

CS6510 Applied Machine Learning

Classifiers

20 Aug 2016

Vineeth N Balasubramanian



Classification Methods

- k-Nearest Neighbors
- Decision Trees
- Naïve Bayes
- Support Vector Machines
- Logistic Regression
- Neural Networks
- Ensemble Methods (Boosting, Random Forests)

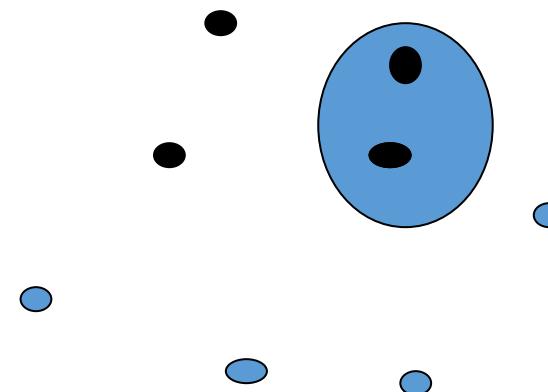
k-Nearest Neighbors

- Eager Learning (**Induction**)
 - Explicit description of target function on the whole training set
- Instance-based Learning (**Transduction**)
 - Learning=storing all training instances
 - Classification=assigning target function to a new instance
 - Referred to as “Lazy” learning

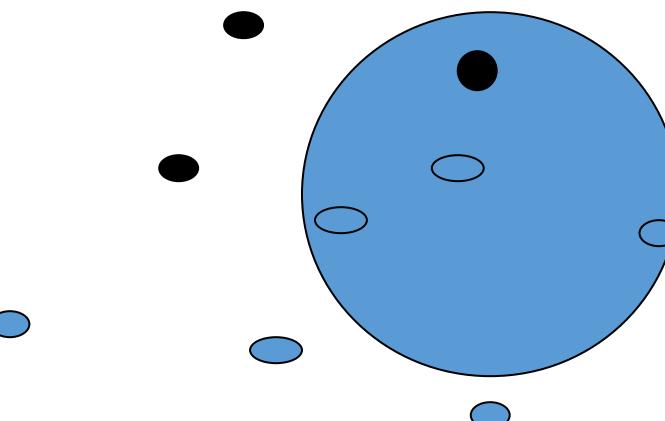
k-Nearest Neighbors

- Instance-based Learning (**Transduction**)

1-Nearest Neighbor



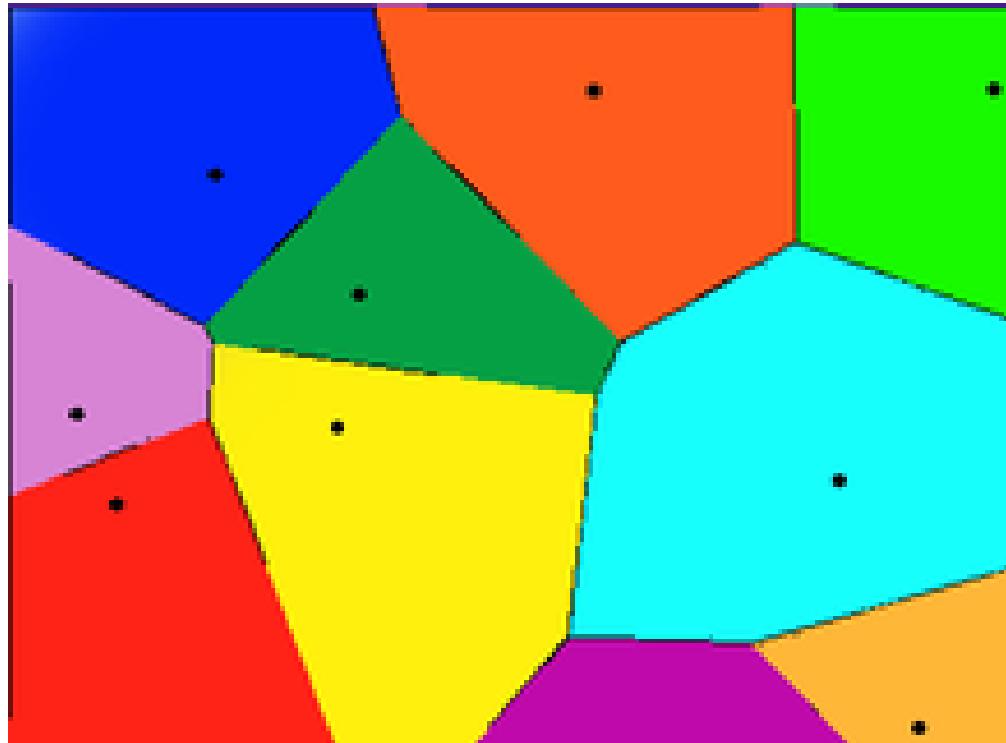
3-Nearest Neighbor



k-Nearest Neighbors

- An arbitrary instance is represented by $(a_1(x), a_2(x), a_3(x), \dots, a_n(x))$
 - $a_i(x)$ denotes features
- Euclidean distance between two instances
 - $d(x_i, x_j) = \sqrt{\sum_{r=1}^n (a_r(x_i) - a_r(x_j))^2}$
- Continuous valued target function
 - Mean value of the k nearest training examples

Voronoi Diagram



Decision surface formed by the training examples!

How to determine k

- Determined experimentally
- Start with $k=1$ and use a test set to validate the error rate of the classifier
- Repeat with $k=k+2$
- Choose the value of k for which the error rate is minimum
- Note: k should be odd number to avoid ties

How to determine k

- Determined experimentally
- Start with $k=1$ and use a test set to validate the error rate of the classifier
- Repeat with $k=k+2$
- Choose the value of k for which the error rate is minimum
- Note: k should be odd number to avoid ties

Improvements

- Distance-Weighted Nearest Neighbors
 - Assign weights to the neighbors based on their ‘distance’ from the query point (E.g., weight ‘may’ be inverse square of the distances)
- Measure “closeness” differently
- Finding “close” examples in a large training set quickly

k-NN: Summary

- **Pros**

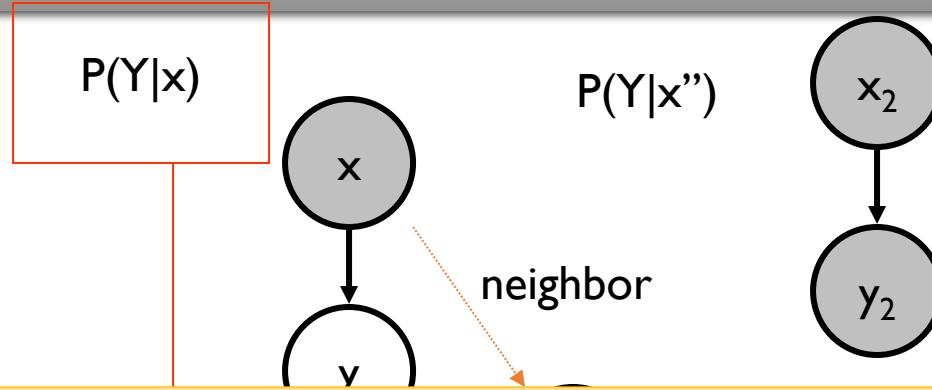
- Highly effective inductive inference method for noisy training data and complex target functions
- Target function for a whole space may be described as a combination of less complex local approximations
- Trains very fast (“Lazy” learner)

- **Cons**

- Curse of dimensionality
 - In higher dimensions, all the data points lie on the surface of the unit hypersphere (the inside is empty!) Check: <http://www.cs.cmu.edu/~venkatg/teaching/CStheory-infoage/chap1-high-dim-space.pdf>
- Storage: all training examples are saved in memory
 - A decision tree or linear classifier is much smaller
- Slow at query time
 - Can be overcome and presorting and indexing training samples

Convergence of 1-NN

$$\begin{aligned} P(\text{knnError}) &= 1 - \Pr(y = y_1) \\ &= 1 - \sum \Pr(Y = y' | x)^2 \end{aligned}$$



Possible to show that: as the size of training data set approaches infinity, the one nearest neighbour classifier guarantees an error rate of no worse than twice the Bayes error rate (the minimum achievable error rate given the distribution of the data). We will see this later.

$\text{let } y = \operatorname{argmax}_y P(y|x)$

Classification Methods

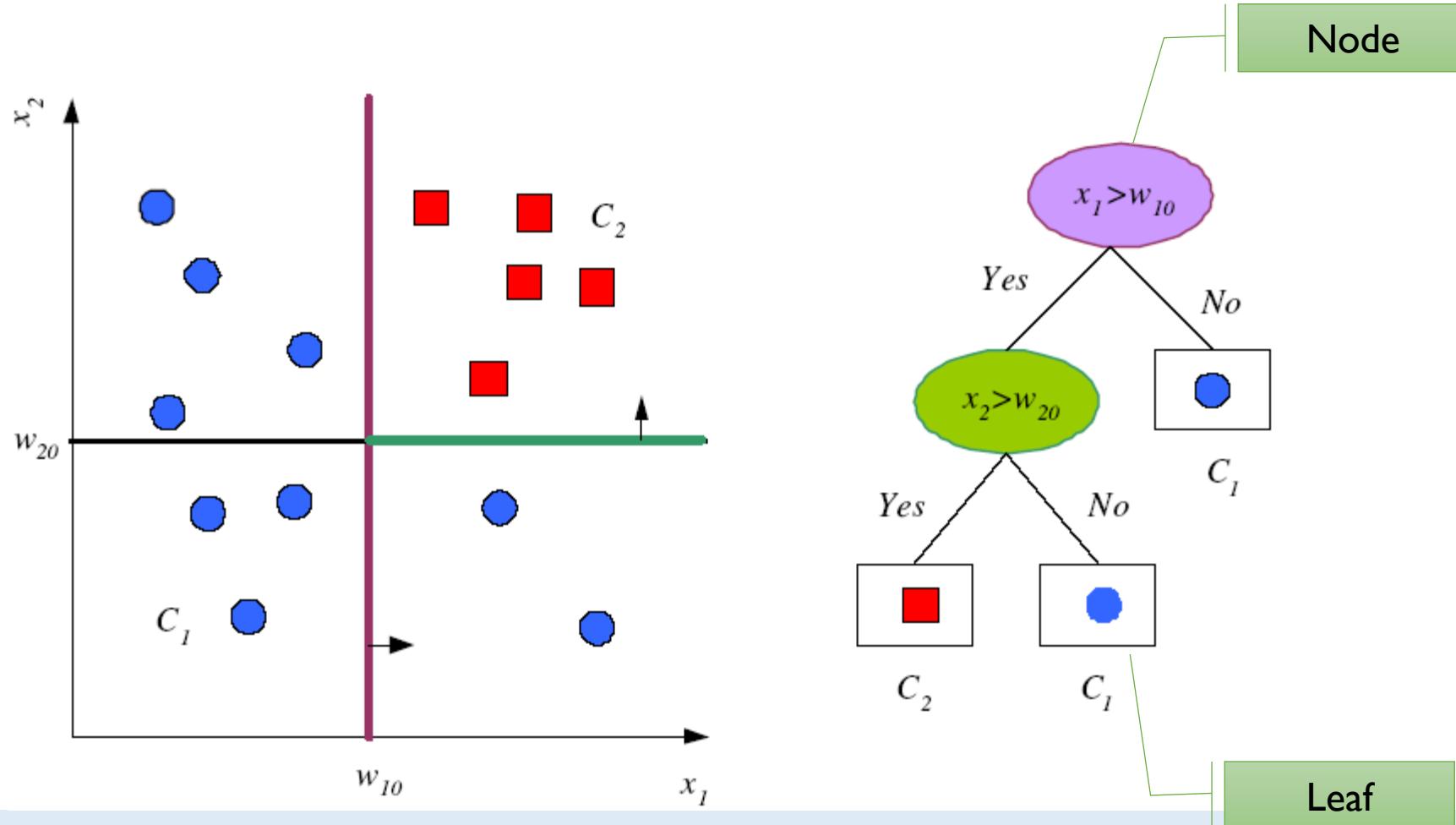
- k-Nearest Neighbors
- Decision Trees
- Naïve Bayes
- Support Vector Machines
- Logistic Regression
- Neural Networks
- Ensemble Methods (Boosting, Random Forests)

Example

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Decision Trees



Decision Trees

- An efficient nonparametric method
- A hierarchical model
- Divided-and-conquer strategy
- Supervised learning

