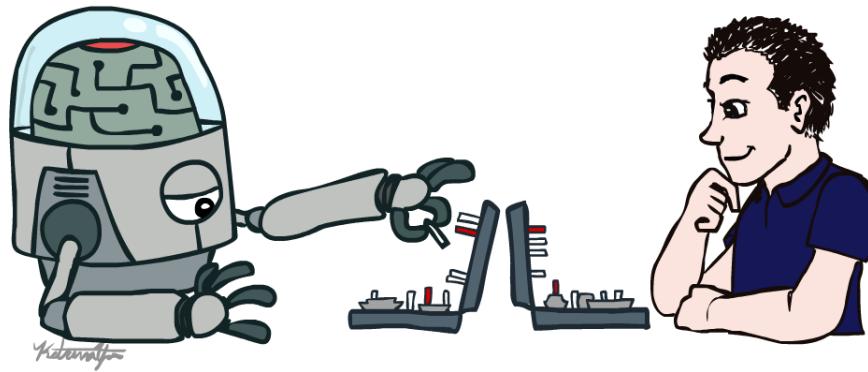


# CSE 3521: Introduction to Artificial Intelligence



[Many slides are adapted from the [UC Berkeley, CS188 Intro to AI](#) at UC Berkeley and previous CSE 3521 course at OSU.]



# Decision Tree

# ID3 Algorithm

---

## **ID3 (S, A, V)**

Let:

S = Learning Set  
A = Attribute Set  
V = Attribute Values

Begin

    Load learning sets and create decision tree root node(rootNode),  
    Add learning set S into root not as its subset

    For rootNode, compute Entropy(rootNode.subset)

    If Entropy(rootNode.subset) == 0 (subset is homogeneous)  
        return a leaf node

    If Entropy(rootNode.subset)!= 0 (subset is not homogeneous)  
        compute Information Gain for each attribute left (not been used for splitting)  
        Find attribute A with Maximum(Gain(S,A))  
        Create child nodes for this root node and add to rootNode in the decision tree

    For each child of the rootNode  
        Apply ID3(S,A,V)  
    Continue until a node with Entropy of 0 or a leaf node is reached

End

# Decision Tree Equations

---

$$Entropy(S) = \sum_{i=1}^c -p_i \log_2 p_i$$

S= Collection of Examples

$p_i$  = proportion of S that belongs to class i

$$Gain(S, A) = Entropy(S) - \sum_{v \in values(A)} \frac{|S_v|}{|S|} Entropy(S_v)$$

A= Attribute

values(A) = set of all possible values for A

$|S|$  = number of example in S

$|S_v|$  = number of example in S where the value of A is v

# Naïve Bayes

# Bayes Rules

---

Bayes Rule tells us how to flip the conditional

Reason about effects to causes

Useful if you assume a generative model for your data

$$P(H|D) = \frac{P(D|H)P(H)}{\sum_h P(D|H)P(H)}$$

The diagram illustrates the components of Bayes' Rule. At the top, the words "Likelihood" and "Prior" are positioned above the equation. Arrows point from these words to the terms  $P(D|H)$  and  $P(H)$  respectively. At the bottom left, the word "Posterior" is placed next to the term  $P(H|D)$ , with an arrow pointing from it. At the bottom right, the word "Normalizer" is placed next to the denominator  $\sum_h P(D|H)P(H)$ , with an arrow pointing from it.

# Bayes Rules

---

Bayes Rule tells us how to flip the conditional  
Reason about effects to causes  
Useful if you assume a generative model for your data

$$P(H|D) \propto P(D|H)P(H)$$

Likelihood                          Prior

Posterior                          Proportional To  
(Doesn't sum to 1)

The diagram illustrates the components of Bayes' Rule. At the top, 'Likelihood' and 'Prior' are shown with arrows pointing downwards towards the central equation  $P(H|D) \propto P(D|H)P(H)$ . From the bottom left, an arrow points upwards to the word 'Posterior'. From the bottom right, an arrow points upwards to the text 'Proportional To (Doesn't sum to 1)'.

# Learning: Naïve Bayes Classifier

---

- First attempt: maximum likelihood estimates
  - simply use the frequencies in the data

$$\hat{P}(c_j) = \frac{\text{doccount}(C = c_j)}{N_{\text{doc}}}$$

$$\hat{P}(w_i | c_j) = \frac{\text{count}(w_i, c_j)}{\sum_{w \in V} \text{count}(w, c_j)}$$

# Smoothing: Naïve Bayes Parameter Estimation

---

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c)}{\sum_{w \in V} \text{count}(w_i, c)}$$

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c) + \alpha}{\sum_{w \in V} \text{count}(w_i, c) + \alpha|V|}$$

## Laplace Smoothing: Naïve Bayes Parameter Estimation

---

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c)}{\sum_{w \in V} \text{count}(w_i, c)}$$

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c) + \alpha}{\sum_{w \in V} \text{count}(w_i, c) + \alpha |V|}$$

$$\hat{P}(w_i|c) = \frac{\text{count}(w_i, c) + 1}{\sum_{w \in V} \text{count}(w_i, c) + |V|}$$

# Naïve Bayes Parameter Learning: Step

---

- Calculate  $P(c_j)$  terms

- For each  $c_j$  in  $C$  do

$docs_j \leftarrow$  all docs with class =  $c_j$

$$P(c_j) \leftarrow \frac{|docs_j|}{|\text{total \# documents}|}$$

- From training corpus, extract Vocabulary

- Calculate  $P(w_k | c_j)$  terms

- $Text_j \leftarrow$  single doc containing all sentences from class =  $c_j$
  - For each word  $w_k$  in *Vocabulary*

$n_k \leftarrow$  # of occurrences of  $w_k$  in  $Text_j$

$n \leftarrow$  # of words in class  $Text_j$

$$P(w_k | c_j) \leftarrow \frac{n_k + \alpha}{n + \alpha |\text{Vocabulary}|}$$

# Naïve Bayes Evaluation: Step

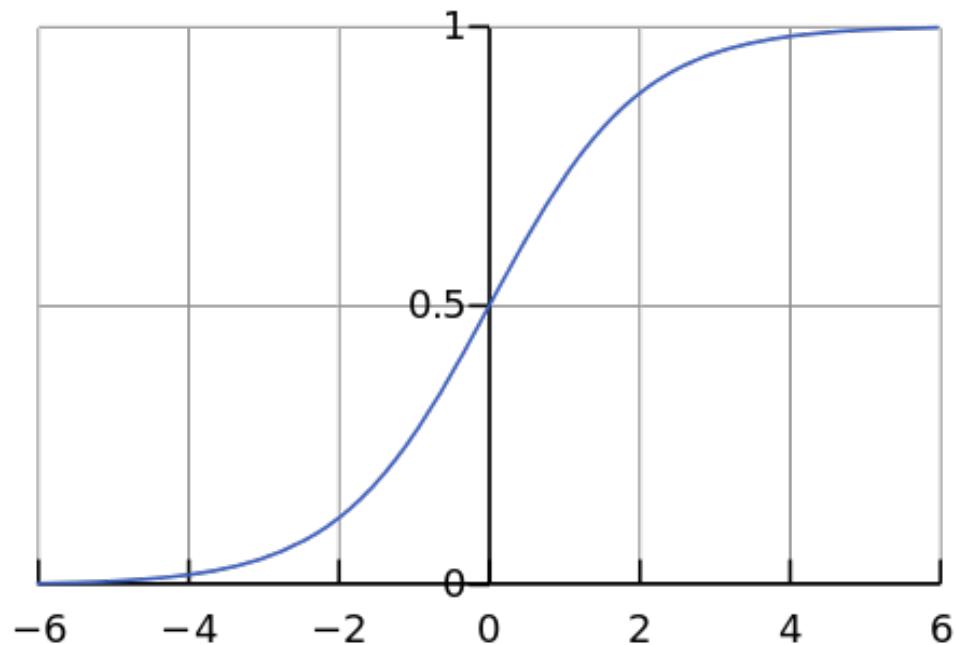
---

- # of correct prediction = 0
- For each  $t_i$  in Test data:
  - max\_prob = 0
  - For each class  $c_j$ 
    - Find  $p_j$  = probability of  $t_j$  to be in  $c_j$
    - If  $p_j > \text{max\_prob}$ :
      - $\text{max\_prob} = p_j$ ,  $\text{max\_prob\_class} = c_j$
  - If  $\text{max\_prob\_class} ==$  actual label of  $t_j$ :
    - # of correct prediction += 1
- Accuracy = 
$$\frac{\text{# of correct prediction}}{\text{# of test data}}$$

# Logistic Regression

# The Logistic function

---



$$P(\text{spam}|x) = \frac{e^z}{e^z + 1} = \frac{1}{1 + e^{-z}}$$

# The Dot Product

---

$$z = \sum_{i=0}^{|X|} w_i x_i$$

- Intuition: weighted sum of features
- All Linear models have this form

## NB vs. LR

---

- Both compute the dot product
- NB: sum of log probabilities
- LR: logistic function

## NB vs. LR:Parameter Learning

---

- Naïve Bayes:
  - Learn conditional probabilities **independently** by counting
- Logistic Regression:
  - Learn weights **jointly**

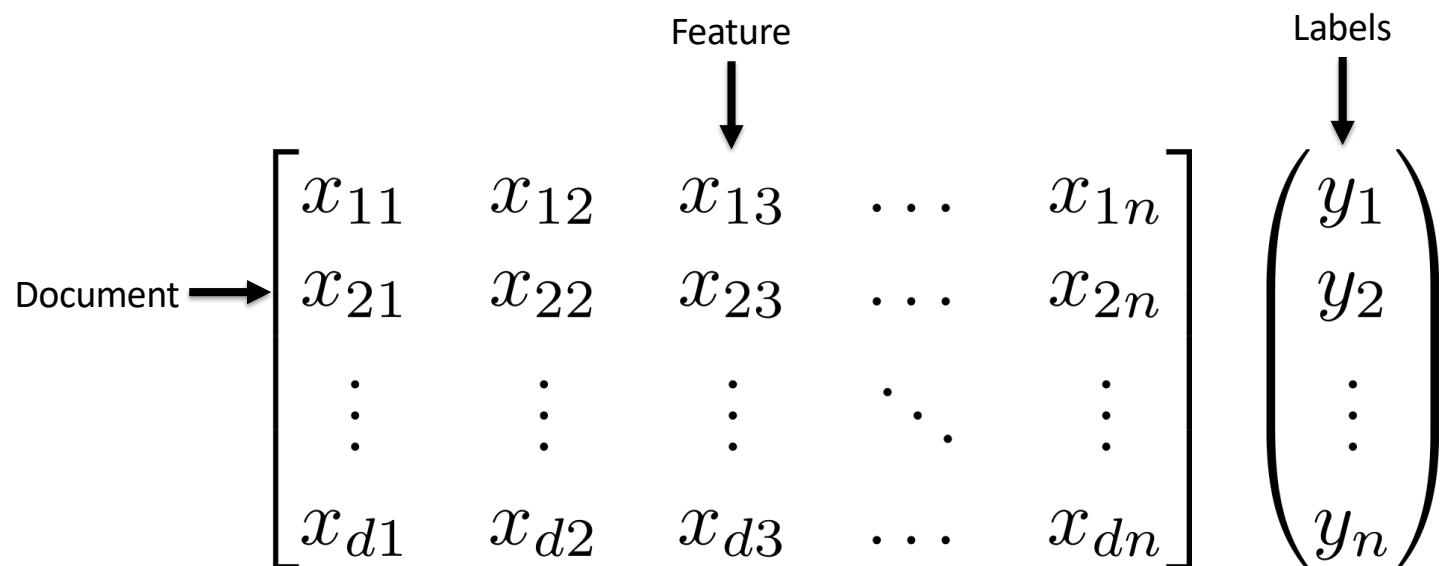
# LR: Learning Weights

---

- Given: a set of feature vectors and labels
- Goal: learn the weights

# Learning Weights

---



## Q: what parameters should we choose?

---

- What is the right value for the weights?
- Maximum Likelihood Principle:
  - Pick the parameters that maximize the probability of the data

# Maximum Likelihood Estimation

---

$$w_{\text{MLE}} = \operatorname{argmax}_w \log P(y_1, \dots, y_d | x_1, \dots, x_d; w)$$

$$= \operatorname{argmax}_w \sum_i \log P(y_i | x_i; w)$$

$$= \operatorname{argmax}_w \sum_i \log \begin{cases} p_i, & \text{if } y_i = 1 \\ 1 - p_i, & \text{if } y_i = 0 \end{cases}$$

$$= \operatorname{argmax}_w \sum_i \log p_i^{\mathbb{I}(y_i=1)} (1 - p_i)^{\mathbb{I}(y_i=0)}$$

# Maximum Likelihood Estimation

---

$$= \operatorname{argmax}_w \sum_i \log p_i^{\mathbb{I}(y_i=1)} (1 - p_i)^{\mathbb{I}(y_i=0)}$$

$$= \operatorname{argmax}_w \sum_i y_i \log p_i + (1 - y_i) \log(1 - p_i)$$

# Maximum Likelihood Estimation

---

- Unfortunately there is no closed form solution
  - (like there was with naïve bayes)
- Solution:
  - Iteratively climb the log-likelihood surface through the derivatives for each weight

# Derivative of Sigmoid

---

$$\begin{aligned}\frac{ds(x)}{dx} &= \frac{d}{dx} \frac{1}{1 + e^{-x}} \\ &= -\left(\frac{1}{1 + e^{-x}}\right)^2 \frac{d}{dx}(1 + e^{-x}) \\ &= -\left(\frac{1}{1 + e^{-x}}\right)^2 e^{-x}(-1) \\ &= s(x)(1 - s(x))\end{aligned}$$

# Gradient ascent

---

Loop While not converged:

For all features  $j$ , compute and add derivatives

$$w_j^{\text{new}} = w_j^{\text{old}} + \eta \frac{\partial}{\partial w_j} \mathcal{L}(w)$$

$\mathcal{L}(w)$ : Training set log-likelihood

$\left( \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots, \frac{\partial \mathcal{L}}{\partial w_n} \right)$  : Gradient vector

# Gradient ascent

---

Loop While not converged:

For all features  $j$ , compute and add derivatives

$$w_j^{\text{new}} = w_j^{\text{old}} + \eta \frac{\partial}{\partial w_j} \mathcal{L}(w)$$

$\mathcal{L}(w)$ : Training set log-likelihood       $\frac{\partial \mathcal{L}}{\partial w_j} = \sum_i (y_i - p_i) x_j$

$\left( \frac{\partial \mathcal{L}}{\partial w_1}, \frac{\partial \mathcal{L}}{\partial w_2}, \dots, \frac{\partial \mathcal{L}}{\partial w_n} \right)$  : Gradient vector

# LR Gradient

---

$$\frac{\partial \mathcal{L}}{\partial w_j} = \sum_i (y_i - p_i) x_j$$

$j \rightarrow$  iterating over features

$i \rightarrow$  iterating over training examples

$Y_i$  = True label

$P_i$  = Predicted Label

# MultiClass Classification

---

- Q: what if we have more than 2 categories?
  - Sentiment: Positive, Negative, Neutral
  - Document topics: Sports, Politics, Business, Entertainment, ...

