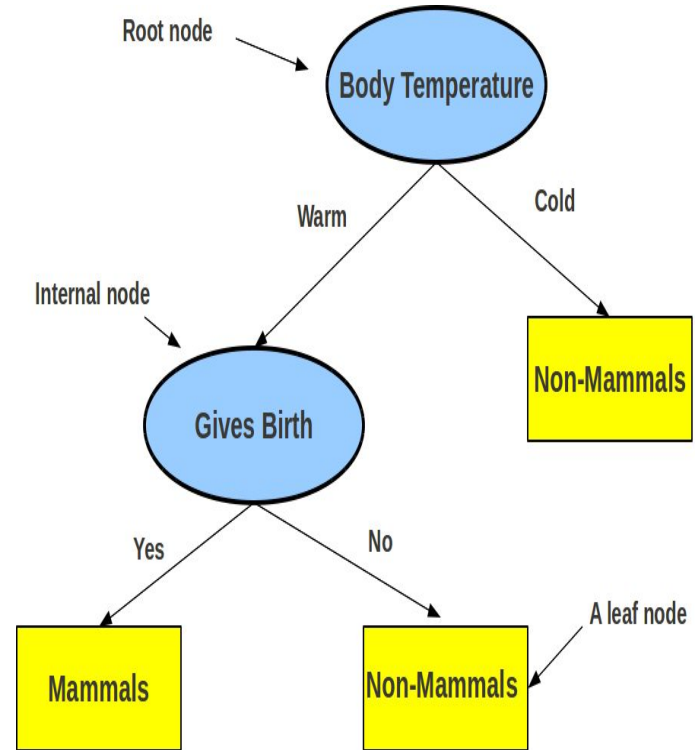
# PREDICTING BY ASKING QUESTIONS



How can we learn this tree using labeled training data?

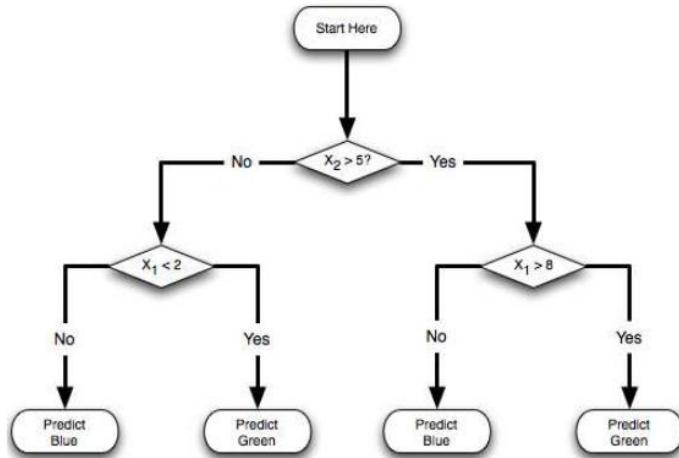
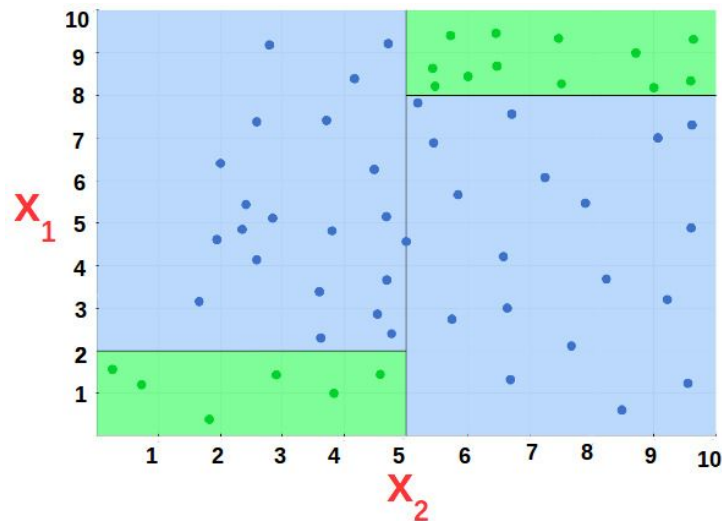
# DECISION TREE

- Defined by a **hierarchy** of rules (in form of a tree).
- Rules form the internal nodes of the tree (topmost internal node = root).
- Each internal node tests the value of some feature and “splits” data across the outgoing branches.
- Note: The tree need not be a binary tree
- (Labeled) Training data is used to construct the Decision Tree (DT)
- The DT can then be used to predict label  $y$  of a test example  $x$



# DECISION TREE: AN EXAMPLE

- Identifying the region - blue or green - a point lies in (binary classification).
  - Each point has 2 features: its coordinates  $\{x_1, x_2\}$  on the 2D plane
  - Left: Training data, Right: A DT constructed using this data