Divide and Conquer

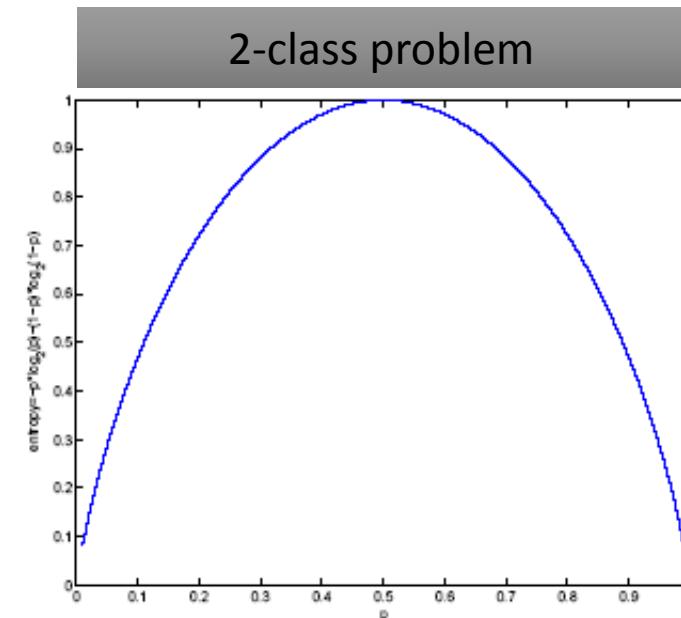
- Internal decision nodes
 - **Univariate:** Uses a single attribute, x_i
 - Numeric x_i :
 - Binary split : $x_i > w_m$
 - Discrete x_i :
 - n -way split for n possible values
 - **Multivariate:** Uses more than one attributes, \mathbf{x}
- Leaves
 - Classification: Class labels, or proportions
 - Regression: Numeric; r average, or local fit
- Learning is **greedy**; find the best split recursively

Classification Trees (C4.5, J48)

- For node m , N_m instances reach m , N_m^i belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m) \equiv p_m^i = \frac{N_m^i}{N_m}$$

- Node m is **pure** if p_m^i is 0 or 1
- Measure of **impurity** is **entropy**



Entropy in information theory specifies the **minimum number of bits** needed to encode the classification accuracy of an instance

Classification Trees

- If node m is pure, generate a leaf and stop, otherwise split and continue recursively
- Impurity after split: N_{mj} of N_m take branch j . N_{mj}^i belong to C_i

$$\hat{P}(C_i | \mathbf{x}, m, j) \equiv p_{mj}^i = \frac{N_{mj}^i}{N_{mj}}$$

$$I'_m = -\sum_{j=1}^n \frac{N_{mj}}{N_m} \sum_{i=1}^K p_{mj}^i \log_2 p_{mj}^i$$

- Choose the *best* attribute(s) to split the remaining instances and make that attribute a decision node

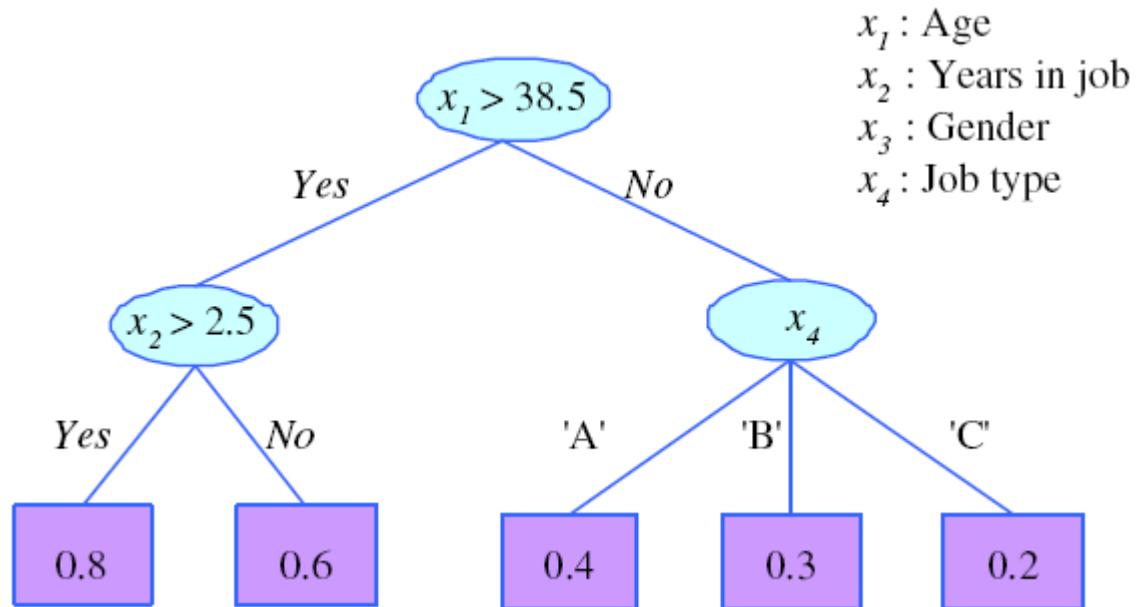
Other Measures of Impurity

- The properties of functions measuring the **impurity** of a split:
 - $\phi(1/2,1/2) \geq \phi(p,1-p)$, for any $p \in [0,1]$
 - $\phi(0,1) = \phi(1,0) = 0$
 - $\phi(p,1-p)$ is increasing in p on $[0, \frac{1}{2}]$
and decreasing in p on $[\frac{1}{2}, 1]$
- Examples (other than entropy)
 - Gini index : $\phi(p,1-p) = 2p(1-p)$
 - Misclassification error : $\phi(p,1-p) = 1 - \max(p,1-p)$

Pruning

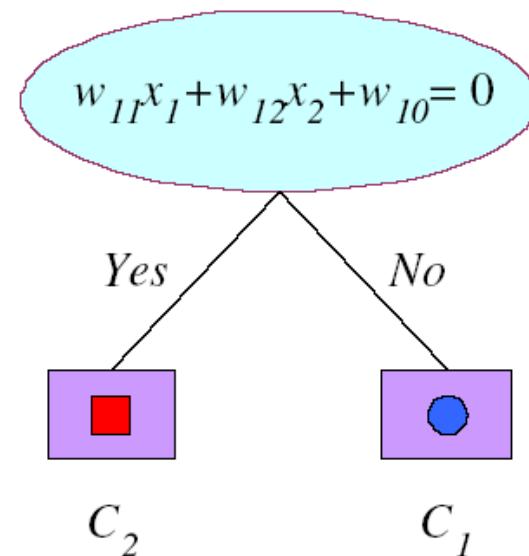
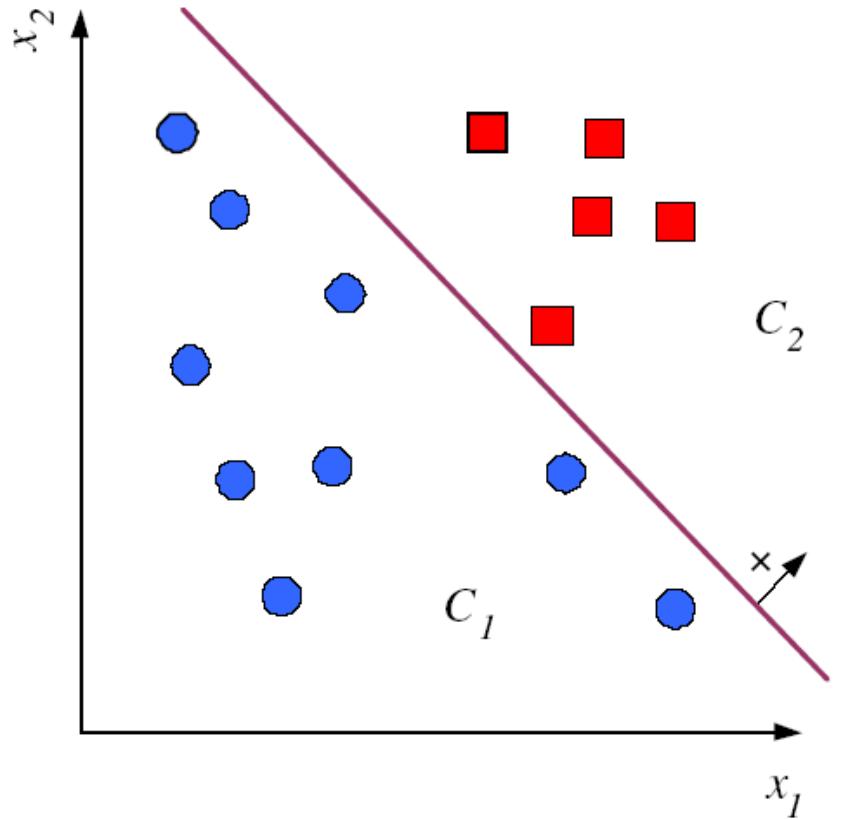
- Remove subtrees for better generalization (decrease variance)
 - Prepruning: Early stopping
 - Postpruning: Grow the whole tree then prune subtrees which overfit on the pruning set
- Prepruning is faster, postpruning is more accurate (requires a separate pruning set)

Rule Extraction from Trees



- R1: IF (age>38.5) AND (years-in-job>2.5) THEN $y = 0.8$
- R2: IF (age>38.5) AND (years-in-job≤2.5) THEN $y = 0.6$
- R3: IF (age≤38.5) AND (job-type='A') THEN $y = 0.4$
- R4: IF (age≤38.5) AND (job-type='B') THEN $y = 0.3$
- R5: IF (age≤38.5) AND (job-type='C') THEN $y = 0.2$

Multivariate Trees



CS6510 Applied Machine Learning

Classifiers

27 Aug 2016

Vineeth N Balasubramanian



Classification Methods

- k-Nearest Neighbors
- Decision Trees
- **Naïve Bayes**
- Support Vector Machines
- Logistic Regression
- Neural Networks
- Ensemble Methods (Boosting, Random Forests)

Probability: Review

Random variable

- Result of tossing a coin is from {Heads,Tails}
- Random var X from {1,0}
- Bernoulli: $P\{X=1\} = p_o^X (1 - p_o)^{1-X}$

Joint and conditional probability

$$P(A|B) = P(A, B)/P(B)$$

Bayes Theorem

$$P(A|B) = P(B|A) P(A)/P(B)$$

Illustration

A	0	0				0
B	0			0		

- $P(A=1) = 3/6 = 1/2, P(A=0) = 3/6 = 1/2.$
- $P(B=1) = 4/6 = 2/3, P(B=0) = 2/6 = 1/3.$
- $P(A=1, B = 1) = 2/6 = 1/3.$
- $P(A=1 | B = 1) = P(A=1, B = 1) / P(B=1) = 1/2.$
- $P(B=1 | A = 1) = P(B=1, A = 1) / P(A=1) = 2/3.$
- $P(A=1 | B = 1) P(B=1)/P(A=1) = 2/3 = P(B=1 | A = 1).$
 - Bayes' Theorem

Naïve Bayes Classifier

- Goal: Learning function $f(x) \rightarrow y$
 - Y : One of k classes (e.g. spam/ham, digit 0-9)
 - $X=X_1, \dots, X_n$: Values of attributes (numeric or categorical)
- Probabilistic classification
 - Most probable class given observation: $\hat{y} = \arg \max_y P(y|x)$
- Bayesian probability of a class

$$P(y|x) = \frac{P(x|y)P(y)}{\sum_{y'} P(x|y')P(y')}$$

class model prior

 $P(x|y)P(y)$
 $\sum_{y'} P(x|y')P(y')$
normalizer $P(x)$

Bayes Theorem

Bayes Theorem: Example

- Given:
 - A doctor knows that meningitis causes stiff neck 50% of the time
 - Prior probability of any patient having meningitis is 1/50,000
 - Prior probability of any patient having stiff neck is 1/20
- If a patient has stiff neck, what's the probability he/she has meningitis?

$$P(M | S) = \frac{P(S | M)P(M)}{P(S)} = \frac{0.5 \times 1/50000}{1/20} = 0.0002$$

What is Naïve about it?

- Consider each attribute and class label as random variables
- Given a record with attributes (A_1, A_2, \dots, A_n)
 - Goal is to predict class C
 - Specifically, we want to find the value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
- Can we estimate $P(C | A_1, A_2, \dots, A_n)$ directly from data?