# MultiClass Classification

---

- Q: what if we have more than 2 categories?
  - Sentiment: Positive, Negative, Neutral
  - Document topics: Sports, Politics, Business, Entertainment, ...
- Could train a separate logistic regression model for each category...
- Pretty clear what to do with Naive Bayes.

# Log-Linear Models

---

$$P(y|x) \propto e^{w \cdot f(d,y)}$$

$$P(y|x) = \frac{1}{Z(w)} e^{w \cdot f(d,y)}$$

# MultiClass Logistic Regression

---

$$P(y|x) \propto e^{w \cdot f(d,y)}$$

$$P(y|x) = \frac{1}{Z(w)} e^{w \cdot f(d,y)}$$

$$P(y|x) = \frac{e^{w \cdot f(d,y)}}{\sum_{y' \in Y} e^{w \cdot f(d,y')}}$$

# MultiClass Logistic Regression

---

- Binary logistic regression:
  - We have one feature vector that matches the size of the vocabulary
- Multiclass in practice:
  - one weight vector for each category

$w_{\text{pos}}$

$w_{\text{neg}}$

$w_{\text{neut}}$

# MultiClass Logistic Regression

---

- Binary logistic regression:
    - We have one feature vector that matches the size of the vocabulary
  - Multiclass in practice:
    - one weight vector
- Can represent this in practice with  
one giant weight vector and  
repeated features for each category.

$w_{\text{pos}}$

$w_{\text{neg}}$

$w_{\text{neut}}$

# Maximum Likelihood Estimation

---

$$w_{\text{MLE}} = \operatorname{argmax}_w \log P(y_1, \dots, y_n | x_1, \dots, x_n; w)$$

$$= \operatorname{argmax}_w \sum_i \log P(y_i | x_i; w)$$

$$= \operatorname{argmax}_w \sum_i \log \frac{e^{w \cdot f(x_i, y_i)}}{\sum_{y' \in Y} e^{w \cdot f(x_i, y_i)}}$$

# Multiclass Learning

---

$$\text{LR : } \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D \sum_{y \in Y} f_j(y, d_i) P(y|d_i)$$

# Perceptron

# Perceptron Algorithm

---

- Algorithm is Very similar to logistic regression
- Not exactly computing gradients

# Perceptron Algorithm

---

- Algorithm is Very similar to logistic regression
- Not exactly computing gradients

```
Initialize weight vector w = 0
Loop for K iterations
    Loop For all training examples xi
        yi' = sign(w * xi)
        if yi' != yi
            w += (yi - yi') * xi
```

# Perceptron Notes

---

- Guaranteed to converge if the data is linearly separable
- Only hyperparameter is maximum number of iterations
- Parameter averaging

## Differences: LR vs. Perceptron

---

- Batch Learning vs. Online learning
- Perceptron doesn't always make updates

# Online Learning (perceptron)

---

- Rather than making a full pass through the data, compute gradient and update parameters after each training example.
- Gradients will be less accurate, but the overall effect is to move in the right direction
- Often works well and converges faster than batch learning

# Multiclass Learning

---

$$\text{Perceptron} : \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D f_j(\arg \max_{y \in Y} P(y|d_i), d_i)$$

# Multiclass Learning

---

$$\text{LR} : \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D \sum_{y \in Y} f_j(y, d_i) P(y|d_i)$$

$$\text{Perceptron} : \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D f_j(\arg \max_{y \in Y} P(y|d_i), d_i)$$

# Multiclass Learning

---

$$\text{LR : } \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D \sum_{y \in Y} f_j(y, d_i) P(y|d_i)$$

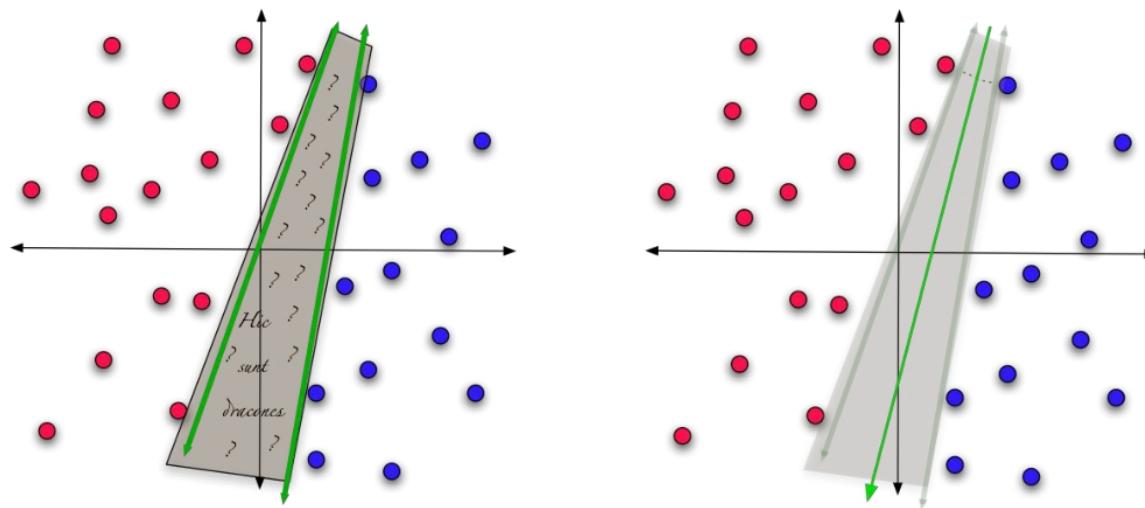
$$\text{Perceptron : } \frac{\partial \mathcal{L}}{\partial w_j} = \sum_{i=1}^D f_j(y_i, d_i) - \sum_{i=1}^D f_j(\arg \max_{y \in Y} P(y|d_i), d_i)$$

# SVM

# Margin

---

- the distance to closest point in the training data
- We tend to get better generalization to **unseen data** if we choose the separating hyperplane which *maximizes the margin*

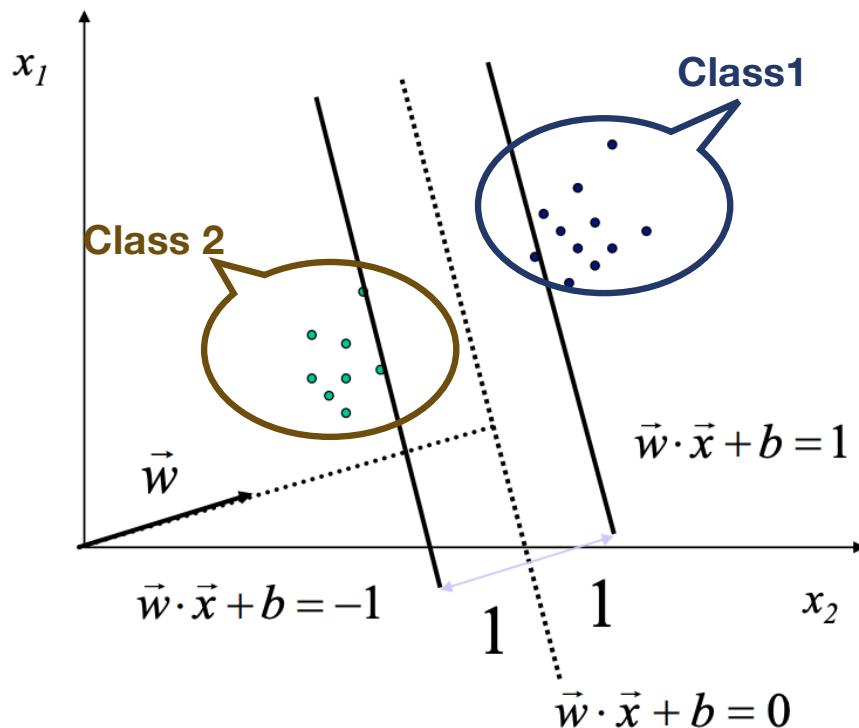


# Support Vector Machines

---

- A learning method which *explicitly calculates the maximum margin hyperplane* by solving a gigantic quadratic programming minimization problem.
- Among the very highest -performing traditional machine learning techniques.
- But it's relatively slow and quite complicated.

# Setting Up the Optimization Problem



The maximum margin can be characterized as a solution to an optimization problem:

$$\begin{aligned} & \max. \frac{2}{\|\vec{w}\|} \\ & s.t. (\vec{w} \cdot \vec{x} + b) \geq 1, \forall x \text{ of class 1} \\ & (\vec{w} \cdot \vec{x} + b) \leq -1, \forall x \text{ of class 2} \end{aligned}$$

Define the margin (whatever it turns out to be) to be *one unit of width*.

# Setting Up the Optimization Problem

---

- If class 1 corresponds to 1 and class 2 corresponds to -1, we can rewrite

$$(w \cdot x_i + b) \geq 1, \quad \forall x_i \text{ with } y_i = 1$$

$$(w \cdot x_i + b) \leq -1, \quad \forall x_i \text{ with } y_i = -1$$

- as

$$y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

- So the problem becomes:

$$\max. \frac{2}{\|w\|}$$

or

$$\min. \frac{1}{2} \|w\|^2$$

$$s.t. y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

$$s.t. y_i(w \cdot x_i + b) \geq 1, \quad \forall x_i$$

# Linear, (Hard-Margin) SVM Formulation

---

- Find  $w, b$  that solves

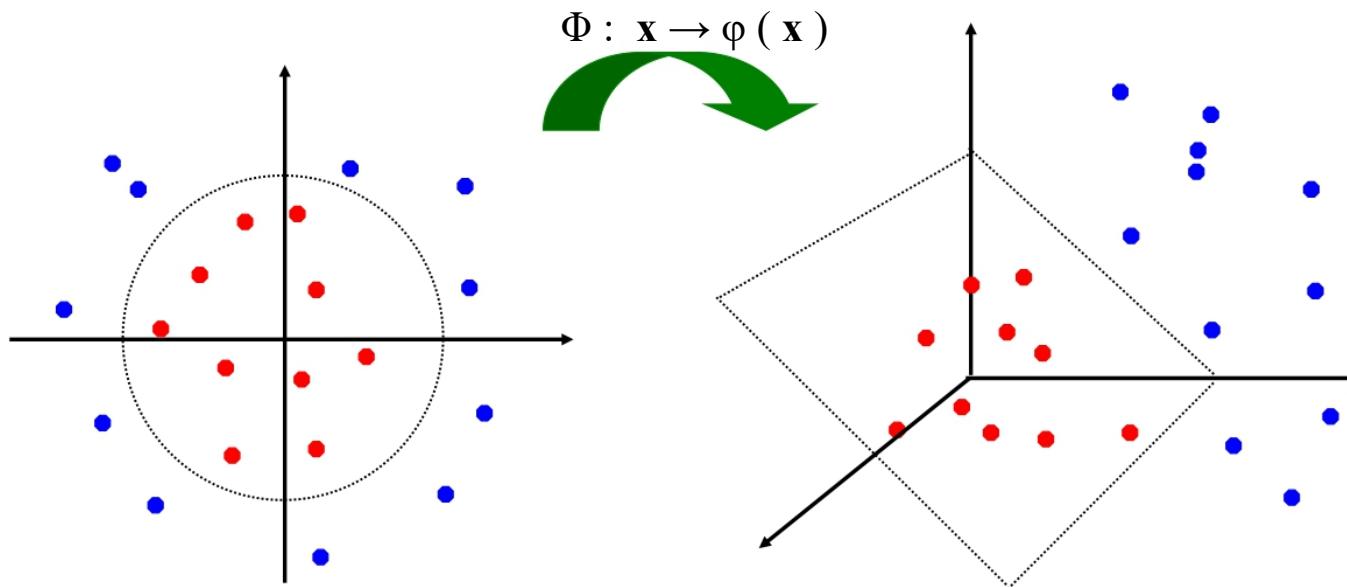
$$\min. \frac{1}{2} \|w\|^2$$

$$s.t. y_i(w \cdot x_i + b) \geq 1, \forall x_i$$

- Problem is **convex**, so there is a unique global minimum value (when feasible)
- There is also a unique minimizer, i.e. weight and  $b$  value that provides the minimum
- Quadratic Programming
  - very efficient computationally with procedures that take advantage of the special structure

# Non - linear SVMs: Feature spaces

- General idea: the original feature space can *always* be mapped to some *higher - dimensional* feature space where the training set is *linearly* separable:



# Kernel Trick

---

- If our data isn't linearly separable, we can define a projection  $\Phi(x_i)$  to map it into a much higher dimensional feature space where it is.
- For SVM where everything can be expressed as the dot products of instances this can be done efficiently using the 'kernel trick':
  - A kernel  $K$  is a function such that:  $K(x_i, x_j) = \Phi(x_i) \cdot \Phi(x_j)$
  - Then, we never need to explicitly map the data into the high-dimensional space to solve the optimization problem – magic!!