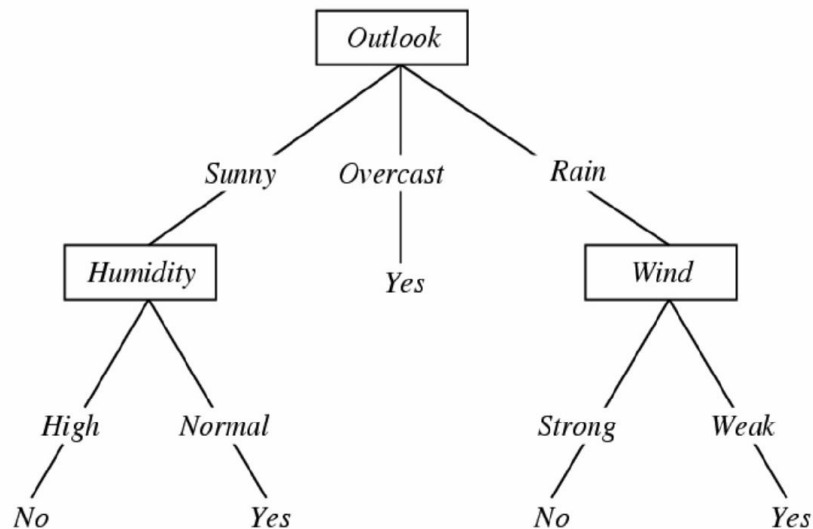


- The DT can be used to predict the region (blue/green) of a new test point
  - By testing the features of the test point
  - In the order defined by the DT (first  $x_2$  and then  $x_1$ )

# DECISION TREE: ANOTHER EXAMPLE

- Deciding whether to play or not to play Tennis on a Saturday
  - A binary classification problem (play vs no-play)
  - Each input (a Saturday) has 4 features: Outlook, Temp., Humidity, Wind
  - Left: Training data, Right: A decision tree constructed using this data

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



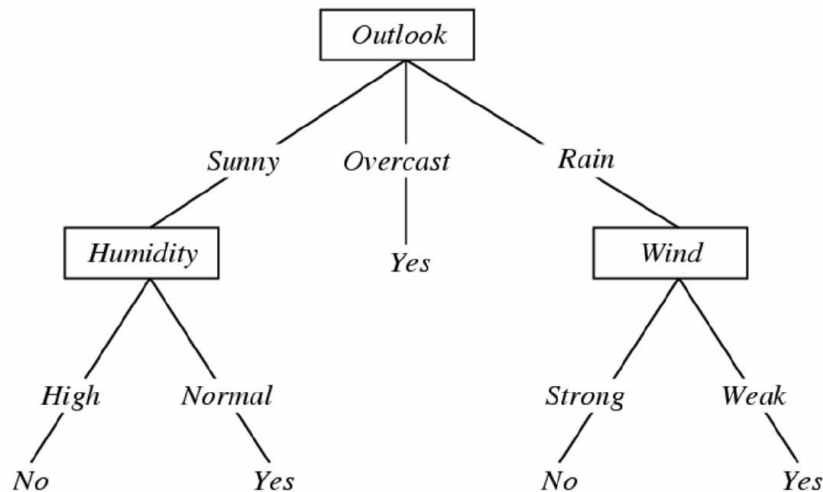
The DT can be used to predict play vs no-play for a new Saturday



# DECISION TREE CONSTRUCTION

- Now let's look at the playing Tennis example

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



- Question:** Why does it make more sense to test the feature “outlook” first?
- Answer:** Of all the 4 features, it's most informative
- We will see shortly how to quantify the informativeness

# ENTROPY

- Entropy is a measure of **randomness/uncertainty** of a set
- Assume our data is a set  $S$  of examples with  $C$  many classes
- $p_c$  is the probability that a random element of  $S$  belongs to class  $c$
- .. basically, the fraction of elements of  $S$  belonging to class  $c$
- Probability vector  $p = [p_1, p_2, \dots, p_C]$  is the **class distribution** of the set  $S$
- Entropy of the set  $S$

$$H(S) = - \sum_{c \in C} p_c \log_2 p_c$$

- If a set  $S$  of examples (or any subset of it) has..
  - Some dominant classes  $\implies$  small entropy of the class distribution
  - Equiprobable classes  $\implies$  high entropy of the class distribution
- We can assess **informativeness of each feature** by looking at **how much it reduces the entropy of the class distribution**

# INFORMATION GAIN

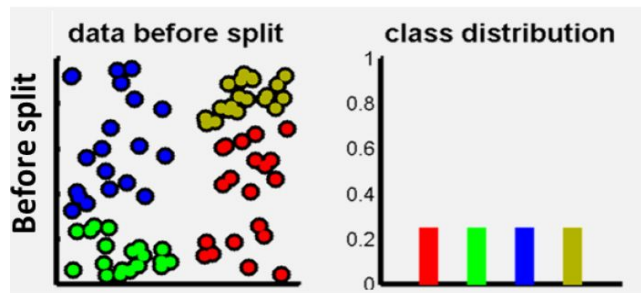
- Let's assume each element of S has a set of features
- Information Gain (IG) on knowing the value of some feature 'F '