What is Naïve about it?

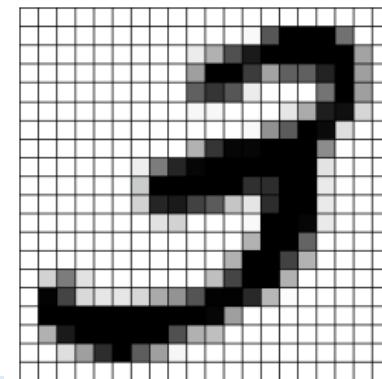
- Approach:
 - compute the posterior probability $P(C | A_1, A_2, \dots, A_n)$ for all values of C using the Bayes theorem
 - Choose value of C that maximizes $P(C | A_1, A_2, \dots, A_n)$
 - Equivalent to choosing value of C that maximizes $P(A_1, A_2, \dots, A_n | C) P(C)$
- How to estimate $P(A_1, A_2, \dots, A_n | C)$?

What is Naïve about it?

$$P(C | A_1 A_2 \dots A_n) = \frac{P(A_1 A_2 \dots A_n | C) P(C)}{P(A_1 A_2 \dots A_n)}$$

How to compute if x is made of multiple attributes?

- 20×20 image of digit = 2^{400} possible combinations!



Naïve Bayes Classifier

- Assume independence among attributes A_i when class is given:
 - $P(A_1, A_2, \dots, A_n | C_j) = P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$
 - Can estimate $P(A_i | C_j)$ for all A_i and C_j .
 - New point is classified to C_j if
 $P(C_j) \prod_i P(A_i | C_j) = P(C_j) P(A_1 | C_j) P(A_2 | C_j) \dots P(A_n | C_j)$
is maximal

Example

PlayTennis: training examples

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No

Example: Learning Phase

Outlook	Play=Yes	Play=No
Sunny	2/9	3/5
Overcast	4/9	0/5
Rain	3/9	2/5

Temperature	Play=Yes	Play=No
Hot	2/9	2/5
Mild	4/9	2/5
Cool	3/9	1/5

Humidity	Play=Yes	Play=No
High	3/9	4/5
Normal	6/9	1/5

Wind	Play=Yes	Play=No
Strong	3/9	3/5
Weak	6/9	2/5

$$P(\text{Play}=\text{Yes}) = 9/14$$

$$P(\text{Play}=\text{No}) = 5/14$$

Example: Test Phase

- Given a new instance,
- $\mathbf{x}' = (\text{Outlook}=\text{Sunny}, \text{Temperature}=\text{Cool}, \text{Humidity}=\text{High}, \text{Wind}=\text{Strong})$
- Look up tables

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{Yes}) = 2/9$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{Yes}) = 3/9$$

$$P(\text{Play}=\text{Yes}) = 9/14$$

$$P(\text{Outlook}=\text{Sunny} \mid \text{Play}=\text{No}) = 3/5$$

$$P(\text{Temperature}=\text{Cool} \mid \text{Play}=\text{No}) = 1/5$$

$$P(\text{Humidity}=\text{High} \mid \text{Play}=\text{No}) = 4/5$$

$$P(\text{Wind}=\text{Strong} \mid \text{Play}=\text{No}) = 3/5$$

$$P(\text{Play}=\text{No}) = 5/14$$

- MAP rule

$$P(\text{Yes} \mid \mathbf{x}') = [P(\text{Sunny} \mid \text{Yes})P(\text{Cool} \mid \text{Yes})P(\text{High} \mid \text{Yes})P(\text{Strong} \mid \text{Yes})]P(\text{Play}=\text{Yes}) = 0.0053$$

$$P(\text{No} \mid \mathbf{x}') = [P(\text{Sunny} \mid \text{No})P(\text{Cool} \mid \text{No})P(\text{High} \mid \text{No})P(\text{Strong} \mid \text{No})]P(\text{Play}=\text{No}) = 0.0206$$

Given the fact $P(\text{Yes} \mid \mathbf{x}') < P(\text{No} \mid \mathbf{x}')$, we label \mathbf{x}' to be "No".

Example: Another

Name	Give Birth	Can Fly	Live in Water	Have Legs	Class
human	yes	no	no	yes	mammals
python	no	no	no	no	non-mammals
salmon	no	no	yes	no	non-mammals
whale	yes	no	yes	no	mammals
frog	no	no	sometimes	yes	non-mammals
komodo	no	no	no	yes	non-mammals
bat	yes	yes	no	yes	mammals
pigeon	no	yes	no	yes	non-mammals
cat	yes	no	no	yes	mammals
leopard shark	yes	no	yes	no	non-mammals
turtle	no	no	sometimes	yes	non-mammals
penguin	no	no	sometimes	yes	non-mammals
porcupine	yes	no	no	yes	mammals
eel	no	no	yes	no	non-mammals
salamander	no	no	sometimes	yes	non-mammals
gila monster	no	no	no	yes	non-mammals
platypus	no	no	no	yes	mammals
owl	no	yes	no	yes	non-mammals
dolphin	yes	no	yes	no	mammals
eagle	no	yes	no	yes	non-mammals

Give Birth	Can Fly	Live in Water	Have Legs	Class
yes	no	yes	no	?

A: attributes

M: mammals

N: non-mammals

$$P(A|M) = \frac{6}{7} \times \frac{6}{7} \times \frac{2}{7} \times \frac{2}{7} = 0.06$$

$$P(A|N) = \frac{1}{13} \times \frac{10}{13} \times \frac{3}{13} \times \frac{4}{13} = 0.0042$$

$$P(A|M)P(M) = 0.06 \times \frac{7}{20} = 0.021$$

$$P(A|N)P(N) = 0.004 \times \frac{13}{20} = 0.0027$$

$P(A|M)P(M) > P(A|N)P(N)$

=> Mammals

Naïve Bayes Classifier: Practical Issues

- Robust to isolated noise points
- Handle missing values by ignoring the instance during probability estimate calculations
- Robust to irrelevant attributes
- Independence assumption may not hold always

Naïve Bayes Classifier: Practical Issues

- For continuous attributes:
 - **Discretize** the range into bins
 - one ordinal attribute per bin
 - violates independence assumption
 - **Two-way split:** $(A < v)$ or $(A > v)$
 - choose only one of the two splits as new attribute
 - **Probability density estimation:**
 - Assume attribute follows a parametrized distribution, e.g. normal distribution
 - Use data to estimate parameters of distribution (e.g., mean and standard deviation) using Maximum Likelihood Estimation
 - Once probability distribution is known, can use it to estimate the conditional probability $P(A_i|c)$

Naïve Bayes Classifier: Example

- Problem Definition
 - Given a set of news articles that are of interest, we would like to learn to classify the articles by topic.
- Naïve Bayes is among the most effective algorithms to perform this task.
- What will be attributes to represent the documents?
 - Vector of words – one attribute per word position in the document
- What is the Target concept
 - Is the document interesting?
 - Topic of the document

Example: 20 Newsgroup

- Given 1000 training documents from each group
- Learn to classify new documents according to the newsgroup it came from
- NBC – 89% accuracy

comp.graphics	misc.forsale
comp.os.ms-windows.misc	rec.autos
comp.sys.ibm.pc.hardware	rec.motorcycles
comp.sys.mac.hardware	rec.sport.baseball
comp.windows.x	rec.sport.hockey
alt.atheism	sci.space
soc.religion.christian	sci.crypt
talk.religion.misc	sci.electronics
talk.politics.mideast	sci.med
talk.politics.misc	talk.politics.guns

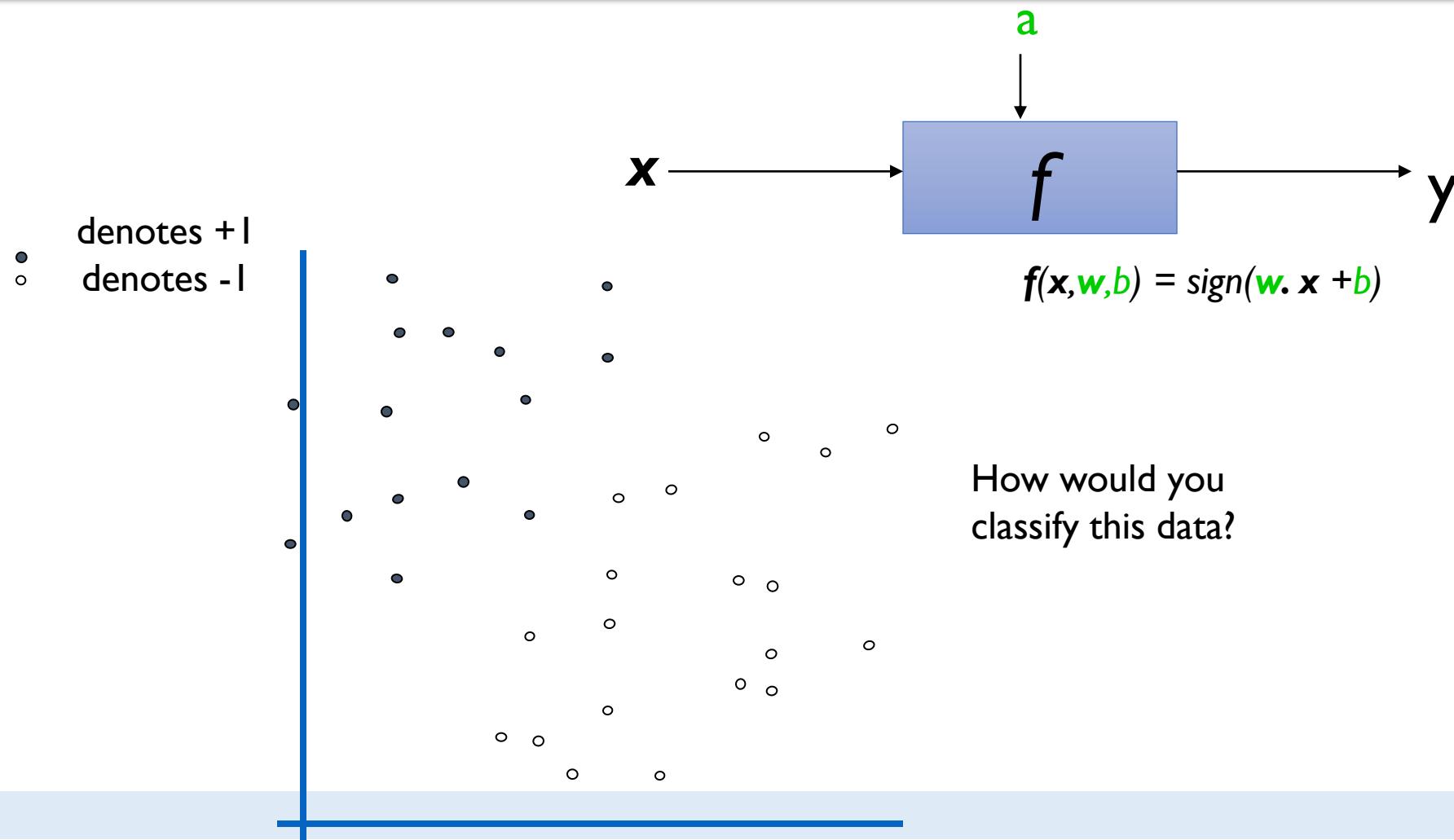
Classification Methods

- k-Nearest Neighbors
- Decision Trees
- Naïve Bayes
- **Support Vector Machines**
- Logistic Regression
- Neural Networks
- Ensemble Methods (Boosting, Random Forests)