# Non-Linear Regression

# Non-linear regression

---

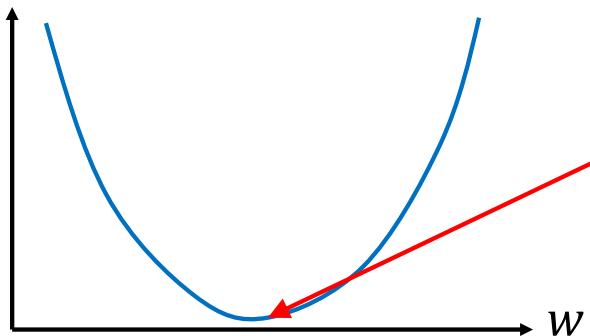
- For many cases when  $f(x; w)$  is not linear w.r.t.  $x$ ,  $f(x; w)$  is not linear w.r.t.  $w$
- Examples:
  - $f(x; w) = e^{wx}$
  - $f(x; w) = \sin(wx)$
- For non-linear regression,  $x$  and  $w$  very likely have different dimensionality
- Example:
  - $f(x; w) = e^{w[1]x} + \sin(w[2]x)$
- We will say  $x$  is D-dimensional and  $w$  is K-dimensional

# When $f(x; w)$ is nonlinear w.r.t. $w$

---

- Squared-error function is no longer parabolic w.r.t. the parameters
  - Can not assume local minimums of error functions are always global minimums
- Often no closed-form solution
  - Need numerical approximation

$$E(w) = \sum_i (wx_i - y_i)^2$$



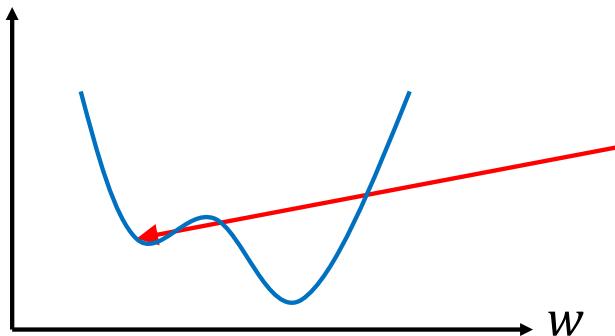
❖  $\frac{\partial E}{\partial w} = 0$  means global minimums  
❖ Closed-forms to find  $w^*$  for  $\frac{\partial E}{\partial w} = 0$

# When $f(x; w)$ is nonlinear w.r.t. $w$

---

- Squared-error function is no longer parabolic w.r.t. the parameters
  - Can not assume local minimums of error functions are always global minimums
- Often no closed-form solution
  - Need numerical approximation

$$E(w) = \sum_i (f(x_i; w) - y_i)^2$$



❖  $\frac{\partial E}{\partial w} = 0$  does not mean global minimums  
❖ No closed-forms to find  $w^*$  for  $\frac{\partial E}{\partial w} = 0$

# Gauss-Newton Algorithm

---

1. Get initial guess for parameters  $\hat{w}^{(1)}$
2. For  $t=1\dots\infty$  (or until we decide to stop)

1. Find error for current guess ( $\hat{w}^{(t)}$ )

$$\Delta y_i = y_i - f(\mathbf{x}_i; \hat{w}^{(t)})$$

2. Find Jacobian around current guess

$$J_i = \frac{df(\mathbf{x}_i; \hat{w}^{(t)})}{dw}$$

3. Calculate change in guess

$$\Delta w = (J^T J)^{-1} J^T \Delta y$$

4. Make new guess

$$\hat{w}^{(t+1)} = \hat{w}^{(t)} + \Delta \hat{w}$$

# Gauss-Newton Algorithm (let $w$ be $K$ -dim)

---

1. Get initial guess for parameters  $\hat{w}^{(1)}$
2. For  $t=1\dots\infty$  (or until we decide to stop)

1. Find error for current guess ( $\hat{w}^{(t)}$ )

$$\Delta y_i = y_i - f(\mathbf{x}_i; \hat{w}^{(t)})$$

2. Find Jacobian around current guess

$$J_{ik} = \frac{\partial f(\mathbf{x}_i; \hat{w}^{(t)})}{\partial w[k]}$$

3. Calculate change in guess

$$\Delta w = (J^T J)^{-1} J^T \Delta y$$

4. Make new guess

$$\hat{w}^{(t+1)} = \hat{w}^{(t)} + \Delta \hat{w}$$

# Gradient Descent (GD)

---

- Simple (and faster?) method for optimization
  - $f(\mathbf{x}; \mathbf{w})$  linear or nonlinear w.r.t. parameters  $\mathbf{w}$
  - Sum or square losses or other losses
- A general method to find local minimum for a multi-parameter error function
  - Search through “parameter space” to find minimum error
- Main idea
  - Evaluate error function for current guess of parameters
  - Determine change in parameters that decrease error
    - From gradient of error function
  - Update parameters in this new change
  - Repeat process until converge (or hit max iterations)

# Algorithm

---

Initialize parameters:

$$\hat{\mathbf{w}}^{(1)} = [w[1], w[2], \dots, w[K]]^T$$

Error function:

$$E(\mathbf{w})$$

---

Compute gradients:

$$\nabla E(\hat{\mathbf{w}}^{(t)}) = \begin{bmatrix} \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[1]} \\ \vdots \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[K]} \end{bmatrix}$$

Update parameters:  $\hat{\mathbf{w}}^{(t+1)} = \hat{\mathbf{w}}^{(t)} - \eta \cdot \nabla E(\hat{\mathbf{w}}^{(t)})$

[Repeat until converges]

Learning rate

## Compute Gradients: Least Squares ( $w$ is 1-dim)

---

Error function:

$$E(\mathbf{w}) = \sum_i [f(\mathbf{x}_i; \mathbf{w}) - y_i]^2$$

Error gradient:

$$\begin{aligned}\nabla E(\mathbf{w}) \\ \frac{dE(\hat{\mathbf{w}}^{(t)})}{dw} &= \frac{\partial}{\partial w} \sum_i [f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i]^2 \\ &= 2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{df(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{dw}\end{aligned}$$

## Compute Gradients: Least Squares ( $w$ is $K$ -dim)

---

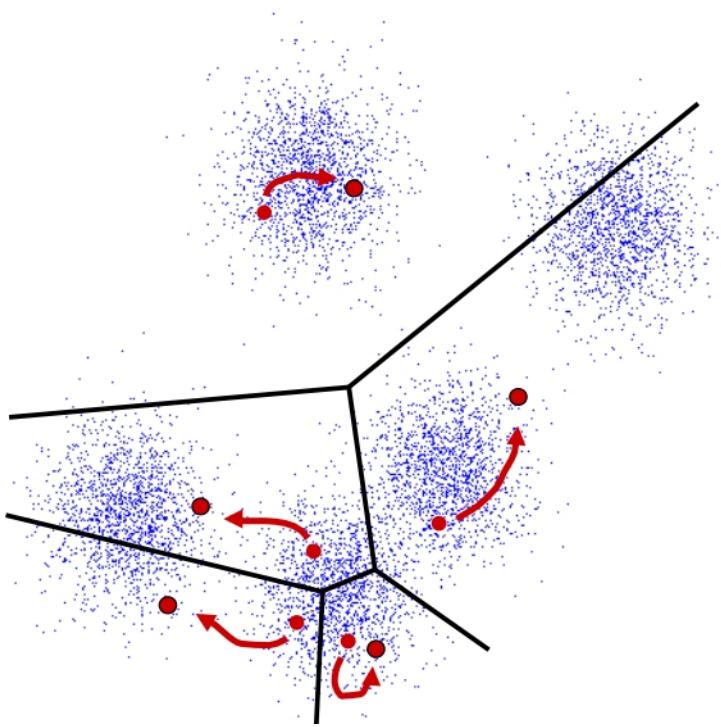
Error function:

$$E(\mathbf{w}) = \sum_i [f(\mathbf{x}_i; \mathbf{w}) - y_i]^2$$

Error gradient:

$$\begin{aligned}\nabla E(\mathbf{w}) \\ \frac{\partial E(\hat{\mathbf{w}}^{(t)})}{\partial w[k]} &= \frac{\partial}{\partial w[k]} \sum_i [f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i]^2 \\ &= 2 \sum_i (f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)}) - y_i) \frac{\partial f(\mathbf{x}_i; \hat{\mathbf{w}}^{(t)})}{\partial w[k]}\end{aligned}$$

# Clustering



THE OHIO STATE UNIVERSITY

# Clustering

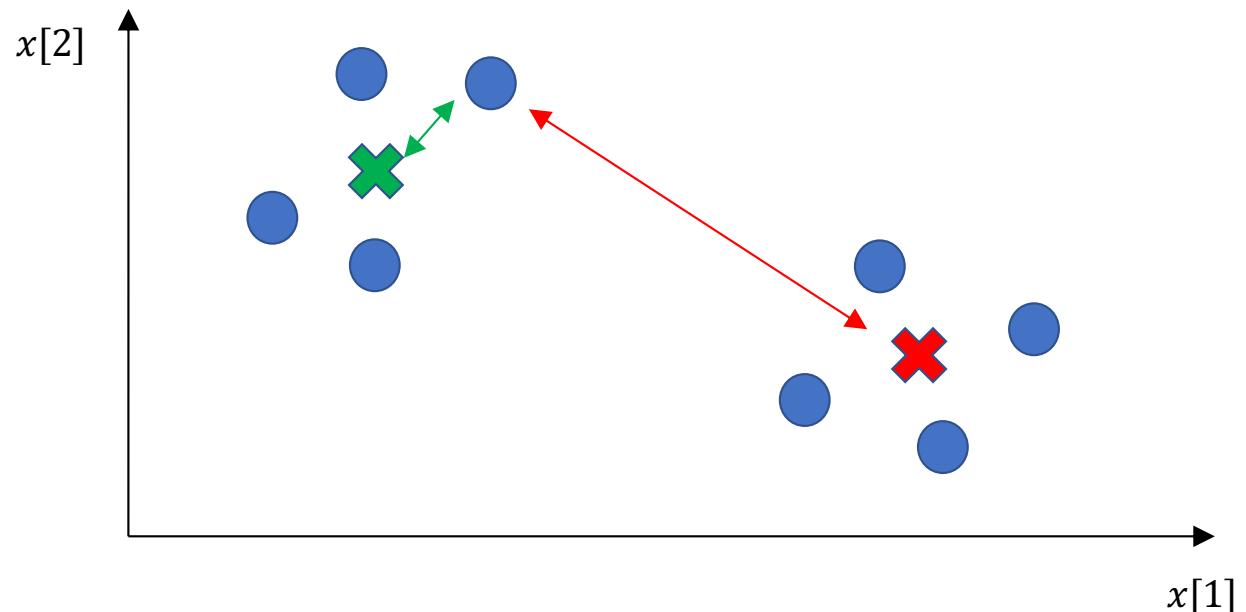
---

- Clustering systems:
  - Unsupervised learning
  - Detect **patterns** in unlabeled data
    - E.g. group emails or search results
    - E.g. find categories of customers
    - E.g. detect anomalous program executions
  - Useful when we don't know what you are looking for
    - Requires data, but no labels
    - May get gibberish