$$IG(S, F) = H(S) - \sum_{f \in F} \frac{|S_f|}{|S|} H(S_f)$$

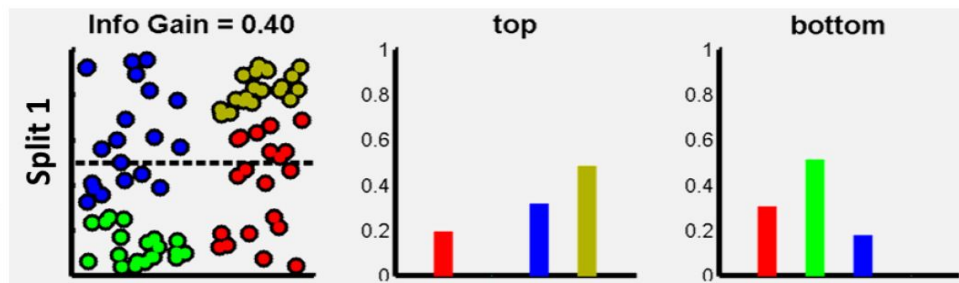
- $S_f$  denotes the subset of elements of S for which feature F has **value** f
- $IG(S, F)$  = entropy of S minus the weighted sum of entropy of its children
- $IG(S, F)$ : **Increase in our certainty** about S once we know the value of F

# ENTROPY AND INFORMATION GAIN: PICTORIALLY

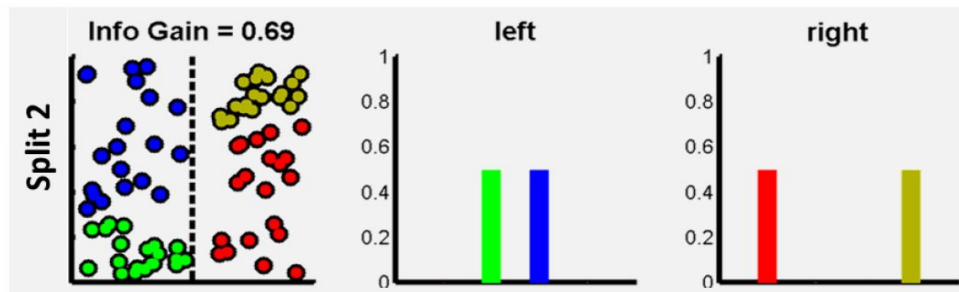
- Assume we have a 4-class problem. Each point has 2 features
- Which feature should we test (i.e., split on) first?



(a)



(b)



(c)

# COMPUTING INFORMATION GAIN

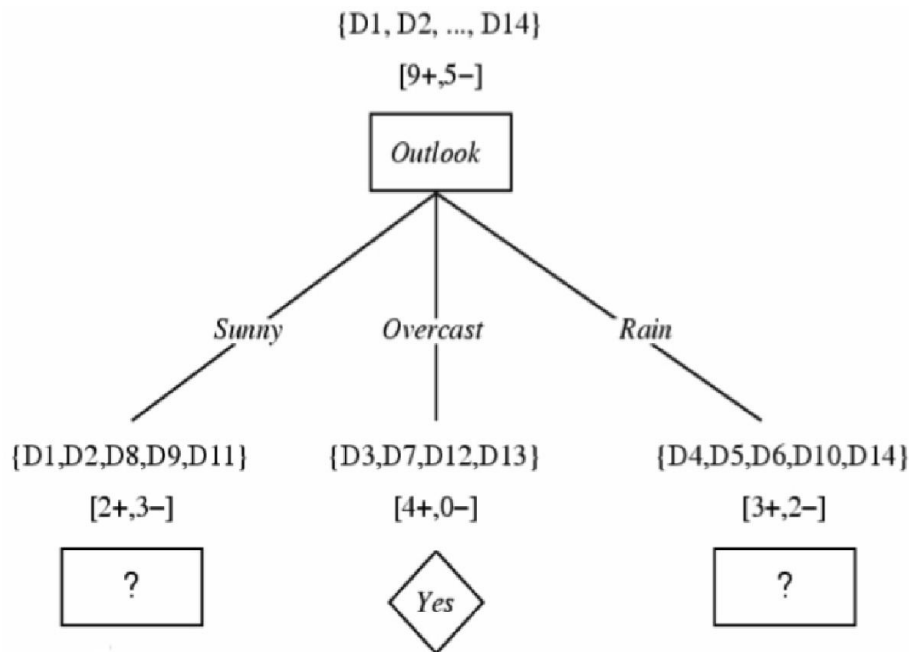
- Coming back to playing tennis..
- Let's begin with the root node of the DT and compute IG of each feature
- Consider feature “wind”  $\in \{\text{weak}, \text{strong}\}$  and its IG w.r.t. the root node
- Root node:  $S = [9+, 5-]$  (all training data: 9 play, 5 no-play)
- Entropy:  $H(S) = -(9/14) \log_2 (9/14) - (5/14) \log_2 (5/14) = 0.94$
- $S_{\text{weak}} = [6+, 2-] \Rightarrow H(S_{\text{weak}}) = 0.811$
- $S_{\text{strong}} = [3+, 3-] \Rightarrow H(S_{\text{strong}}) = 1$