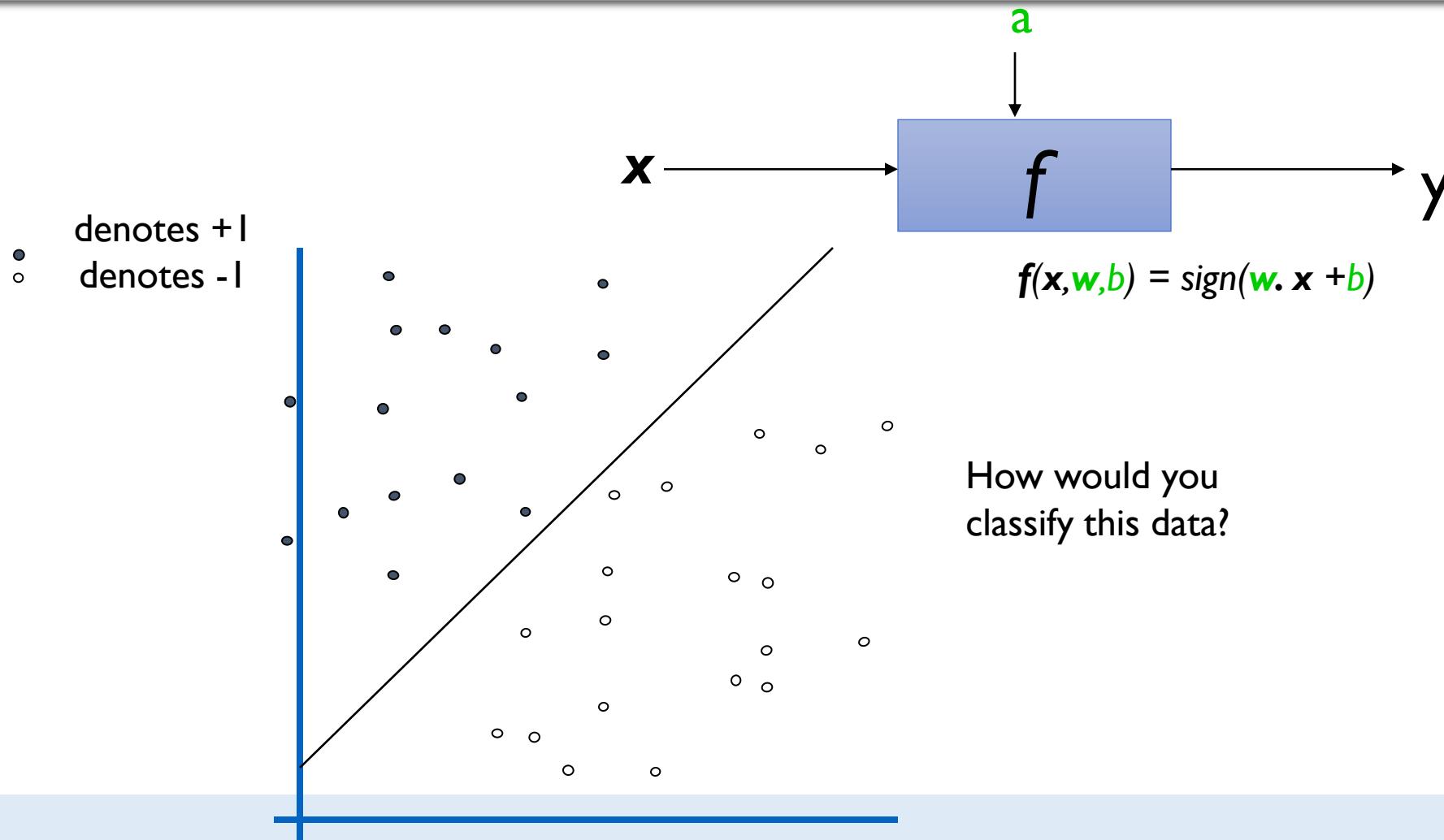
SVM: History

- SVM is inspired from statistical learning theory
- SVM was first introduced in 1992
- SVM becomes popular because of its success in handwritten digit recognition
- SVM is now regarded as an important example of “kernel methods”, arguably the hottest area in machine learning. <http://www.kernel-machines.org/>