# K-Means

---

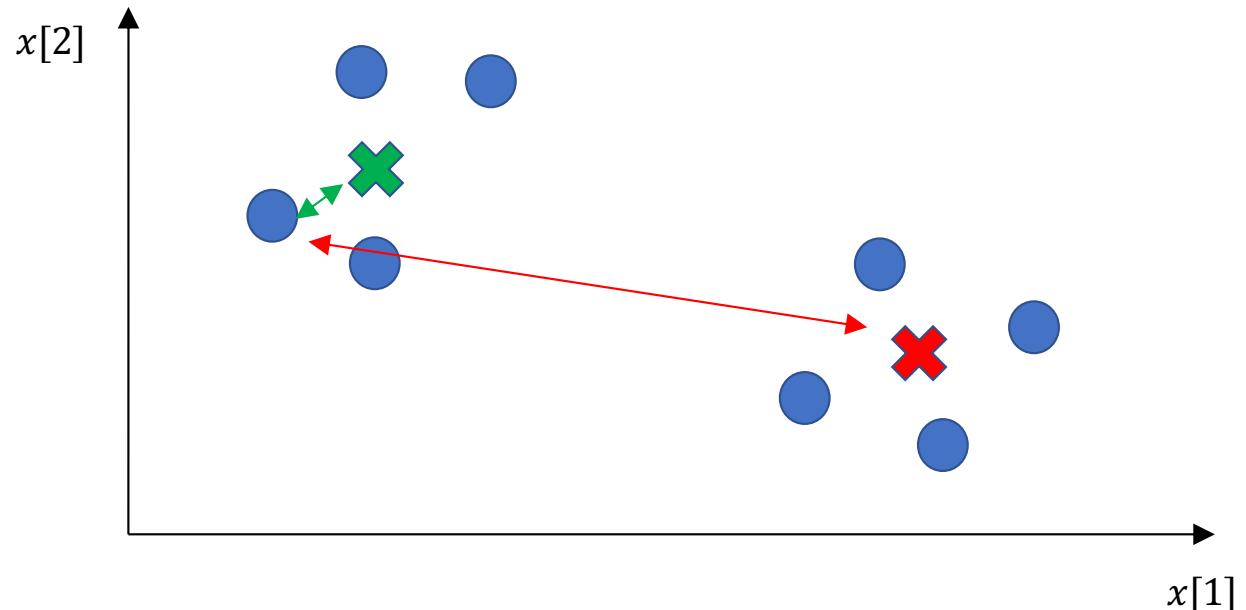
- Find K centers and group data instances into the K centers
  - Example: K = 2
  - Assign each data instance to the closest (most similar, smallest distant) center



# K-Means

---

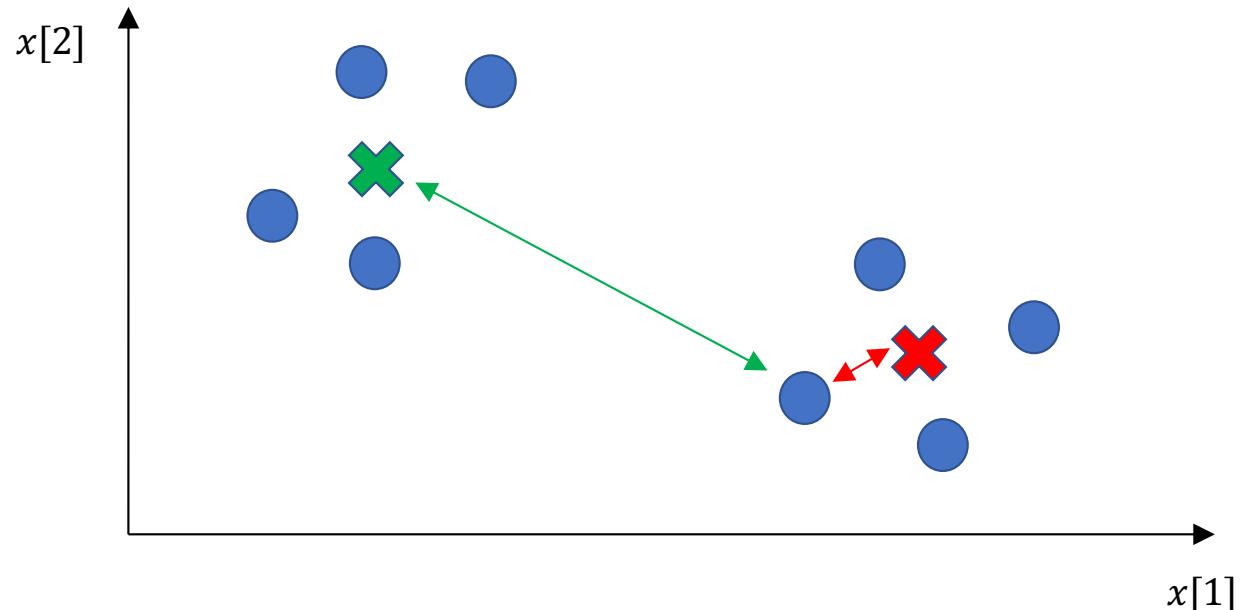
- Find K centers and group data instances into the K centers
  - Example: K = 2
  - Assign each data instance to the closest (most similar, smallest distant) center



# K-Means

---

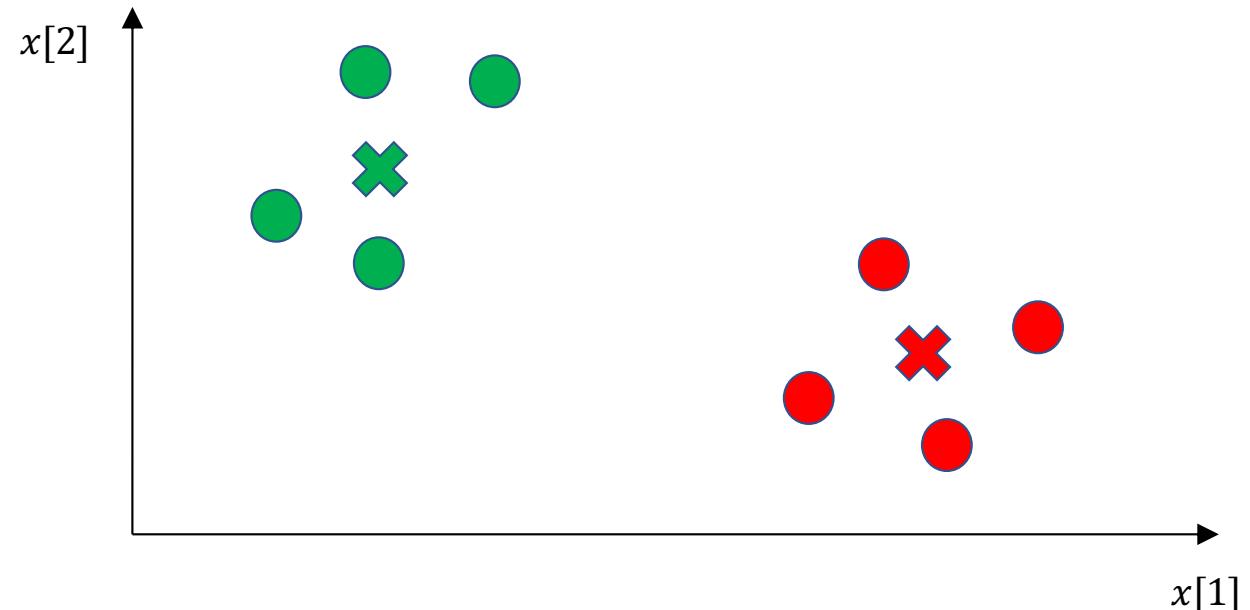
- Find K centers and group data instances into the K centers
  - Example: K = 2
  - Assign each data instance to the closest (most similar, smallest distant) center



# K-Means

---

- Find K centers and group data instances into the K centers
  - Example: K = 2
  - Assign each data instance to the closest (most similar, smallest distant) center



# K-Means

---

- Formulate as an optimization problem (total distance to the centers)

$$\operatorname{argmin}_{\{y_i \in \{1, \dots, K\}\}_{i=1}^N, \{\mathbf{c}_k \in \mathbb{R}^D\}_{k=1}^K} \sum_{k=1}^K \sum_{\{i \mid y_i=k\}} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

All data instances assigned to group k

- To find **centers (K of them)** and **group assignments (N data instances)** simultaneously

# K-Means

---

- Formulate as an optimization problem (total distance to the centers)

$$\operatorname{argmin}_{\{y_i \in \{1, \dots, K\}\}_{i=1}^N, \{\mathbf{c}_k \in \mathbb{R}^D\}_{k=1}^K} \sum_{k=1}^K \sum_{\{i \mid y_i=k\}} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

All data instances assigned to group k

- To find **centers (K of them)** and **group assignments (N data instances)** simultaneously
- Given centers, then for each data point  $i$ , the assignment is independent

$$\operatorname{argmin}_{y_i \in \{1, \dots, K\}} \sum_{k=1}^K \sum_{i|y_i=k} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2 \Rightarrow y_i = \operatorname{argmin}_k \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

# K-Means

---

- Formulate as an optimization problem (total distance to the centers)

$$\operatorname{argmin}_{\{y_i \in \{1, \dots, K\}\}_{i=1}^N, \{\mathbf{c}_k \in \mathbb{R}^D\}_{k=1}^K} \sum_{k=1}^K \sum_{\{i \mid y_i=k\}} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

All data instances assigned to group k

- To find **centers (K of them)** and **group assignments (N data instances)** simultaneously
- Given assignments, then finding the center for each group  $k$  is independent

$$\operatorname{argmin}_{\mathbf{c}_k \in \mathbb{R}^D} \sum_{\{i \mid y_i=k\}} \|\mathbf{x}_i - \mathbf{c}_k\|_2^2 \Rightarrow \mathbf{c}_k = \frac{1}{|\{i \mid y_i=k\}|} \sum_{\{i \mid y_i=k\}} \mathbf{x}_i$$

# K-Means: an iterative clustering algorithm

---

- Pick  $K$  random points as cluster centers (means):  $c_1 \dots c_K$
- Alternate (the error is monotonically decreasing):
  - Assign each example  $x_i$  to the mean  $c_k$  that is closest to it

$$\operatorname{argmin}_{y_i \in \{1, \dots, K\}} \sum_{k=1}^K \sum_{i | y_i = k} \|\textcolor{brown}{x}_i - \textcolor{violet}{c}_k\|_2^2 \Rightarrow y_i = \operatorname{argmin}_k \|\textcolor{brown}{x}_i - c_k\|_2^2$$

- Set each mean  $c_k$  to the average of its assigned points

$$\operatorname{argmin}_{c_k \in \mathbb{R}^D} \sum_{\{i | y_i = k\}} \|\textcolor{brown}{x}_i - \textcolor{violet}{c}_k\|_2^2 \Rightarrow c_k = \frac{1}{|\{i | y_i = k\}|} \sum_{\{i | y_i = k\}} x_i$$

- Stop when no points assignments change

# K-Means

---

- Data:  $\{\mathbf{x}_1, \dots, \mathbf{x}_N\}$
- An iterative clustering algorithm
  - Pick K random cluster centers,  $\mathbf{c}_1 \dots \mathbf{c}_K$
  - For  $t=1 \dots T$ : [or, stop if assignments don't change]
    - for  $i = 1 \dots N$ : [update cluster assignments]

$$y_i = \operatorname{argmin}_k \|\mathbf{x}_i - \mathbf{c}_k\|_2^2$$

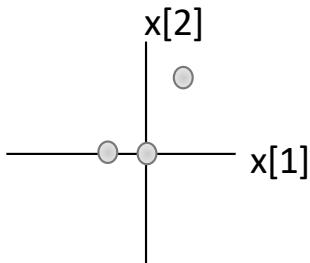
- for  $k=1 \dots K$ : [update cluster centers]

$$\mathbf{c}_k = \frac{1}{|\{i \mid y_i = k\}|} \sum_{\{i \mid y_i = k\}} \mathbf{x}_i$$

# K-Means

---

|     | x[1] | x[2] |
|-----|------|------|
| i=1 | -1   | 0    |
| i=2 | 0    | 0    |
| i=3 | 2    | 2    |



Random cluster means:

$$c_1 = [-1, 0], c_2 = [0, 0]$$

$t = 0$ :

|       | i=1 | i=2 | i=3 |
|-------|-----|-----|-----|
| $c_1$ | 0   | 1   | 13  |
| $c_2$ | 1   | 0   | 8   |

- $y_1 = \operatorname{argmin}_k \operatorname{dist}(x_1, c_k) = 1$
- $y_2 = \operatorname{argmin}_k \operatorname{dist}(x_2, c_k) = 1$
- $y_3 = \operatorname{argmin}_k \operatorname{dist}(x_3, c_k) = 2$
  
- $c_1 = (1/1) * [-1, 0] = [-1, 0]$
- $c_2 = ((1/2) * ([0, 0] + [2, 2])) = [1, 1]$

$$c_1 = [-1, 0], c_2 = [1, 1]$$

$t = 1$ :

|       | i=1 | i=2 | i=3 |
|-------|-----|-----|-----|
| $c_1$ | 0   | 1   | 13  |
| $c_2$ | 4   | 4   | 18  |

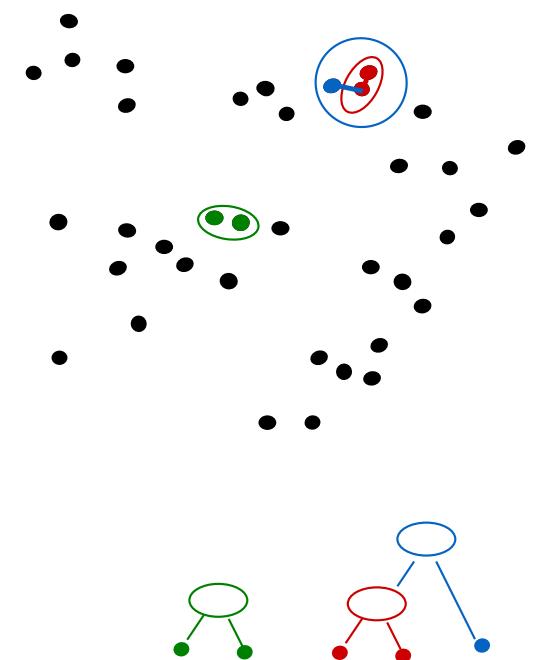
- $y_1 = \operatorname{argmin}_k \operatorname{dist}(x_1, c_k) = 1$
- $y_2 = \operatorname{argmin}_k \operatorname{dist}(x_2, c_k) = 1$
- $y_3 = \operatorname{argmin}_k \operatorname{dist}(x_3, c_k) = 2$
  
- $c_1 = (1/2) * ([-1, 0] + [0, 0]) = [-0.5, 0]$
- $c_2 = (1/1) * ([2, 2]) = [2, 1]$

- $t=2$ :
- **Stop!!**
- $y_i$  won't change (verify!)