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no

$$\begin{aligned} IG(S, \text{wind}) &= H(S) - \frac{|S_{\text{weak}}|}{|S|} H(S_{\text{weak}}) - \frac{|S_{\text{strong}}|}{|S|} H(S_{\text{strong}}) \\ &= 0.94 - 8/14 * 0.811 - 6/14 * 1 \\ &= 0.048 \end{aligned}$$

# CHOOSING THE MOST INFORMATIVE FEATURE

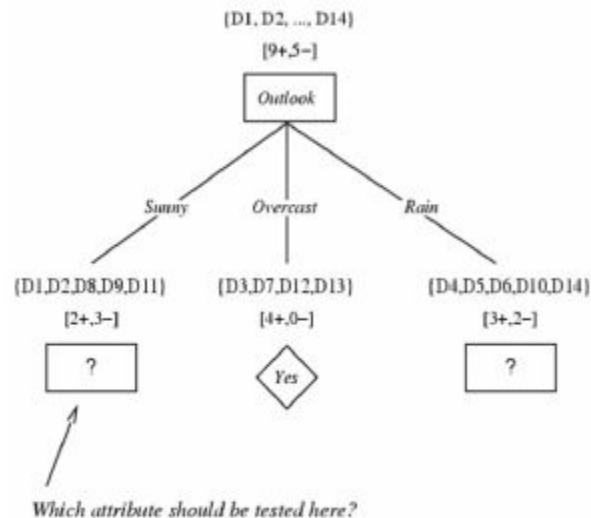
- At the root node, the information gains are:
  - $IG(S, \text{wind}) = 0.048$  (we already saw)
  - $IG(S, \text{outlook}) = 0.246$
  - $IG(S, \text{humidity}) = 0.151$
  - $IG(S, \text{temperature}) = 0.029$
- “outlook” has the maximum  $IG \implies$  chosen as the root node



# GROWING THE TREE

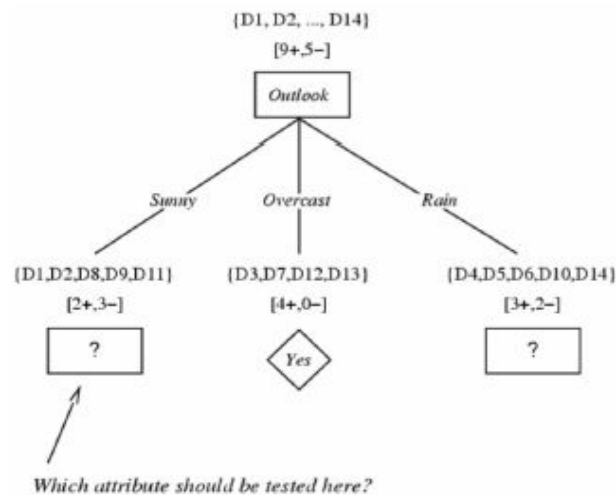
- How to decide which feature to test next ?
- **Rule:** Iterate - for each child node, select the feature with the highest IG

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



- For level-2, left node:  $S = [2+; 3-]$  (days 1,2,8,9,11)
- Compute the Information Gain for each feature (except outlook)
- The feature with the highest Information Gain should be chosen for this node

day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



For this node ( $S = [2+, 3-]$ ), the  $IG$  for the feature temperature:

$$IG(S, \text{temperature}) = H(S) - \sum_{v \in \{\text{hot}, \text{mild}, \text{cool}\}} \frac{|S_v|}{|S|} H(S_v)$$

$$S = [2+, 3-] \Rightarrow H(S) = -(2/5) * \log_2(2/5) - (3/5) * \log_2(3/5) = 0.971$$

$$S_{\text{hot}} = [0+, 2-] \Rightarrow H(S_{\text{hot}}) = -0 * \log_2(0) - (2/2) * \log_2(2/2) = 0$$

$$S_{\text{mild}} = [1+, 1-] \Rightarrow H(S_{\text{mild}}) = -(1/2) * \log_2(1/2) - (1/2) * \log_2(1/2) = 1$$