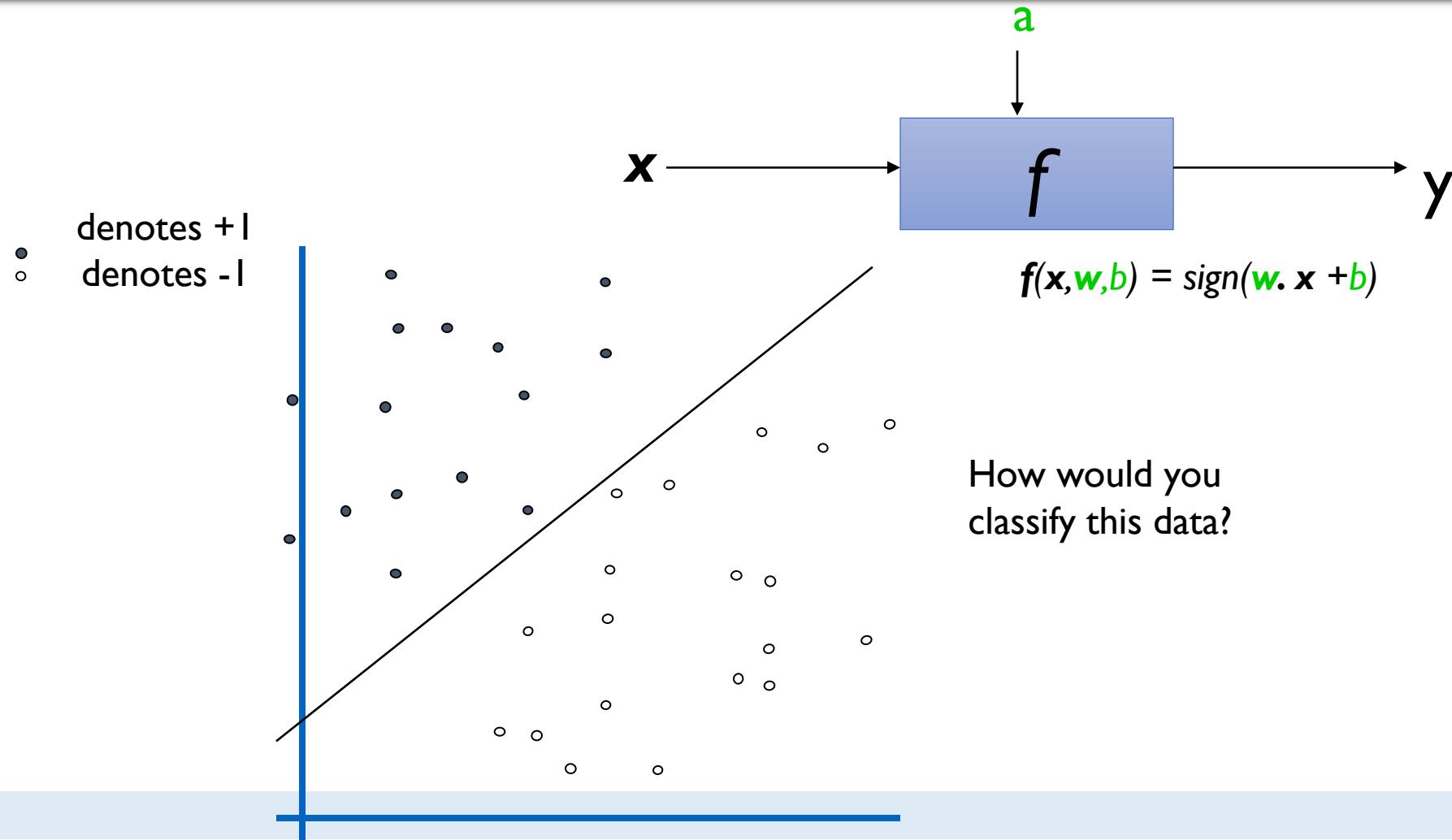
Linear Classifiers



Linear Classifiers

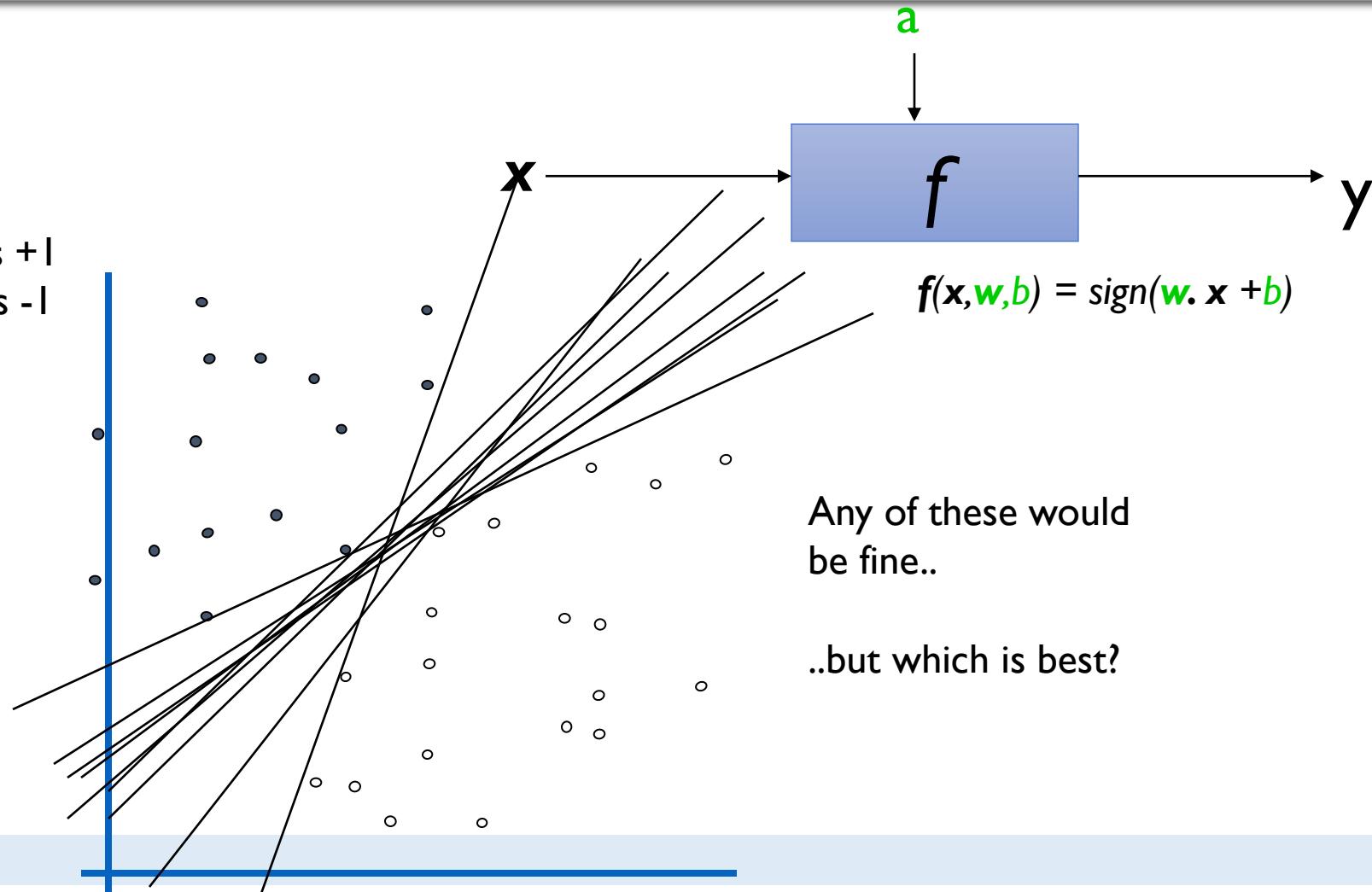


Linear Classifiers



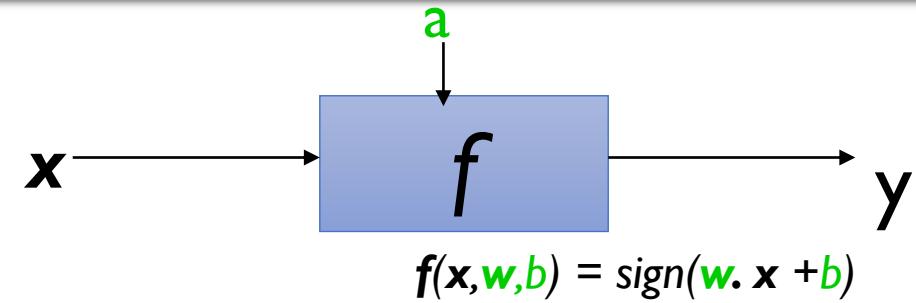
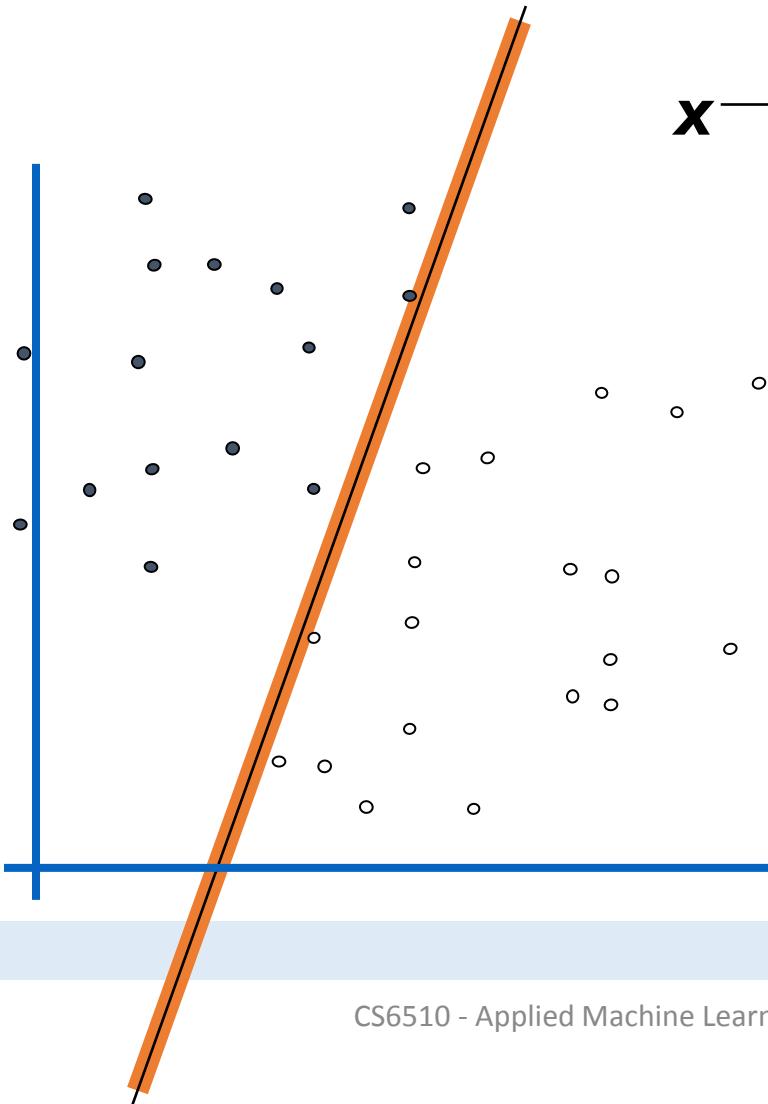
Linear Classifiers

- denotes +1
- denotes -1



Linear Classifiers

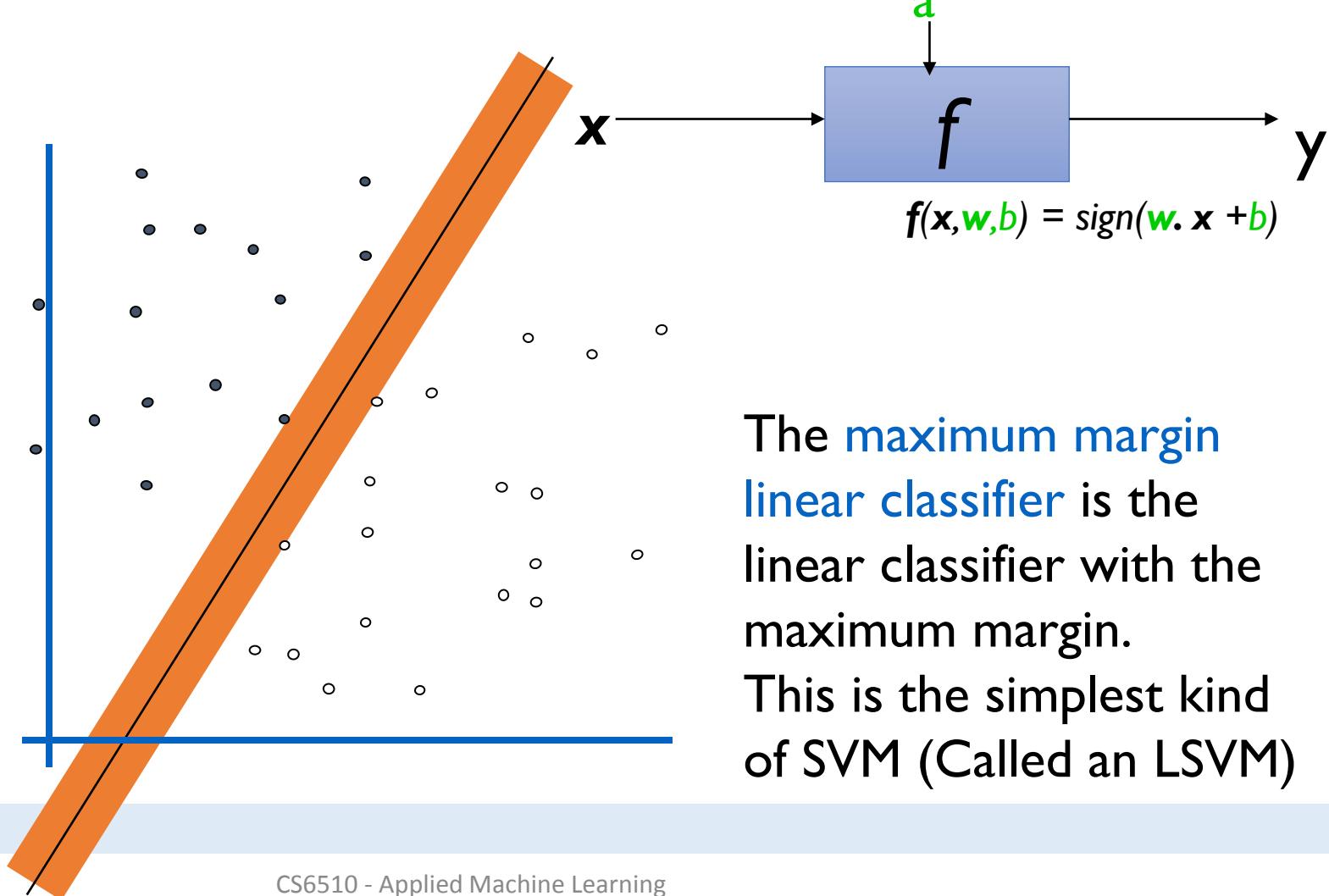
- denotes +1
- denotes -1



Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

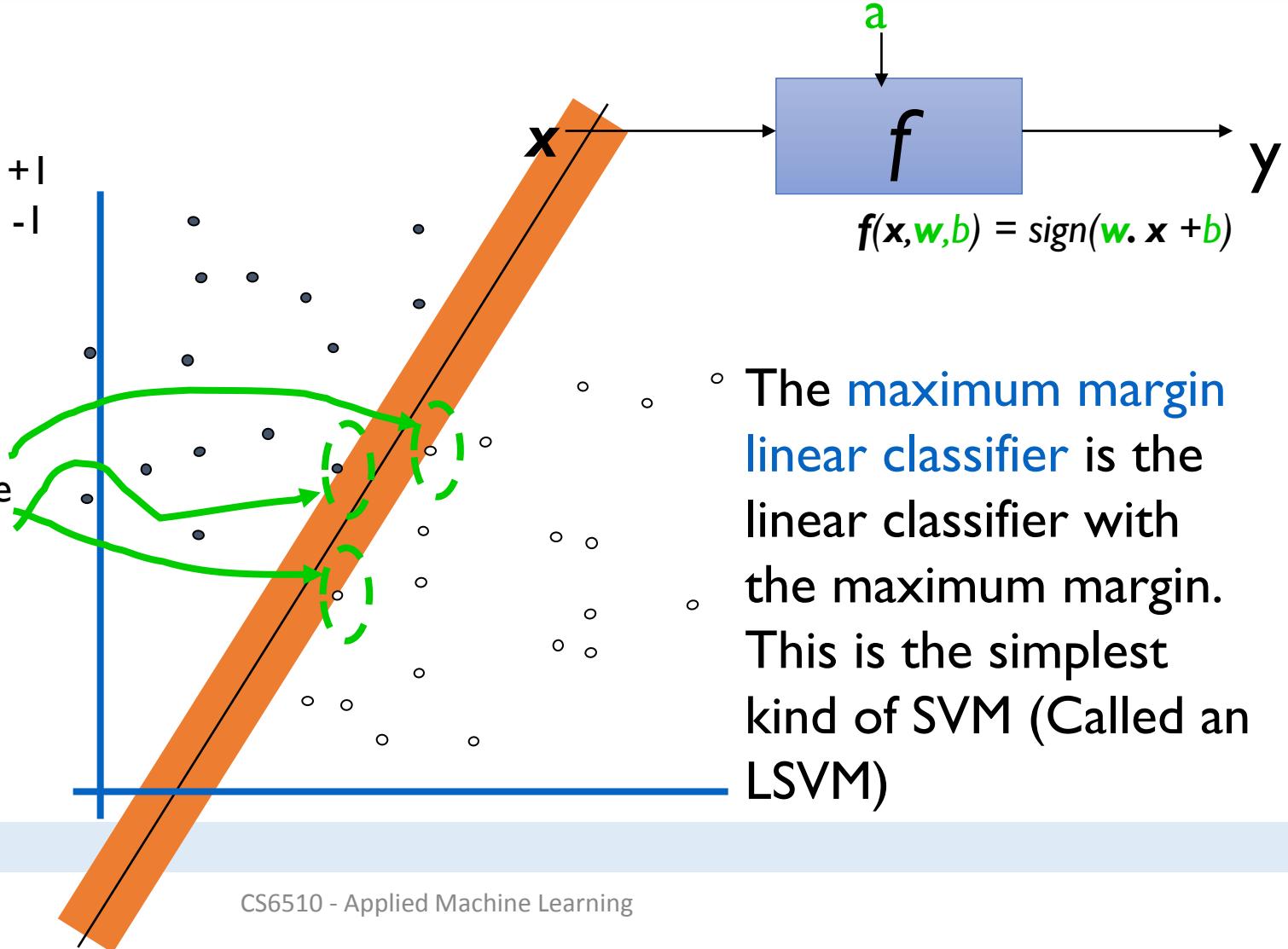
Linear Classifiers

- denotes +1
- denotes -1



Maximum Margin Classifier

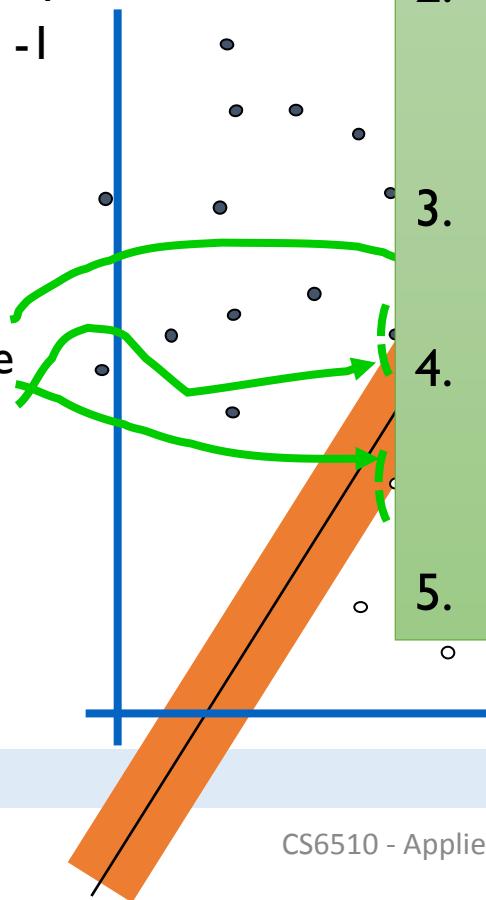
- denotes +1
 - denotes -1
- Support Vectors are those data points that the margin pushes up against



Why Maximum Margin?