# Agglomerative Clustering

---

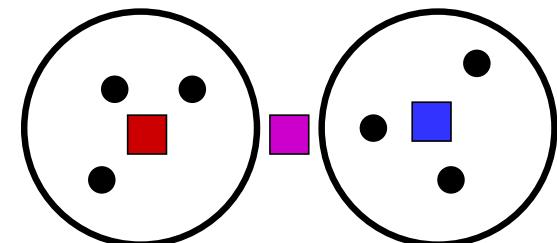
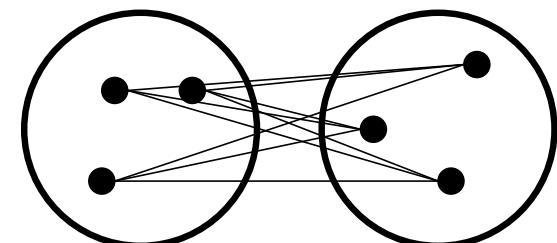
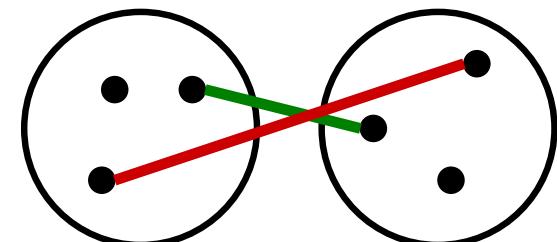
- Agglomerative clustering:
  - First merge very similar instances
  - Incrementally build larger clusters out of smaller clusters
- Algorithm:
  - Maintain a set of clusters
  - Initially, each instance in its own cluster
  - Repeat:
    - Pick the two **closest** clusters
    - Merge them into a new cluster
    - Stop when there is only one cluster left
- Produces not one clustering, but a family of clusters represented by a **dendrogram**



# Agglomerative Clustering

---

- How should we define “closest” for clusters with multiple elements?
- Many options
  - Closest pair (single-link clustering)
  - Farthest pair (complete-link clustering)
  - Average of all pairs
  - Ward’s method (min variance, like k-means)
- Different choices create different clustering behaviors