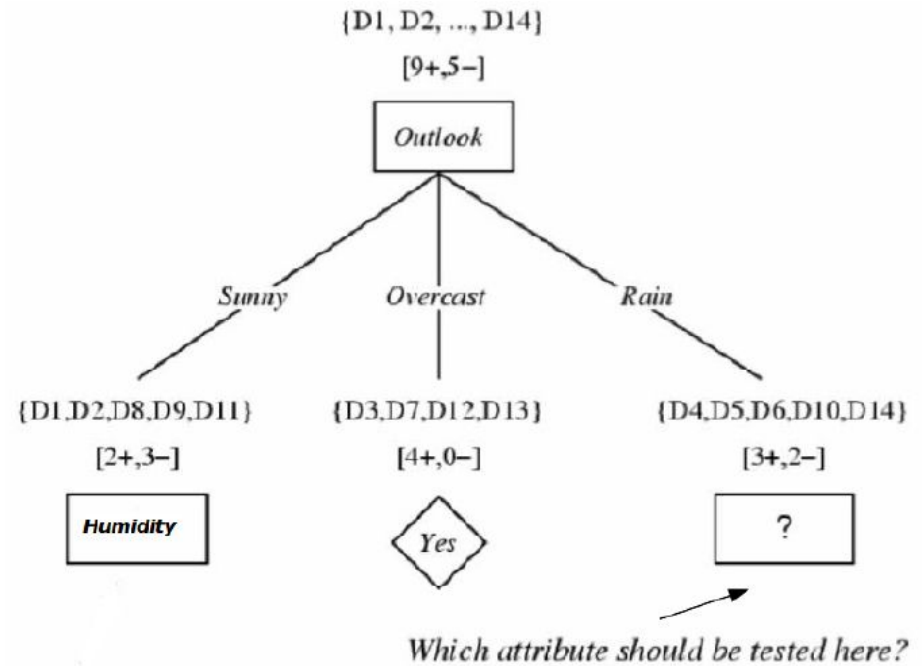
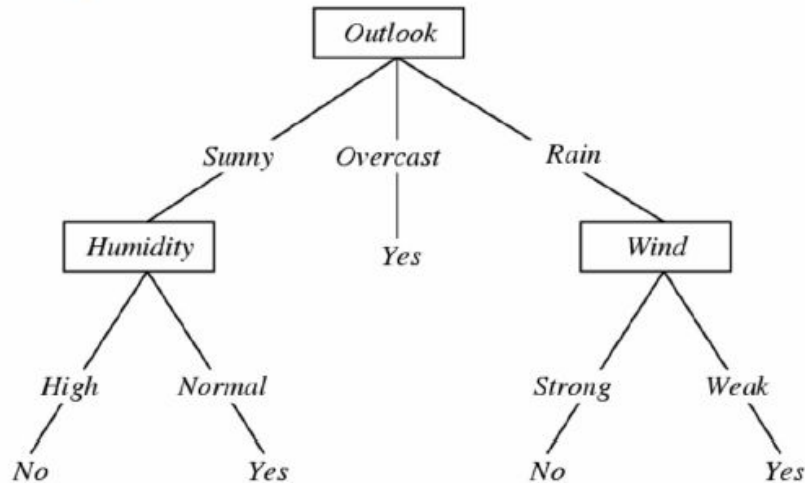
$$S_{\text{cool}} = [1+, 0-] \Rightarrow H(S_{\text{cool}}) = -(1/1) * \log_2(1/1) - (0/1) * \log_2(0/1) = 0$$

$$IG(S, \text{temperature}) = 0.971 - 2/5 * 0 - 2/5 * 1 - 1/5 * 0 = 0.570$$

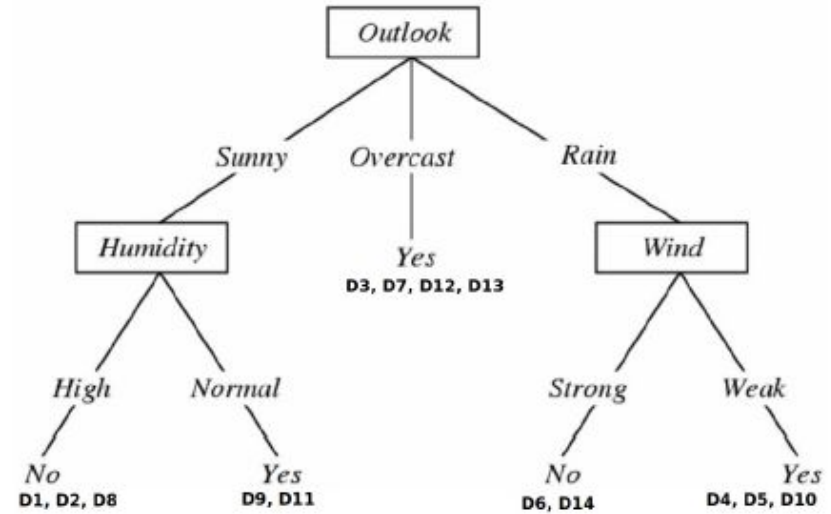
Likewise we can compute:  $IG(S, \text{humidity}) = 0.970$  ,  $IG(S, \text{wind}) = 0.019$



- Level-2, middle node: no need to grow (already a leaf)
- Level-2, right node: repeat the same exercise!
  - Compute IG for each feature (except outlook)
- Level-2 expansion gives us the following tree:



day	outlook	temperature	humidity	wind	play
1	sunny	hot	high	weak	no
2	sunny	hot	high	strong	no
3	overcast	hot	high	weak	yes
4	rain	mild	high	weak	yes
5	rain	cool	normal	weak	yes
6	rain	cool	normal	strong	no
7	overcast	cool	normal	strong	yes
8	sunny	mild	high	weak	no
9	sunny	cool	normal	weak	yes
10	rain	mild	normal	weak	yes
11	sunny	mild	normal	strong	yes
12	overcast	mild	high	strong	yes
13	overcast	hot	normal	weak	yes
14	rain	mild	high	strong	no



- Stop expanding a node further when:
  - It consist of examples all having the same label (the node becomes “**pure**”)
  - Or we run out of features to test!