- denotes +1
- denotes -1

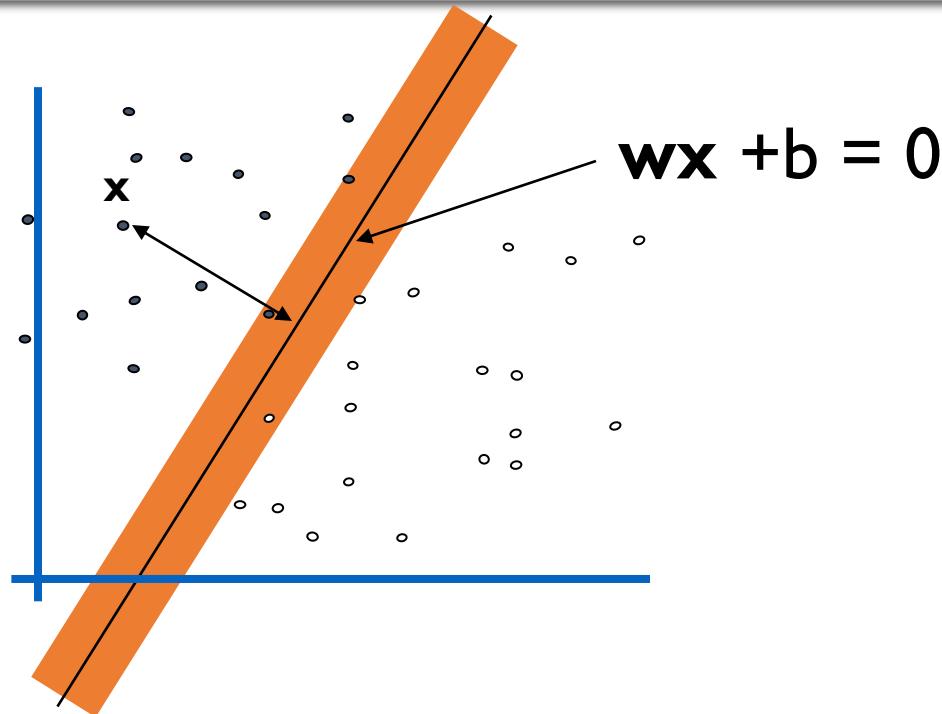
Support Vectors are those data points that the margin pushes up against



1. Intuitively this feels safest.
2. If we've made a small error in the location of the boundary this gives us least chance of causing a misclassification.
3. The model is immune to removal of any non-support-vector datapoints.
4. There's some theory (using VC dimension) that is related to (but not the same as) the proposition that this is a good thing.
5. Empirically it works very very well.

Estimating the Margin

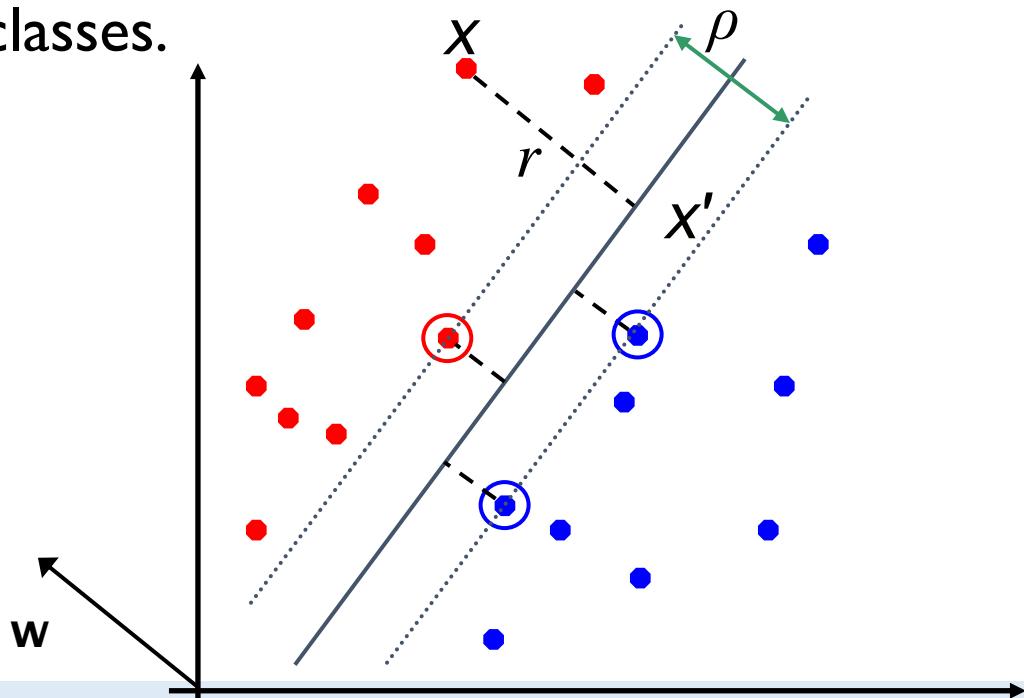
- denotes +1
- denotes -1



- What is the distance expression for a point x to a line $wx+b=0$?

Estimating the Margin

- Distance from example to the separator is $r = y \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|}$
- Examples closest to the hyperplane are **support vectors**.
- Margin ρ** of the separator is the width of separation between support vectors of classes.



Derivation of finding r :

Dotted line $\mathbf{x}' - \mathbf{x}$ is perpendicular to decision boundary so parallel to \mathbf{w} .
Unit vector is $\mathbf{w}/|\mathbf{w}|$, so line is $r\mathbf{w}/|\mathbf{w}|$.

$$\mathbf{x}' = \mathbf{x} - yr\mathbf{w}/|\mathbf{w}|.$$

$$\mathbf{x}' \text{ satisfies } \mathbf{w}^T \mathbf{x}' + b = 0.$$

$$\text{So } \mathbf{w}^T(\mathbf{x} - yr\mathbf{w}/|\mathbf{w}|) + b = 0$$

$$\text{Recall that } |\mathbf{w}| = \sqrt{\mathbf{w}^T \mathbf{w}}.$$

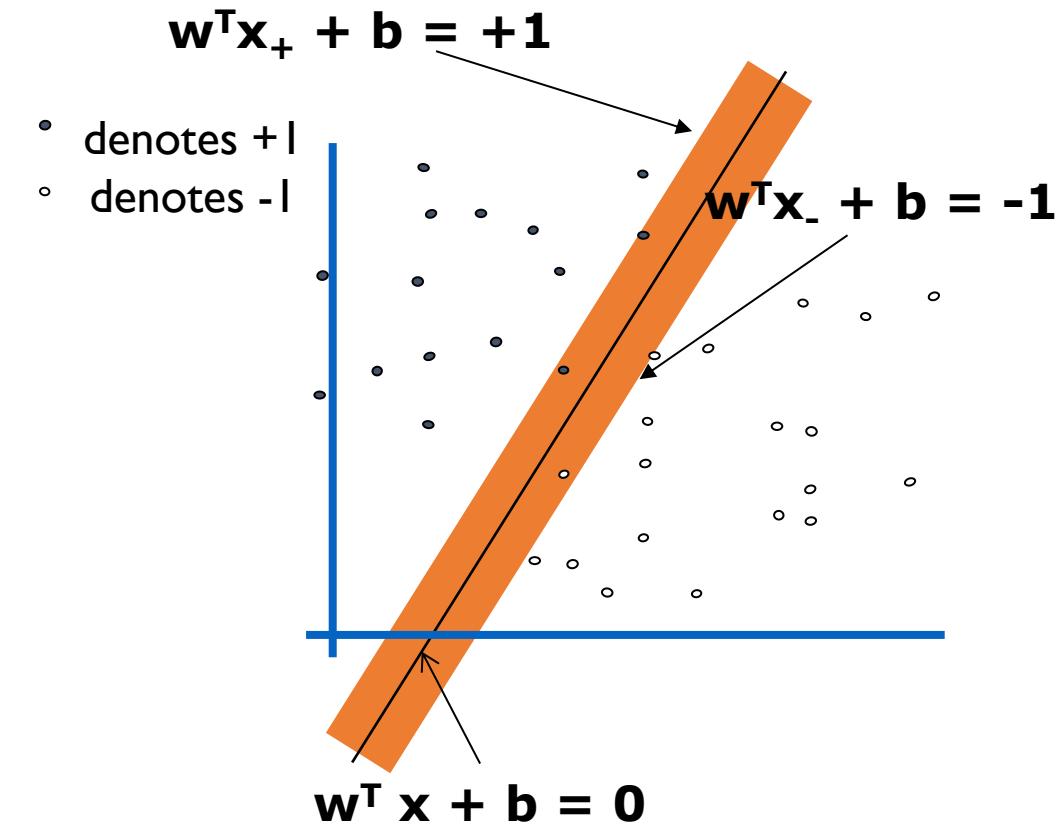
$$\text{So } \mathbf{w}^T \mathbf{x} - yr|\mathbf{w}| + b = 0$$

So, solving for r gives:

$$r = y(\mathbf{w}^T \mathbf{x} + b)/|\mathbf{w}|$$

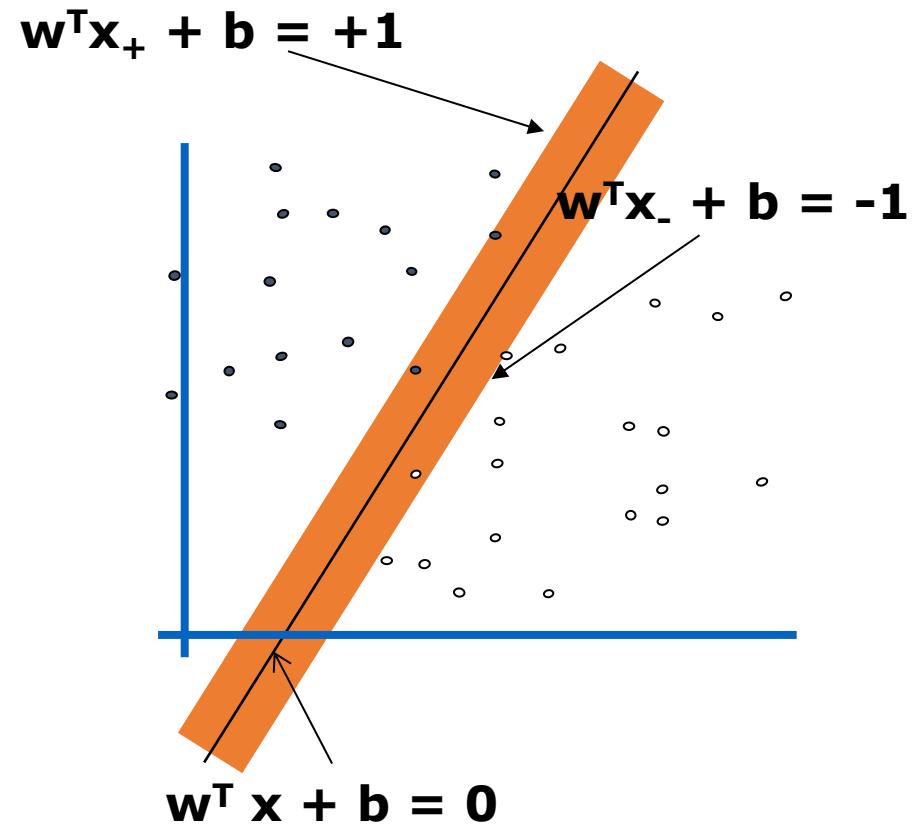
Estimating the Margin

- Since $\mathbf{w}^T \mathbf{x} + b = 0$ and $c(\mathbf{w}^T \mathbf{x} + b) = 0$ define the same plane, we have the freedom to choose the normalization of \mathbf{w} (i.e. c)
- Let us choose normalization such that $\mathbf{w}^T \mathbf{x}_+ + b = +1$ and $\mathbf{w}^T \mathbf{x}_- + b = -1$ for the positive and negative support vectors respectively