# Agglomerative Clustering

---

complete link



single link



*Each discrete value along the x-axis is treated as one data instance.*

# Agglomerative Clustering

---

complete link



single link

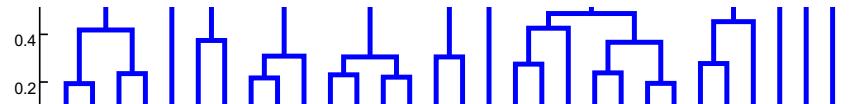


*Each discrete value along the x-axis is treated as one data instance.*

# Agglomerative Clustering

---

complete link



single link

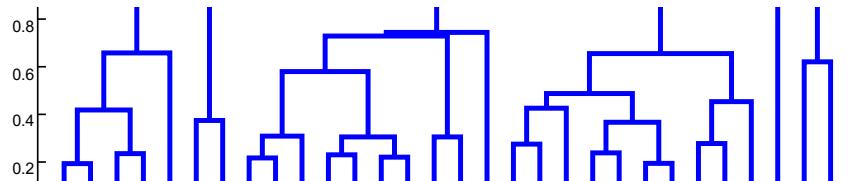


*Each discrete value along the x-axis is treated as one data instance.*

# Agglomerative Clustering

---

complete link



single link

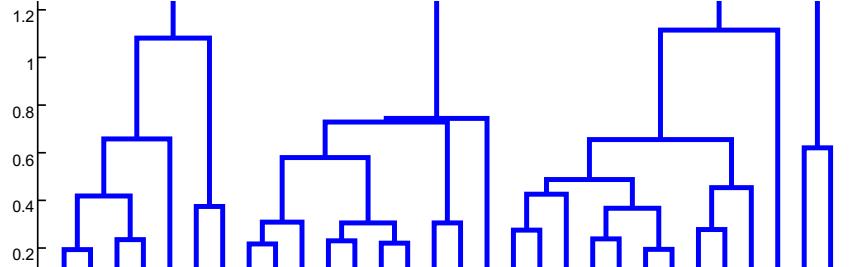


*Each discrete value along the x-axis is treated as one data instance.*

# Agglomerative Clustering

---

complete link



single link

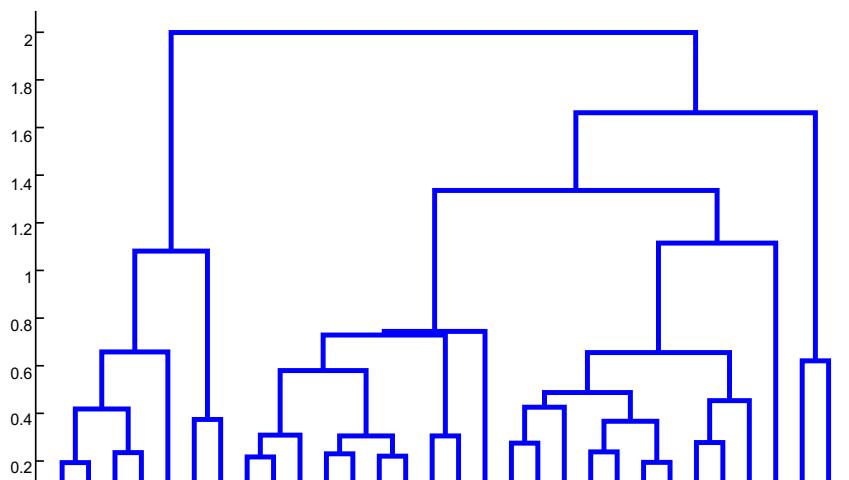


*Each discrete value along the x-axis is treated as one data instance.*

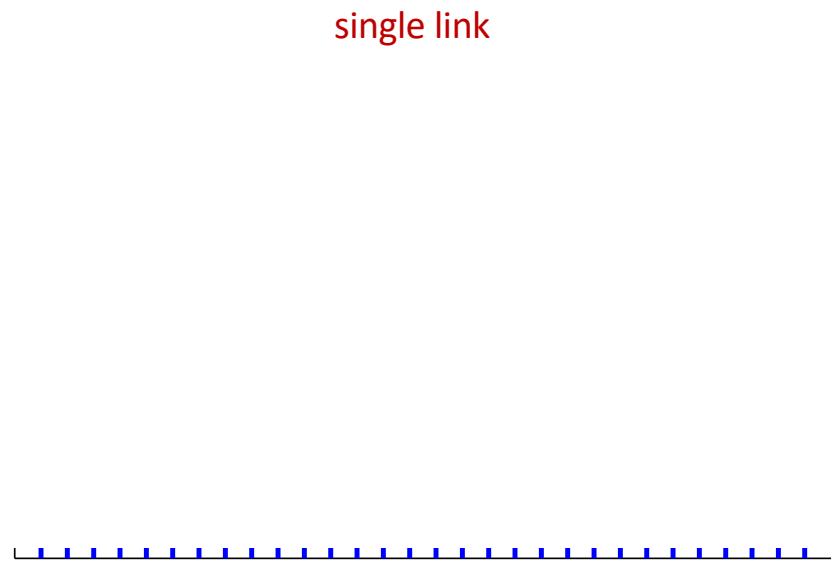
# Agglomerative Clustering

---

complete link



single link

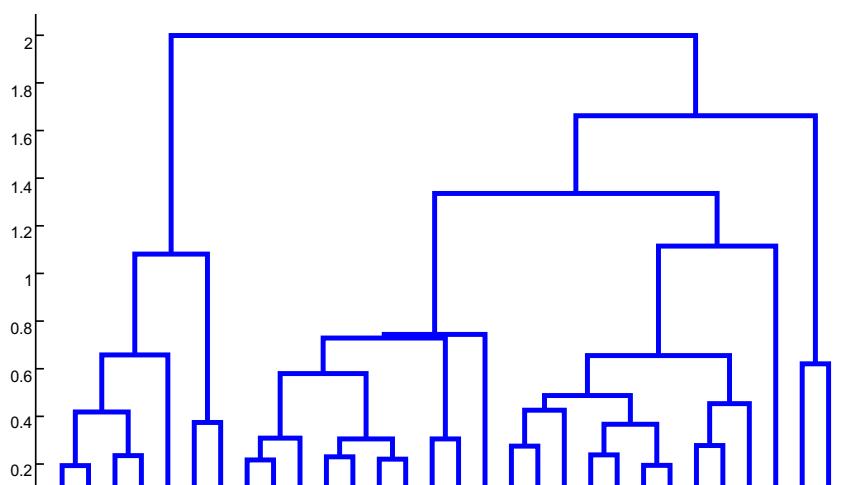


*Each discrete value along the x-axis is treated as one data instance.*

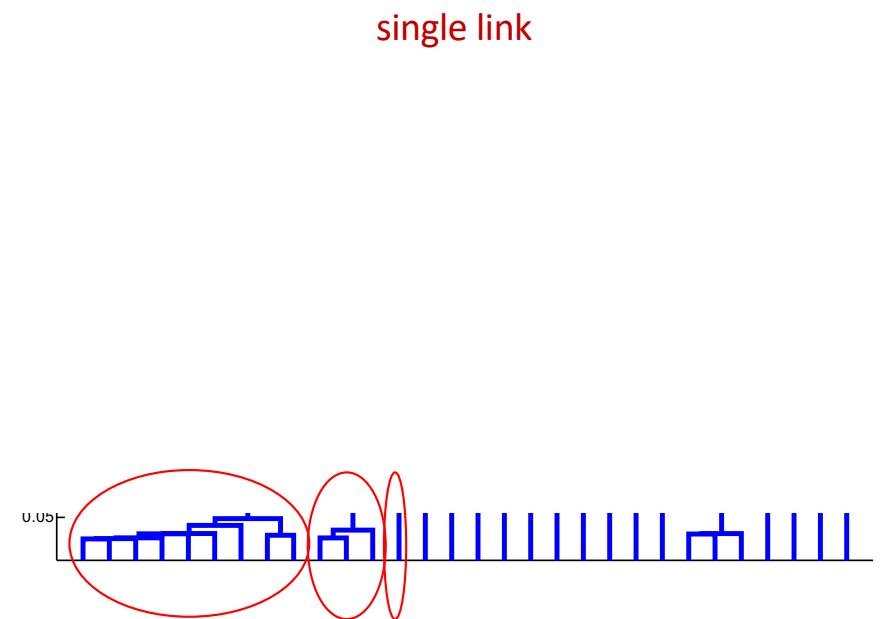
# Agglomerative Clustering

---

complete link



single link

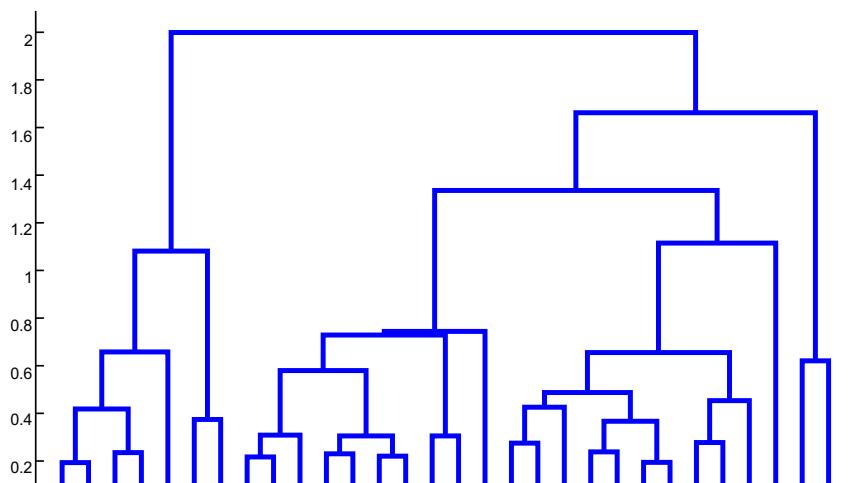


*Each discrete value along the x-axis is treated as one data instance.*

# Agglomerative Clustering

---

complete link



single link

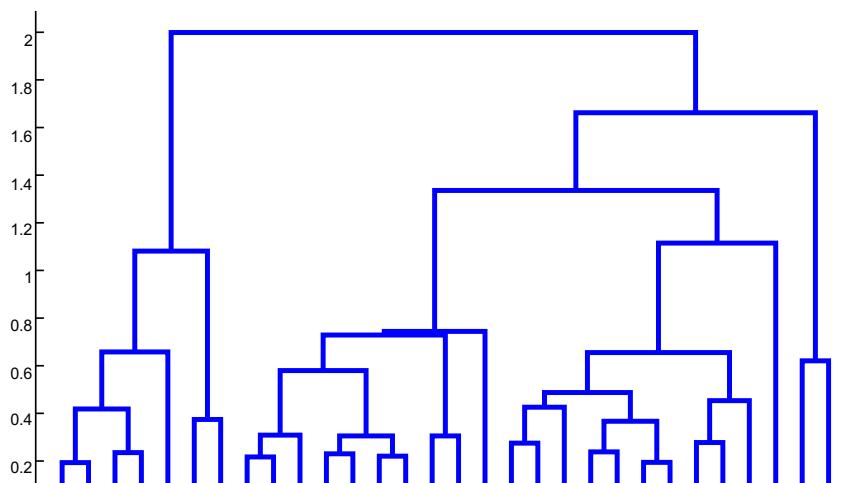


*Each discrete value along the x-axis is treated as one data instance.*

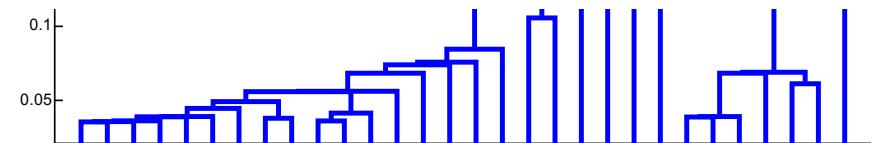
# Agglomerative Clustering

---

complete link



single link

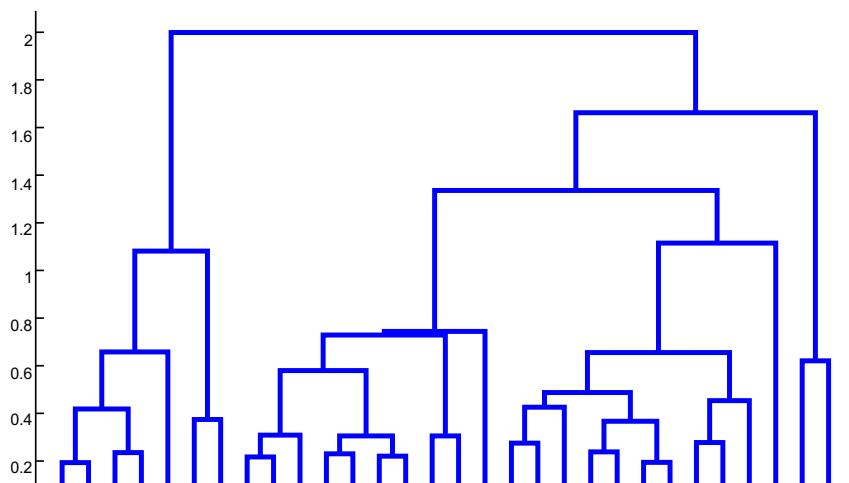


*Each discrete value along the x-axis is treated as one data instance.*

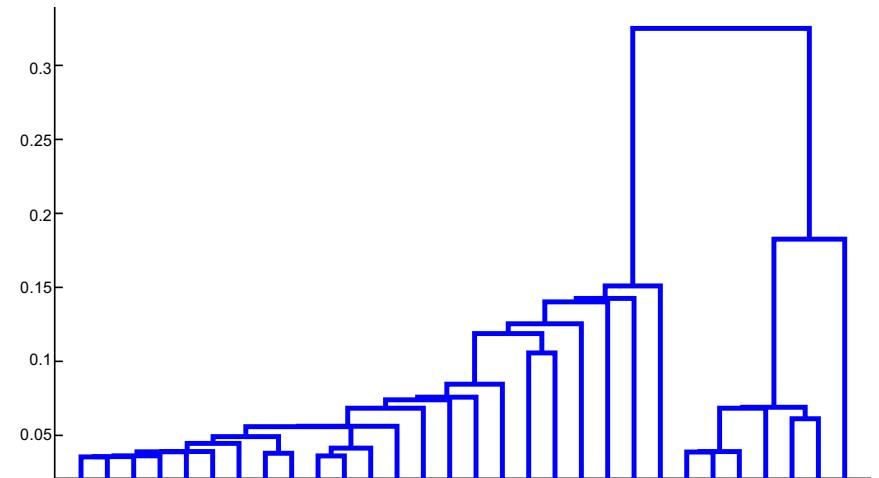
# Agglomerative Clustering

---

complete link



single link

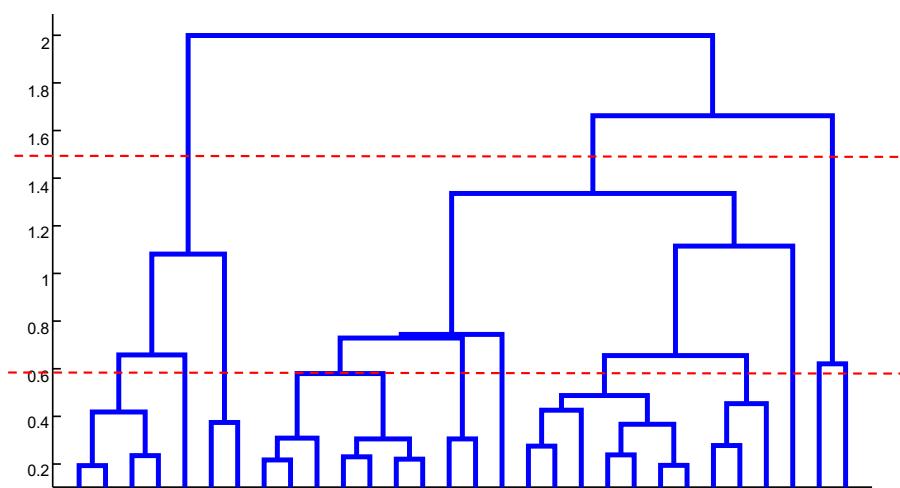


*Each discrete value along the x-axis is treated as one data instance.*

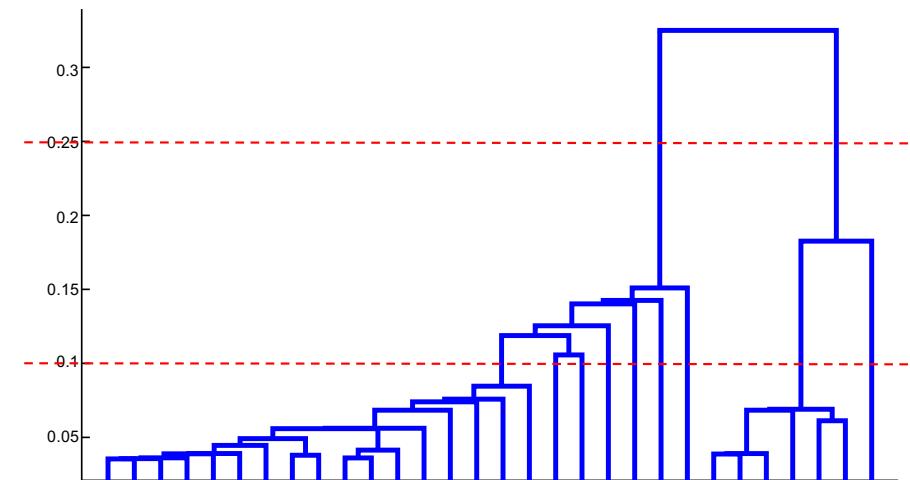
# Agglomerative Clustering

---

complete link



single link

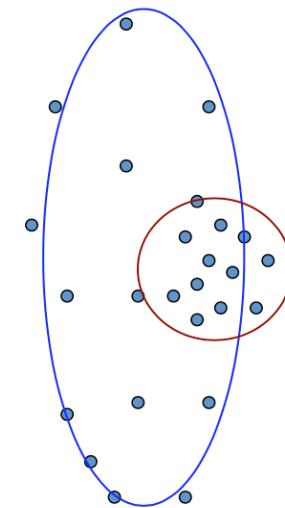


*Each discrete value along the x-axis is treated as one data instance.*

## (One) bad case for “hard assignments”?

---

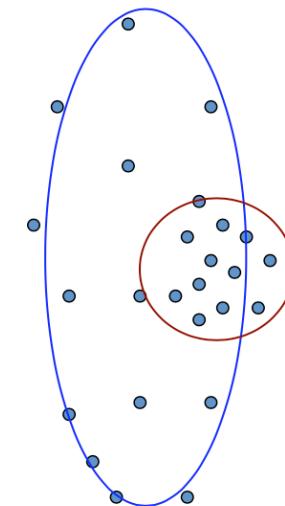
- Clusters may overlap
- Some clusters may be “wider” than others



# (One) bad case for “hard assignments”?

---

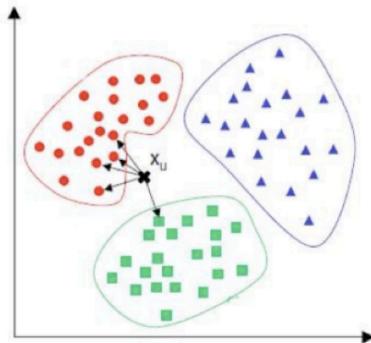
- Clusters may overlap
- Some clusters may be “wider” than others
- *Distances can be deceiving!*



KNN

# K-Nearest Neighbor

- Given training data  $\mathcal{D} = \{(\mathbf{x}_1, \mathbf{y}_1), \dots, (\mathbf{x}_N, \mathbf{y}_N)\}$  and a test point
- Prediction Rule: Look at the  $K$  most similar training examples



- For classification: assign the majority class label (**majority voting**)
- For regression: assign the **average response**
- The algorithm requires:
  - Parameter  $K$ : number of nearest neighbors to look for
  - Distance function**: To compute the similarities between examples

# K-Nearest Neighbor Algorithm

- Compute the test point's distance from each training point
- Sort the distances in ascending (or descending) order
- Use the sorted distances to select the  $K$  nearest neighbors
- Use **majority rule** (for classification) or **averaging** (for regression)

**Note:**  $K$ -Nearest Neighbors is called a *non-parametric* method

- Unlike other supervised learning algorithms,  $K$ -Nearest Neighbors doesn't learn an explicit mapping  $f$  from the training data
- It simply uses the training data at the test time to make predictions

# K-NN: Compute Distance

---

- The  $K$ -NN algorithm requires computing distances of the test example from each of the training examples
- Several ways to compute distances
- The choice depends on the **type of the features** in the data
- Real-valued features ( $\mathbf{x}_i \in \mathbb{R}^D$ ): **Euclidean distance** is commonly used

$$d(\mathbf{x}_i, \mathbf{x}_j) = \sqrt{\sum_{m=1}^D (x_{im} - x_{jm})^2} = \sqrt{\|\mathbf{x}_i\|^2 + \|\mathbf{x}_j\|^2 - 2\mathbf{x}_i^T \mathbf{x}_j}$$

- Generalization of the distance between points in 2 dimensions
- $\|\mathbf{x}_i\| = \sqrt{\sum_{m=1}^D x_{im}^2}$  is called the **norm** of  $\mathbf{x}_i$ 
  - Norm of a vector  $\mathbf{x}$  is also its **length**

# K-NN: Feature Normalization

---

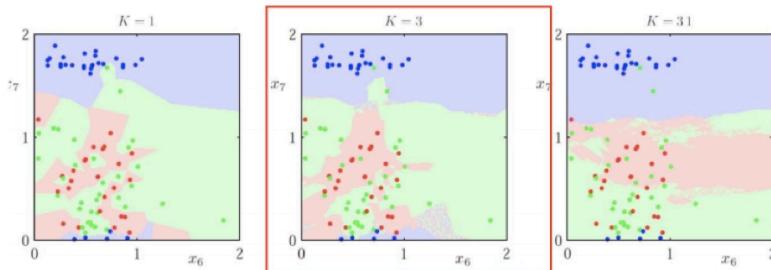
- Note: Features should be on the same scale
- Example: if one feature has its values in millimeters and another has in centimeters, we would need to normalize
- One way is:
  - Replace  $x_{im}$  by  $z_{im} = \frac{(x_{im} - \bar{x}_m)}{\sigma_m}$  (make them zero mean, unit variance)
  - $\bar{x}_m = \frac{1}{N} \sum_{i=1}^N x_{im}$ : empirical mean of  $m^{th}$  feature
  - $\sigma_m^2 = \frac{1}{N} \sum_{i=1}^N (x_{im} - \bar{x}_m)^2$ : empirical variance of  $m^{th}$  feature

# K-NN: Other Distance Measure

---

- Binary-valued features
  - Use Hamming distance:  $d(x_i, x_j) = \sum_{m=1}^D \mathbb{I}(x_{im} \neq x_{jm})$
  - Hamming distance counts the number of features where the two examples disagree
- Mixed feature types (some real-valued and some binary-valued)?
  - Can use mixed distance measures
  - E.g., Euclidean for the real part, Hamming for the binary part
- Can also assign **weights** to features:  $d(x_i, x_j) = \sum_{m=1}^D w_m d(x_{im}, x_{jm})$

# K-NN: Choice of K



- Small  $K$ 
  - Creates many small regions for each class
  - May lead to non-smooth) decision boundaries and overfit
- Large  $K$ 
  - Creates fewer larger regions
  - Usually leads to smoother decision boundaries (caution: too smooth decision boundary can underfit)
- Choosing  $K$ 
  - Often data dependent and heuristic based
  - Or using [cross-validation](#) (using some **held-out data**)
  - In general, a  $K$  too small or too big is bad!