# DECISION TREE LEARNING ALGORITHM

**A recursive algorithm:**

$DT(\text{Examples}, \text{Labels}, \text{Features})$ :

If all examples are positive, return a single node tree  $Root$  with label = +

If all examples are negative, return a single node tree  $Root$  with label = -

If all features exhausted, return a single node tree  $Root$  with majority label

Otherwise, let  $F$  be the feature having the highest information gain

$Root \leftarrow F$

For each possible value  $f$  of  $F$

- Add a tree branch below  $Root$  corresponding to the test  $F = f$
- Let  $Examples_f$  be the set of examples with feature  $F$  having value  $f$
- Let  $Labels_f$  be the corresponding labels
- If  $Examples_f$  is empty, add a leaf node below this branch with label = most common label in  $Examples$
- Otherwise, add the following subtree below this branch:

$DT(Examples_f, Labels_f, Features - \{F\})$

- Note:  $Features - \{F\}$  removes feature  $F$  from the feature set  $Features$

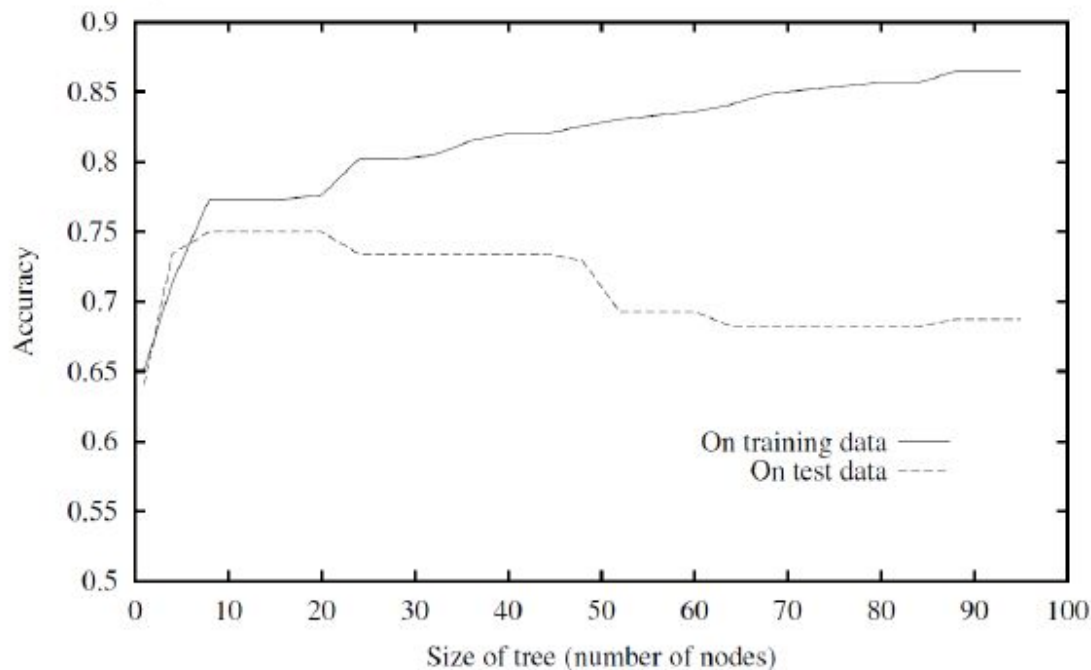
Activate Windows

Go to Settings to activate Wi

# OVERFITTING IN DECISION TREES

## Overtting Illustration

- High training accuracy doesn't necessarily imply high test accuracy



# AVOIDING OVERFITTING: DECISION TREE PRUNING

- Desired: a DT that is not too big in size, yet ts the training data reasonably
- Mainly two approaches
  - Prune while building the tree (**stopping early**)
  - Prune after building the tree (**post-pruning**)
- Criteria for judging which nodes could potentially be pruned
  - Use a **validation set** (separate from the training set)
    - Prune each possible node that doesn't hurt the accuracy on the validation set
    - **Greedily remove** the node that improves the validation accuracy the most
    - Stop when the validation set accuracy starts worsening
- **Minimum Description Length** (MDL): more details when we cover Model Selection

# DECISION TREE EXTENSIONS

- Real-valued features can be dealt with using thresholding
- Real-valued labels (**Regression Trees**) by re-defining entropy or using other criteria (how similar to each other are the  $y$ 's at any node)
- Other criteria for judging feature informativeness
  - Gini-index, misclassification rate
- More sophisticated decision rules at the internal nodes (anything that splits the data into homogeneous groups; e.g., a machine learning classifier)
- Handling features with differing costs

# SOME ASPECTS ABOUT DECISION TREES

## Some key strengths:

- Simple and easy to interpret
- Do not make any assumption about distribution of data
- Easily handle different types of features (real, categorical/nominal, etc.)
- Very fast at test time (just need to check the features, starting the root node and following the DT until you reach a leaf node)
- Multiple DTs can be combined via ensemble methods (e.g., Decision Forest)
- Each DT can be constructed using a (random) small subset of features.