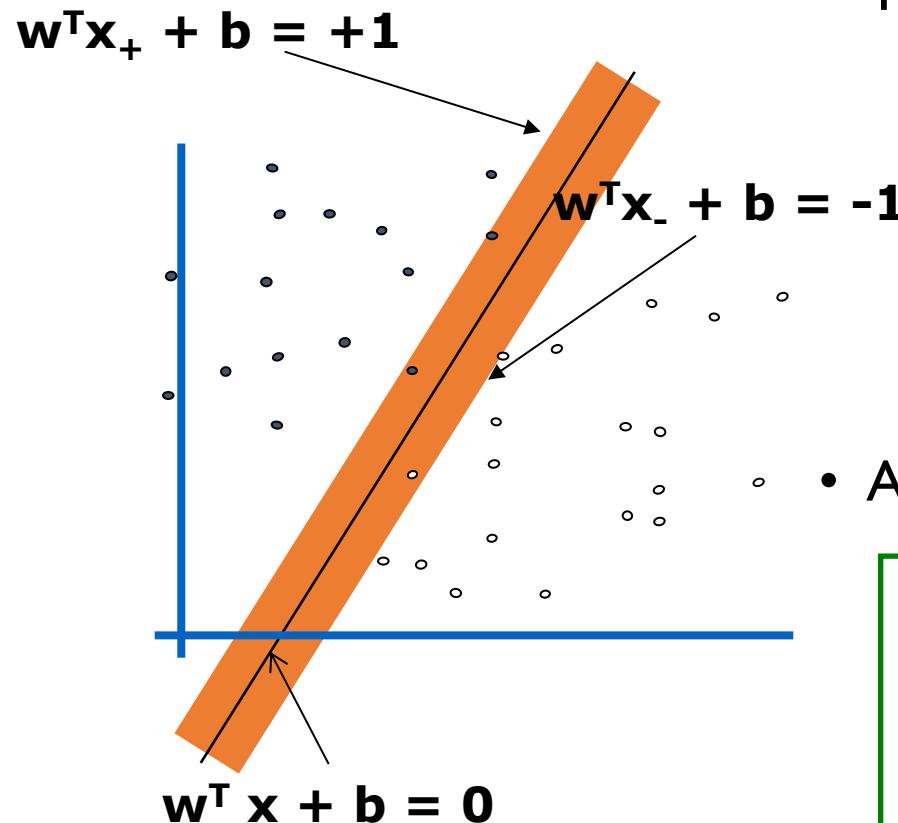


Estimating the Margin

- Since $\mathbf{w}^T \mathbf{x} + b = 0$ and $c(\mathbf{w}^T \mathbf{x} + b) = 0$ define the same plane, we have the freedom to choose the normalization of \mathbf{w} (i.e. c)
- Let us choose normalization such that $\mathbf{w}^T \mathbf{x}_+ + b = +1$ and $\mathbf{w}^T \mathbf{x}_- + b = -1$ for the positive and negative support vectors respectively
- Hence, margin now is: $2 * \frac{\mathbf{w}^T \mathbf{x} + b}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$



Maximizing the Margin



- Then we can formulate the *quadratic optimization problem*:

Find w and b such that

$$r = \frac{2}{\|w\|} \text{ is maximized; and for all } \{(x_i, y_i)\}$$

$$w^T x_i + b \geq 1 \text{ if } y_i = 1; \quad w^T x_i + b \leq -1 \quad \text{if } y_i = -1$$

- A better formulation ($\min \|w\| = \max 1/\|w\|$):

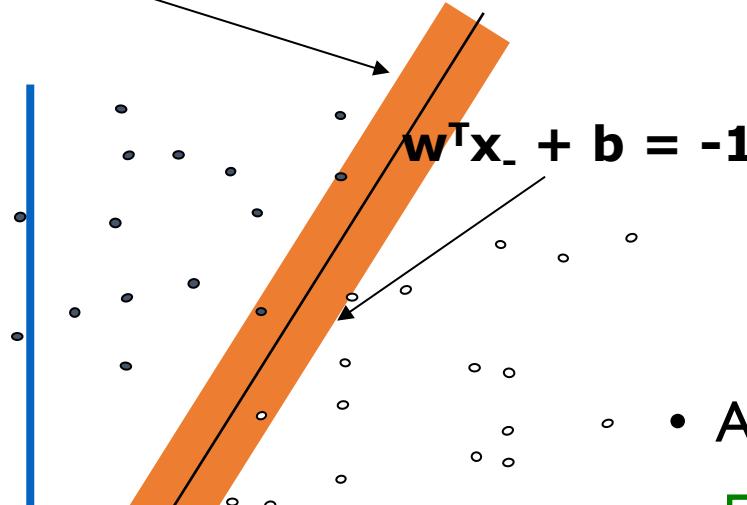
Find w and b such that

$$\Phi(w) = \frac{1}{2} w^T w \text{ is minimized;}$$

$$\text{and for all } \{(x_i, y_i)\}: \quad y_i (w^T x_i + b) \geq 1$$

Maximizing the Margin

$$\mathbf{w}^T \mathbf{x}_+ + b = +1$$



How to solve?

Quadratic Programming

- Then we can formulate the *quadratic optimization problem*:

Find \mathbf{w} and b such that

$$r = \frac{2}{\|\mathbf{w}\|} \text{ is maximized; and for all } \{(\mathbf{x}_i, y_i)\}$$
$$\mathbf{w}^T \mathbf{x}_i + b \geq 1 \text{ if } y_i = 1; \quad \mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{if } y_i = -1$$

- A better formulation ($\min \|\mathbf{w}\| = \max 1/\|\mathbf{w}\|$):

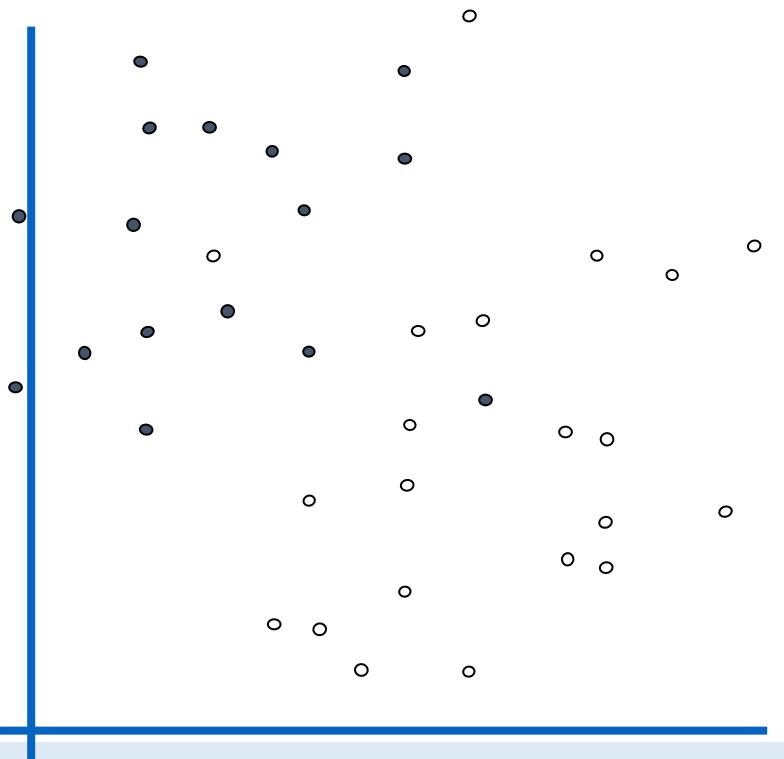
Find \mathbf{w} and b such that

$$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \text{ is minimized;}$$

and for all $\{(\mathbf{x}_i, y_i)\}$: $y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$

Non-separable Data

- denotes +1
- denotes -1



This is going to be a problem!
What should we do?

SVM for Noisy Data

$$\{\vec{w}^*, b^*\} = \min_{\vec{w}, b} \sum_{i=1}^d w_i^2 + c \sum_{j=1}^N \varepsilon_j$$

$$y_1 (\vec{w} \cdot \vec{x}_1 + b) \geq 1 - \varepsilon_1, \varepsilon_1 \geq 0$$

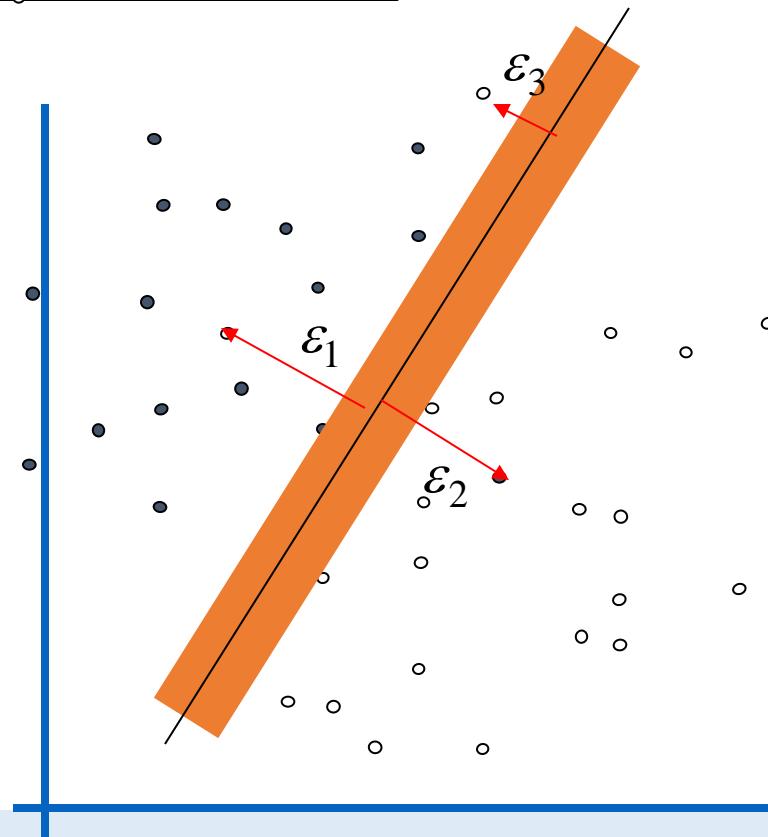
$$y_2 (\vec{w} \cdot \vec{x}_2 + b) \geq 1 - \varepsilon_2, \varepsilon_2 \geq 0$$

...

$$y_N (\vec{w} \cdot \vec{x}_N + b) \geq 1 - \varepsilon_N, \varepsilon_N \geq 0$$

- Balance the trade off between margin and classification errors

• denotes +1
○ denotes -1

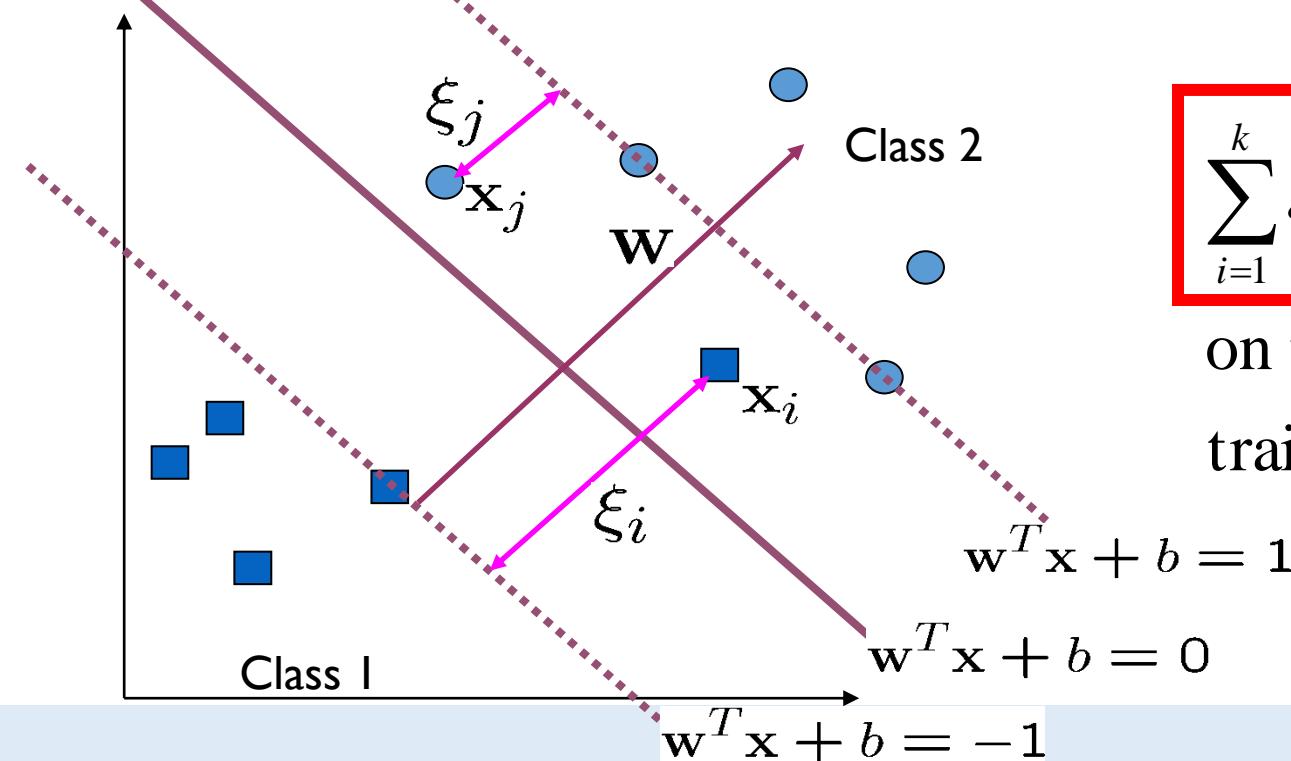


SVM for Noisy Data

$\varepsilon_i \geq 1 \Leftrightarrow y_i(\mathbf{w}^T \mathbf{x}_i + b) < 0$, i.e., misclassification

$0 < \varepsilon_i < 1 \Leftrightarrow x_i$ is correctly classified, but lies inside the margin

$\varepsilon_i = 0 \Leftrightarrow x_i$ is classified correctly, and lies outside the margin



$\sum_{i=1}^k \varepsilon_i$ is an upperbound

on the number of
training errors.