# K-NN: Pseudocode

1. Calculate “ $d(x, x_i)$ ”  $i = 1, 2, \dots, n$ ; where  $d$  denotes the Euclidean distance between the points.
2. Arrange the calculated  $n$  Euclidean distances in non-decreasing order.
3. Let  $k$  be a +ve integer, take the first  $k$  distances from this sorted list.
4. Find those  $k$ -points corresponding to these  $k$ -distances.
5. Let  $k_i$  denotes the number of points belonging to the  $i^{\text{th}}$  class among  $k$  points i.e.  $k \geq 0$
6. If  $k_i > k_j \forall i \neq j$  then put  $x$  in class  $i$ .

Let  $(X_i, C_i)$  where  $i = 1, 2, \dots, n$  be data points.  $X_i$  denotes feature values &  $C_i$  denotes labels for  $X_i$  for each  $i$ .

Assuming the number of classes as ‘c’  
 $c_i \in \{1, 2, 3, \dots, c\}$  for all values of  $i$

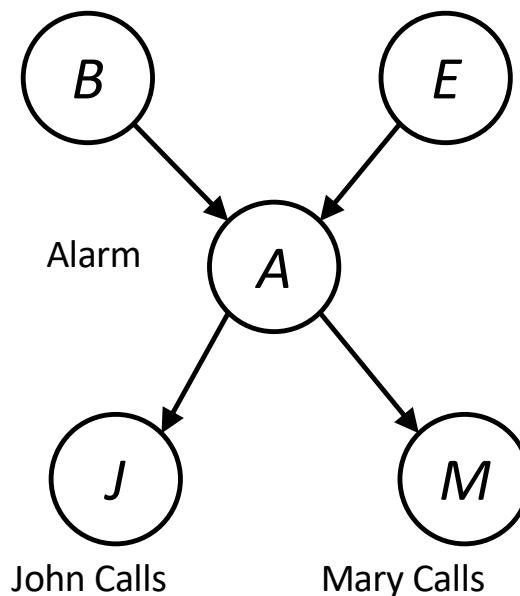
Let  $x$  be a point for which label is not known, and we would like to find the label class using  $k$ -nearest neighbor algorithms.

# Bayes Network

# Inference on PGMs

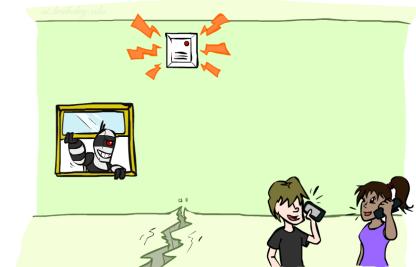
Burglary

| B  | P(B)  |
|----|-------|
| +b | 0.001 |
| -b | 0.999 |



| E  | P(E)  |
|----|-------|
| +e | 0.002 |
| -e | 0.998 |

Earthquake



| A  | J  | P(J A) |
|----|----|--------|
| +a | +j | 0.9    |
| +a | -j | 0.1    |
| -a | +j | 0.05   |
| -a | -j | 0.95   |

| A  | M  | P(M A) |
|----|----|--------|
| +a | +m | 0.7    |
| +a | -m | 0.3    |
| -a | +m | 0.01   |
| -a | -m | 0.99   |

| B  | E  | A  | P(A B,E) |
|----|----|----|----------|
| +b | +e | +a | 0.95     |
| +b | +e | -a | 0.05     |
| +b | -e | +a | 0.94     |
| +b | -e | -a | 0.06     |
| -b | +e | +a | 0.29     |
| -b | +e | -a | 0.71     |
| -b | -e | +a | 0.001    |
| -b | -e | -a | 0.999    |

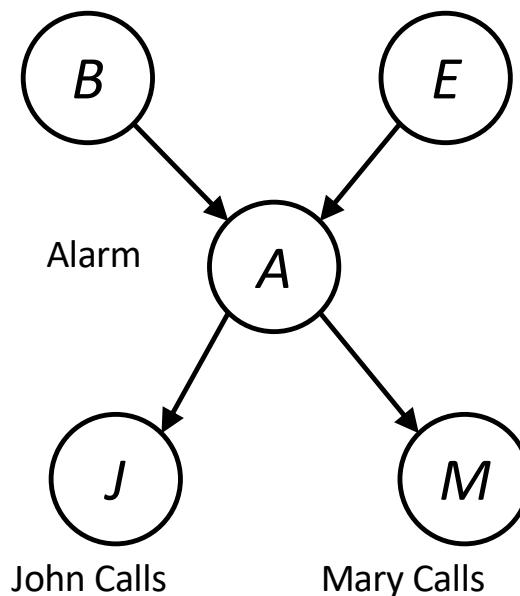
$$P(+b, -e, +a, -j, +m) =$$

$$P(+b)P(-e)P(+a|+b, -e)P(-j|+a)P(+m|+a) =$$

# Inference on PGMs

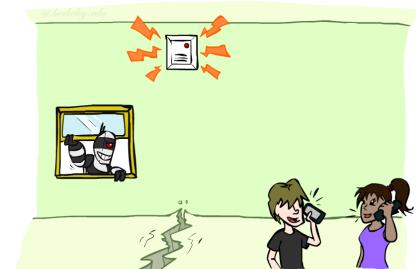
Burglary

| B  | P(B)  |
|----|-------|
| +b | 0.001 |
| -b | 0.999 |



| E  | P(E)  |
|----|-------|
| +e | 0.002 |
| -e | 0.998 |

Earthquake



| A  | J  | P(J A) |
|----|----|--------|
| +a | +j | 0.9    |
| +a | -j | 0.1    |
| -a | +j | 0.05   |
| -a | -j | 0.95   |

| A  | M  | P(M A) |
|----|----|--------|
| +a | +m | 0.7    |
| +a | -m | 0.3    |
| -a | +m | 0.01   |
| -a | -m | 0.99   |

$$P(+b, -e, +a, -j, +m) =$$

$$P(+b)P(-e)P(+a|+b, -e)P(-j|+a)P(+m|+a) =$$

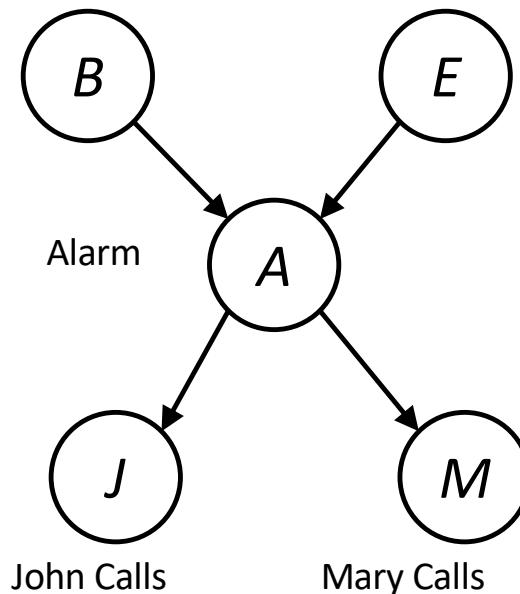
$$0.001 \times 0.998 \times 0.94 \times 0.1 \times 0.7$$

| B  | E  | A  | P(A B,E) |
|----|----|----|----------|
| +b | +e | +a | 0.95     |
| +b | +e | -a | 0.05     |
| +b | -e | +a | 0.94     |
| +b | -e | -a | 0.06     |
| -b | +e | +a | 0.29     |
| -b | +e | -a | 0.71     |
| -b | -e | +a | 0.001    |
| -b | -e | -a | 0.999    |

# Inference on PGMs

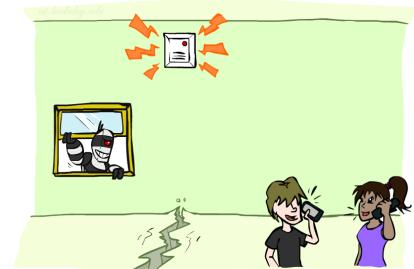
Burglary

| B  | P(B)  |
|----|-------|
| +b | 0.001 |
| -b | 0.999 |



| E  | P(E)  |
|----|-------|
| +e | 0.002 |
| -e | 0.998 |

Earthquake



| A  | J  | P(J A) |
|----|----|--------|
| +a | +j | 0.9    |
| +a | -j | 0.1    |
| -a | +j | 0.05   |
| -a | -j | 0.95   |

| A  | M  | P(M A) |
|----|----|--------|
| +a | +m | 0.7    |
| +a | -m | 0.3    |
| -a | +m | 0.01   |
| -a | -m | 0.99   |

| B  | E  | A  | P(A B,E) |
|----|----|----|----------|
| +b | +e | +a | 0.95     |
| +b | +e | -a | 0.05     |
| +b | -e | +a | 0.94     |
| +b | -e | -a | 0.06     |
| -b | +e | +a | 0.29     |
| -b | +e | -a | 0.71     |
| -b | -e | +a | 0.001    |
| -b | -e | -a | 0.999    |

$$P(B | +j, +m) = \frac{P(B, +j, +m)}{P(+j, +m)} = ?$$

# Markov Blanket

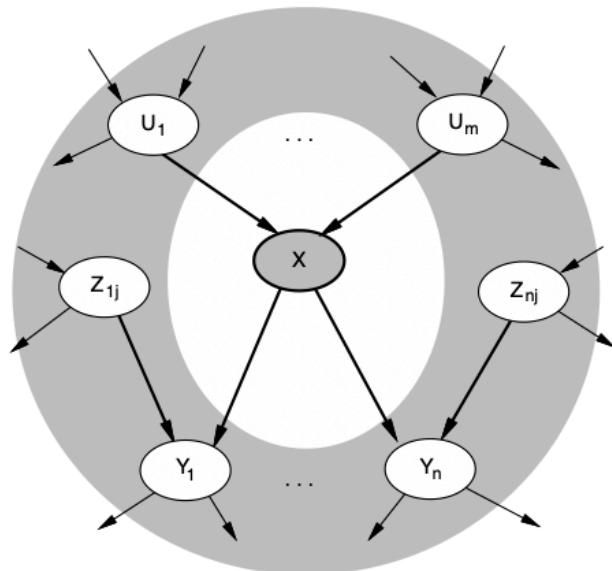
---

- Definition:
  - The smallest set of “observed” nodes that renders a node  $t$  conditionally independent of all the other nodes in the graph.
- Markov blanket in DAG is:
  - Parents
  - Children
  - Co-parents (other nodes that are also parents of the children)

# Markov Blanket

---

- Each node is conditionally independent of all others given its Markov blanket:
  - parents + children + children’s parents



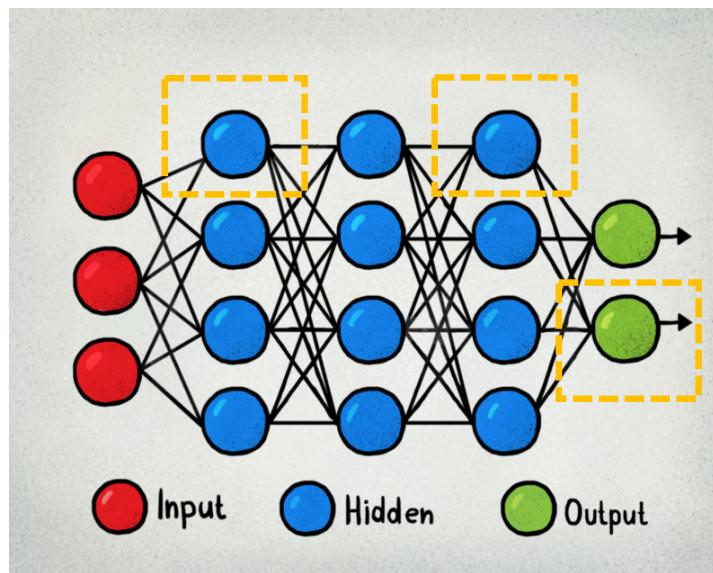
# Neural Network

# “Artificial” Neuron

---

- A basic “computational unit” of neural networks

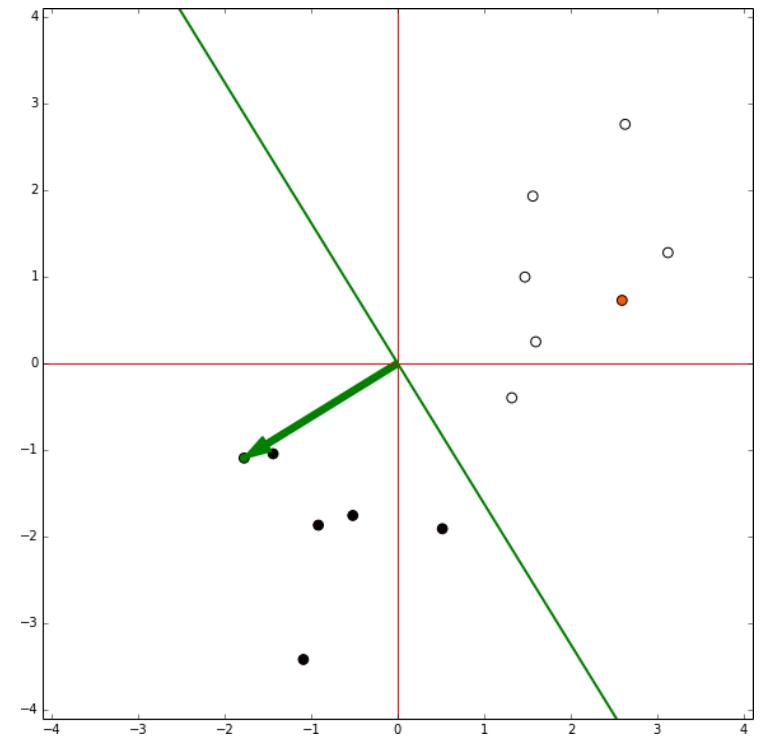
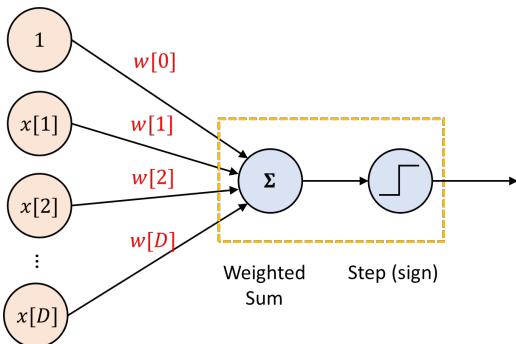
$$x \quad z^{(1)} \quad z^{(2)} \quad z^{(3)} \quad y$$