# SOME ASPECTS ABOUT DECISION TREES

## Some key weaknesses:

- Learning the optimal DT is NP-Complete. The existing algorithms are heuristics (e.g., greedy selection of features)
- Can be unstable if some labeled examples are noisy
- Can sometimes become very complex unless some pruning is applied



# LOGISTIC REGRESSION

**Not really a regression....**

# LOGISTIC REGRESSION: THE MODEL

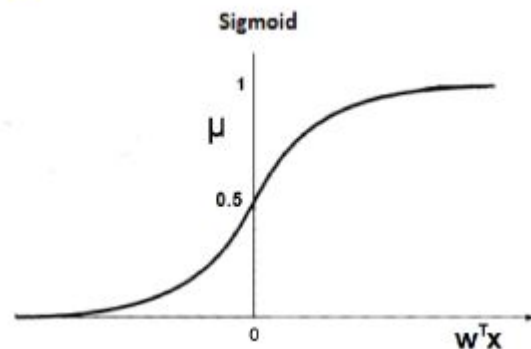
- A model for doing *probabilistic* binary classification
- Predicts **label probabilities** rather than a hard value of the label

$$\begin{aligned}p(y_n = 1 | \mathbf{x}_n, \mathbf{w}) &= \mu_n \\p(y_n = 0 | \mathbf{x}_n, \mathbf{w}) &= 1 - \mu_n\end{aligned}$$

- The model's prediction is a probability defined using the **sigmoid function**

$$f(\mathbf{x}_n) = \mu_n = \sigma(\mathbf{w}^\top \mathbf{x}_n) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x}_n)} = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$$

- The sigmoid first computes a real-valued “score”  $\mathbf{w}^\top \mathbf{x} = \sum w_d x_d$  and “squashes” it between (0,1) to turn this score into a **probability score**.
- Model parameter is the unknown  $\mathbf{w}$ . Need to learn it from training data.



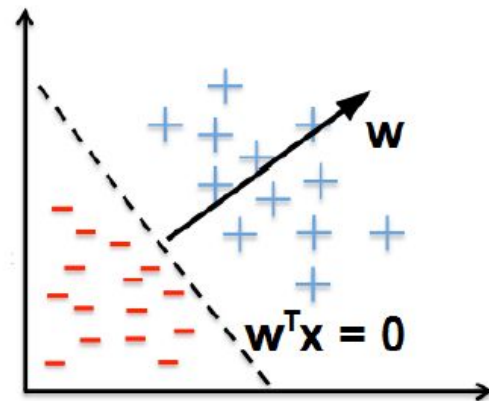
# LOGISTIC REGRESSION: AN INTERPRETATION

- Recall that the logistic regression model defines

$$p(y = 1|\mathbf{x}, \mathbf{w}) = \mu = \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(-\mathbf{w}^\top \mathbf{x})} = \frac{\exp(\mathbf{w}^\top \mathbf{x})}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

$$p(y = 0|\mathbf{x}, \mathbf{w}) = 1 - \mu = 1 - \sigma(\mathbf{w}^\top \mathbf{x}) = \frac{1}{1 + \exp(\mathbf{w}^\top \mathbf{x})}$$

- Thus if  $\mathbf{w}^\top \mathbf{x} > 0$  then the positive class is more probable.
- A linear classification model. Separates the two classes via a hyperplane (similar to other linear classification models such as Perceptron and SVM)



# LOGISTIC REGRESSION: THE LOSS FUNCTION

- What loss function to use? One option is to use the squared loss

$$\ell(y_n, f(\mathbf{x}_n)) = (y_n - f(\mathbf{x}_n))^2 = (y_n - \mu_n)^2 = (y_n - \sigma(\mathbf{w}^\top \mathbf{x}_n))^2$$

- This is **non-convex** and not easy to optimize.
- Consider the following loss function

$$\ell(y_n, f(\mathbf{x}_n)) = \begin{cases} -\log(\mu_n) & y_n = 1 \\ -\log(1 - \mu_n) & y_n = 0 \end{cases}$$

- This loss function makes intuitive sense
  - If  $y_n = 1$  but  $\mu_n$  is close to 0 (model makes error) then loss will be high
  - If  $y_n = 0$  but  $\mu_n$  is close to 1 (model makes error) then loss will be high