# Perceptron Algorithm

- Let  $x \in \mathbb{R}^D$  and  $w \in \mathbb{R}^D$  ( $b$  is merged into  $w$ ),  $y \in \{-1,1\}$

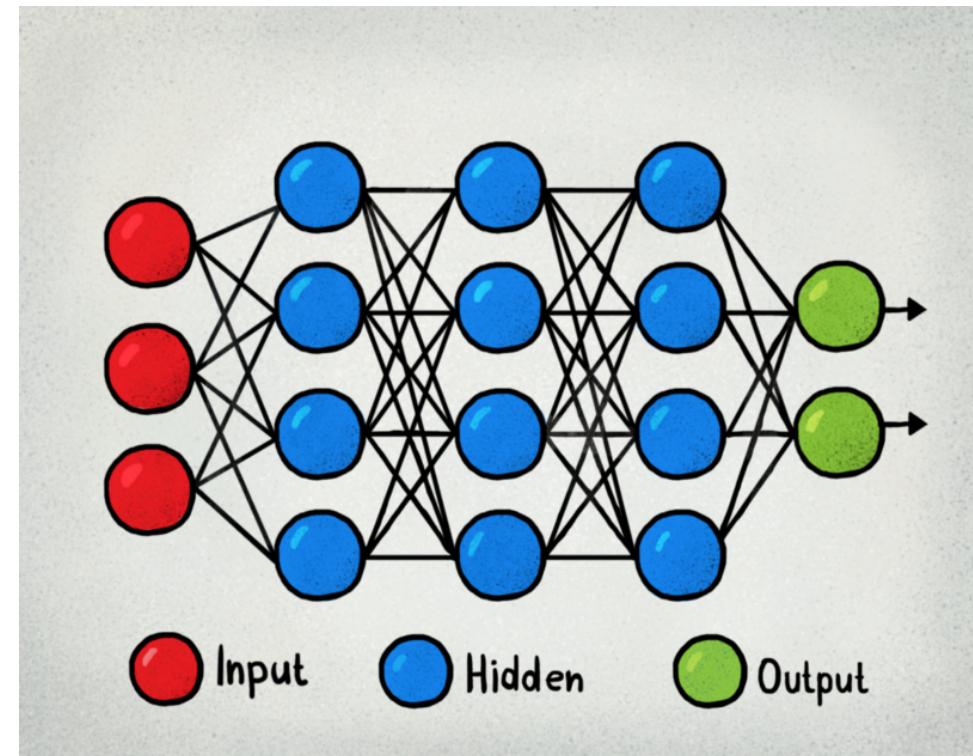
- Initialize weight vector  $w = \mathbf{0}$
- Loop for T iterations
  - Loop for all training examples  $x_n$  (random order!)
  - Predict  $\hat{y}_n = \text{sign}(w^T x_n)$
  - If  $\hat{y}_n \neq y_n$ 
    - Update:  $w \leftarrow w + \eta(y_n x_n)$



# Multi-layer perceptron

---

- How to learn?
- Perceptron algorithm?
  - Only look at one instance a time
  - Can update only a single perceptron
  - Not converge for not linear separable data



# Training a neural network: gradient descent

---

- Let  $x \in \mathbb{R}^D$ ;  $y$  as true label, and  $\{(x_n, y_n)\}$  as the training data
- $\theta$  as all parameters,  $\hat{y} = f_\theta(x)$  as the neural network's prediction
- Initialize  $\theta$  with [with some specific methods]
- Loop for T “epochs”
  - $\nabla_\theta L = \frac{1}{N} \sum_{n=1}^N \nabla_\theta l(y_n, \hat{y}_n)$
  - $\theta \leftarrow \theta - \eta \times \nabla_\theta L$

# Stochastic gradient descent

---

- Let  $\mathbf{x} \in \mathbb{R}^D$ ;  $y$  as true label, and  $\{(\mathbf{x}_n, y_n)\}$  as the training data
- $\boldsymbol{\theta}$  as all parameters,  $\hat{y} = f_{\boldsymbol{\theta}}(\mathbf{x})$  as the neural network's prediction
- Initialize  $\boldsymbol{\theta}$  with [with some specific methods]
- Loop for T “epochs”
  - Loop for all training examples  $\mathbf{x}_n$  (random order!)
  - $\nabla_{\boldsymbol{\theta}} L = \nabla_{\boldsymbol{\theta}} l(y_n, \hat{y}_n)$
  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \times \nabla_{\boldsymbol{\theta}} L$

# “Mini-batch” Stochastic gradient descent

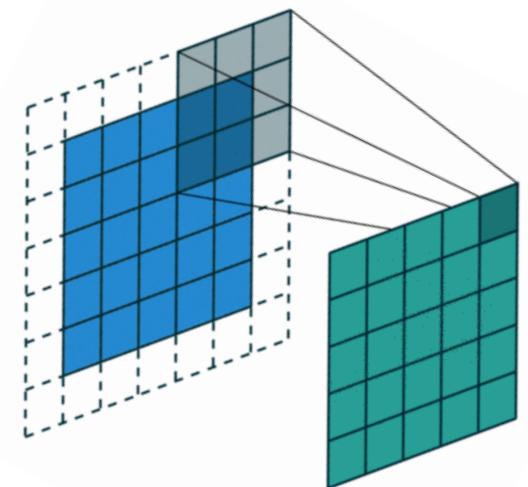
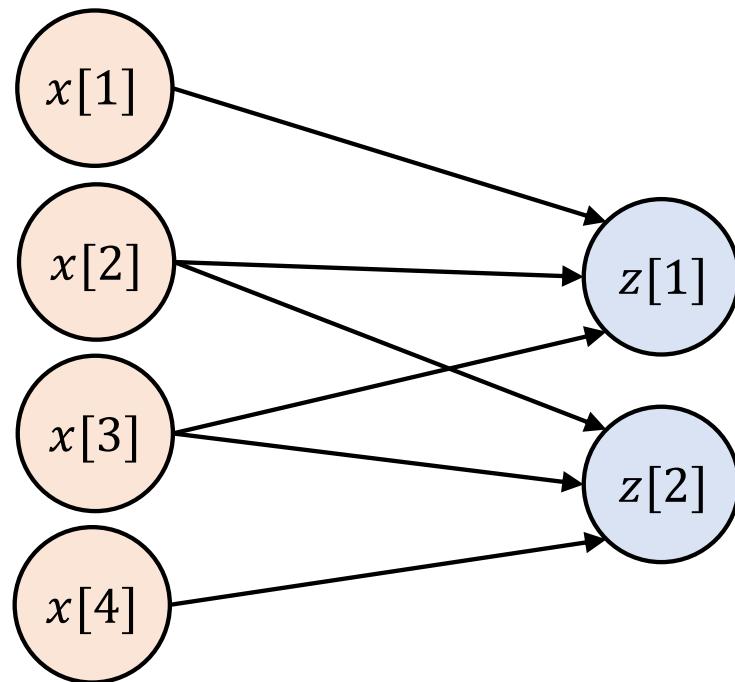
---

- Let  $\mathbf{x} \in \mathbb{R}^D$ ;  $y$  as true label, and  $\{(\mathbf{x}_n, y_n)\}$  as the training data
- $\boldsymbol{\theta}$  as all parameters,  $\hat{y} = f_{\boldsymbol{\theta}}(\mathbf{x})$  as the neural network’s prediction
- Initialize  $\boldsymbol{\theta}$  with [with some specific methods]
- Loop for T “epochs”
  - Loop for “B sampled examples from  $\{(\mathbf{x}_n, y_n)\}$ ” (called batch) without replacement (**random order!**)
  - $\nabla_{\boldsymbol{\theta}} L = \frac{1}{N} \sum_{b=1}^B \nabla_{\boldsymbol{\theta}} l(y_b, \hat{y}_b)$
  - $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \times \nabla_{\boldsymbol{\theta}} L$

# Convolutional Neural Networks

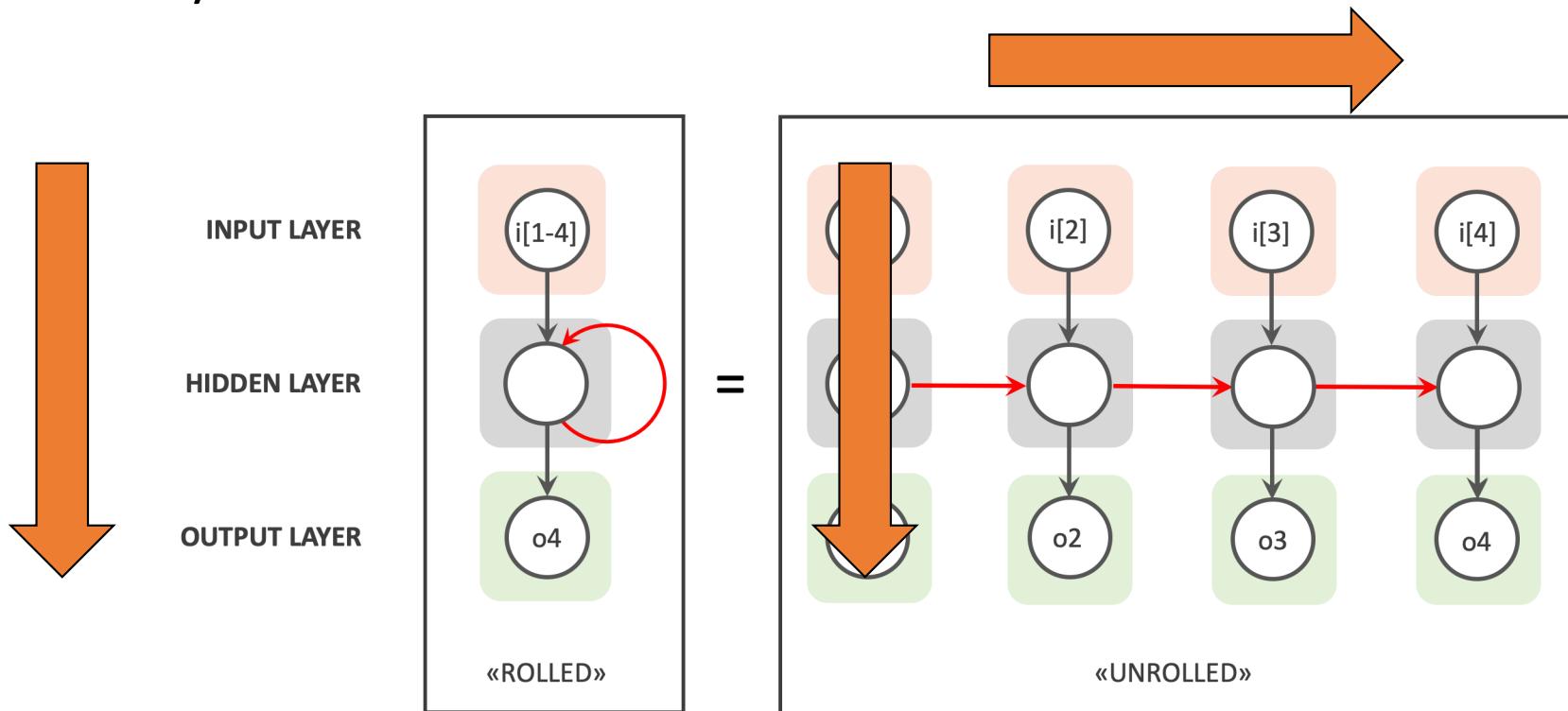
---

- A current node is not directly affected by “all nodes in the previous layer”
- The network “weights” on the edges can be “re-used”



# Recurrent Neural Network

- Other styles: Recurrent neural network



A simple cartoon illustration of a white cat with black outlines. The cat has a small tuft of hair on its head and is wearing a thick, orange and white striped scarf. It is standing upright with its front paws slightly apart.

GOOD LUCK  
ON YOUR EXAMS!  
I BELIEVE IN YOU,  
AND I THINK YOU  
CAN DO REALLY WELL!