- The above loss function can be combined and written more compactly as

$$\ell(y_n, f(\mathbf{x}_n)) = -y_n \log(\mu_n) - (1 - y_n) \log(1 - \mu_n)$$

- This is a function of the unknown parameter  $w$  since  $\mu_n = (\mathbf{w}^T \mathbf{x}_n)$

The loss function over the entire training data

$$L(\mathbf{w}) = \sum_{n=1}^N \ell(y_n, f(\mathbf{x}_n)) = \sum_{n=1}^N [-y_n \log(\mu_n) - (1 - y_n) \log(1 - \mu_n)]$$

This is also known as the **cross-entropy loss**

- Sum of the cross-entropies b/w true label  $y_n$  and predicted label prob.  $\mu_n$

Plugging in  $\mu_n = \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{1 + \exp(\mathbf{w}^\top \mathbf{x}_n)}$  and chugging, we get (verify yourself)

$$L(\mathbf{w}) = - \sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

We can add a regularizer (e.g., squared  $\ell_2$  norm of  $\mathbf{w}$ ) to prevent overfitting

$$L(\mathbf{w}) = - \sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) + \lambda \|\mathbf{w}\|^2$$

# ESTIMATING THE WEIGHT VECTOR $\mathbf{w}$

Loss function/NLL for logistic regression (ignoring the regularizer term)

$$L(\mathbf{w}) = - \sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n)))$$

The loss function is convex in  $\mathbf{w}$  (thus has a unique minimum)

The gradient/derivative of  $L(\mathbf{w})$  w.r.t.  $\mathbf{w}$  (let's ignore the regularizer)

$$\begin{aligned} \mathbf{g} = \frac{\partial L(\mathbf{w})}{\partial \mathbf{w}} &= \frac{\partial}{\partial \mathbf{w}} \left[ - \sum_{n=1}^N (y_n \mathbf{w}^\top \mathbf{x}_n - \log(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))) \right] \\ &= - \sum_{n=1}^N \left( y_n \mathbf{x}_n - \frac{\exp(\mathbf{w}^\top \mathbf{x}_n)}{(1 + \exp(\mathbf{w}^\top \mathbf{x}_n))} \mathbf{x}_n \right) \\ &= - \sum_{n=1}^N (y_n - \mu_n) \mathbf{x}_n = \mathbf{X}^\top (\boldsymbol{\mu} - \mathbf{y}) \end{aligned}$$

Can't get a closed form solution for  $\mathbf{w}$  by setting the derivative to zero

- Need to use iterative methods (e.g., gradient descent) to solve for  $\mathbf{w}$

# GRADIENT DESCENT FOR LOGISTIC REGRESSION

We can use gradient descent (GD) to solve for  $\mathbf{w}$  as follows:

- Initialize  $\mathbf{w}^{(1)} \in \mathbb{R}^D$  randomly.
- Iterate the following until convergence

$$\underbrace{\mathbf{w}^{(t+1)}}_{\text{new value}} = \underbrace{\mathbf{w}^{(t)}}_{\text{previous value}} - \eta \underbrace{\sum_{n=1}^N (\mu_n^{(t)} - y_n) \mathbf{x}_n}_{\text{gradient at previous value}}$$

where  $\eta$  is the **learning rate** and  $\mu^{(t)} = \sigma(\mathbf{w}^{(t)\top} \mathbf{x}_n)$  is the predicted label probability for  $\mathbf{x}_n$  using  $\mathbf{w} = \mathbf{w}^{(t)}$  from the previous iteration

Note that the updates give larger weights to those examples on which the current model makes larger mistakes, as measured by  $(\mu_n^{(t)} - y_n)$



# MORE ON GRADIENT DESCENT...

GD can converge slowly and is also sensitive to the step size

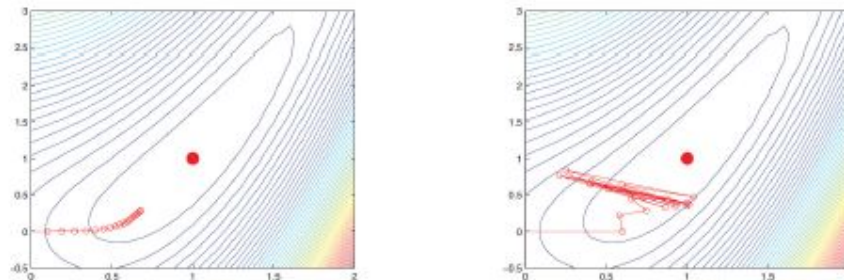


Figure: Left: small step sizes. Right: large step sizes

Several ways to remedy this<sup>1</sup>. E.g.,

- Choose the optimal step size  $\eta_t$  (different in each iteration) by **line-search**
- Add a **momentum term** to the updates

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta_t \mathbf{g}^{(t)} + \alpha_t (\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)})$$

- Use **second-order methods** (e.g., **Newton's method**) to exploit the **curvature** of the loss function  $L(\mathbf{w})$ : Requires computing the **Hessian matrix**

# THAT'S ALL FOLKS!

Checkout the Python Notebook again now.