SVM for Noisy Data

Use the Lagrangian formulation for the optimization problem.

Introduce a positive Lagrangian multiplier for each inequality constraint.

$$y_i(x_i \bullet w + b) - 1 + \varepsilon_i \geq 0, \text{ for all } i.$$

$\varepsilon_i \geq 0, \text{ for all } i.$

α_i β_i

Lagrangian multipliers

The diagram illustrates the decomposition of the SVM margin equation. It shows the original equation $y_i(x_i \bullet w + b) - 1 + \varepsilon_i \geq 0$ at the top. Below it, the term $y_i(x_i \bullet w + b) - 1$ is highlighted with a blue arrow pointing to the label "Lagrangian multipliers". To the right, another blue arrow points from the term $-\sum_i \alpha_i y_i x_i$ to the label β_i . Further down, a brown arrow points from the term $+\sum_i \varepsilon_i$ to the label " $\varepsilon_i \geq 0, \text{ for all } i.$ ". Finally, a brown arrow points from the term -1 to the label "for all i."

Get the following Lagrangian:
$$L_p = \|w\|^2 + c \sum_i \varepsilon_i - \sum_i \alpha_i \{y_i(x_i \bullet w + b) - 1 + \varepsilon_i\} - \sum_i \beta_i \varepsilon_i$$

SVM for Noisy Data

$$L_p = \|w\|^2 + c \sum_i \varepsilon_i - \sum_i \alpha_i \{y_i(x_i \bullet w + b) - 1 + \varepsilon_i\} - \sum_i \beta_i \varepsilon_i$$

$$\frac{\partial L_p}{\partial w} = 2w - \sum_i \alpha_i y_i x_i = 0 \Rightarrow w = \frac{1}{2} \sum_i \alpha_i y_i x_i$$

$$\frac{\partial L_p}{\partial b} = -\frac{1}{2} \sum_i \alpha_i y_i = 0 \Rightarrow \sum_i \alpha_i y_i = 0$$

$$\frac{\partial L_p}{\partial \varepsilon_i} = c - \beta_i - \alpha_i = 0 \Rightarrow c = \beta_i + \alpha_i$$

Take the derivatives of L_p with respect to w , b , and ε_i .

Karush-Kuhn-Tucker Conditions

$$0 \leq \alpha_i \leq c \quad \forall i$$

$$L_D = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j (x_i \bullet x_j)$$

Both ε_i and its multiplier β_i are not involved in the function.

Dual Form

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\mathbf{x}_k \cdot \mathbf{x}_l)$

Subject to these constraints:

$$0 \leq \alpha_k \leq c \quad \forall k \quad \sum_{k=1}^R \alpha_k y_k = 0$$

Then define:

$$\mathbf{w} = \frac{1}{2} \sum_{k=1}^R \alpha_k y_k \mathbf{x}_k$$

$$y_i(\mathbf{x}_i \bullet \mathbf{w} + b) - 1 = 0$$

$$b = -y_i(\mathbf{x}_i \bullet \mathbf{w}) - 1$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, b) = sign(\mathbf{w} \cdot \mathbf{x} + b)$$

Dual Form

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\mathbf{x}_k \cdot \mathbf{x}_l)$

Subject to these constraints:

$$0 \leq \alpha_k \leq c \quad \forall k$$

$$\sum_{k=1}^R \alpha_k y_k = 0$$

Then define:

$$\mathbf{w} = \frac{1}{2} \sum_{k=1}^R \alpha_k y_k \mathbf{x}_k$$

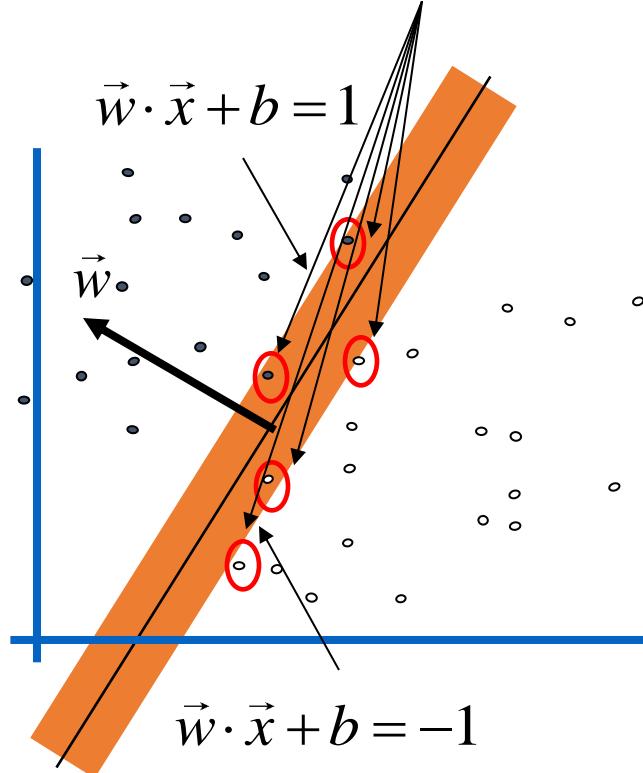
$$\begin{aligned} y_i(\mathbf{x}_i \bullet \mathbf{w} + b) - 1 &= 0 \\ b &= -y_i(y_i(\mathbf{x}_i \bullet \mathbf{w}) - 1) \end{aligned}$$

Datapoints with $\alpha_k > 0$ will be the support vectors

Then ..so this sum only needs to be over the support vectors.
 $f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{x} + b)$

Support Vectors

- denotes +1
- denotes -1



Support Vectors

$$\forall i: \alpha_i (y_i (\vec{w} \cdot \vec{x}_i + b) - (1 - \varepsilon_i)) = 0$$

$\alpha_i = 0$ for non-support vectors
 $\alpha_i \neq 0$ for support vectors

$$\mathbf{w} = \frac{1}{2} \sum_{k=1}^R \alpha_k y_k \mathbf{x}_k$$

Decision boundary is determined
only by those support vectors !

CS6510 Applied Machine Learning

Classifiers

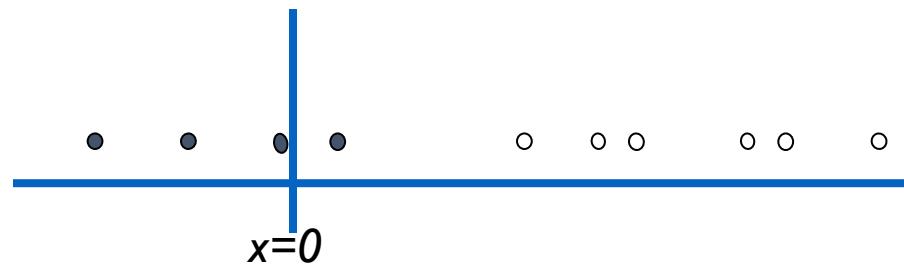
3 Sep 2016

Vineeth N Balasubramanian



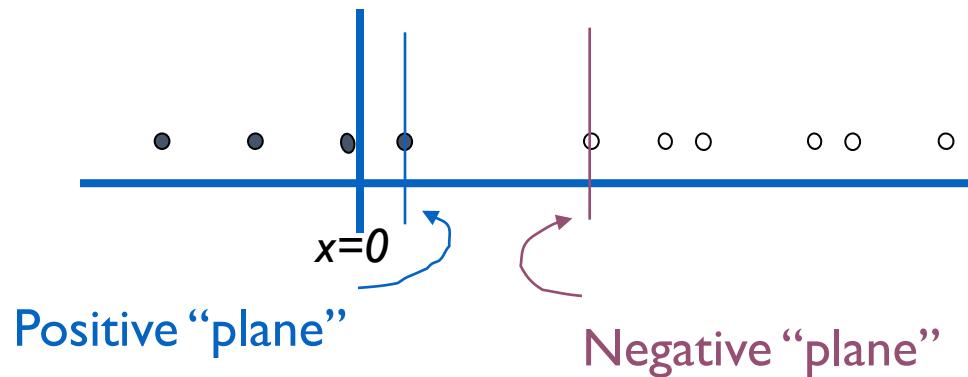
Assume we are in 1-dimension

What would SVMs
do with this
data?



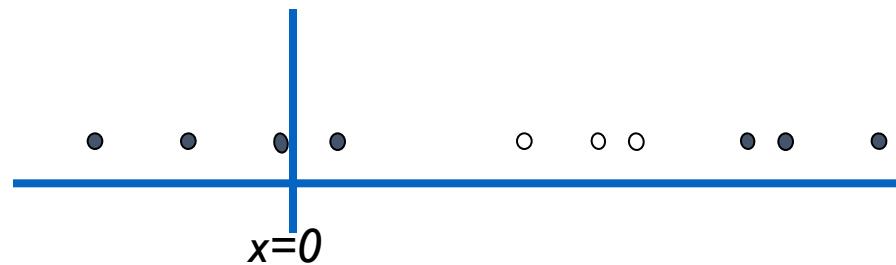
Assume we are in 1-dimension

Not a big surprise

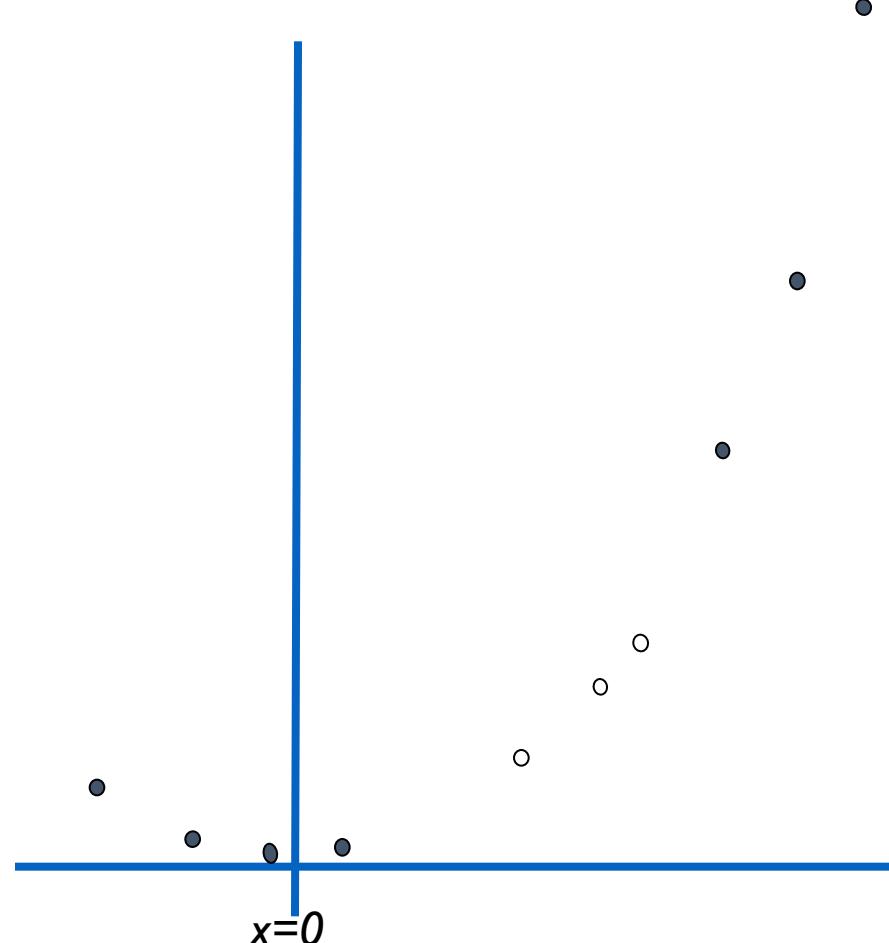


Harder 1-dimensional Dataset

What can be done
about this?



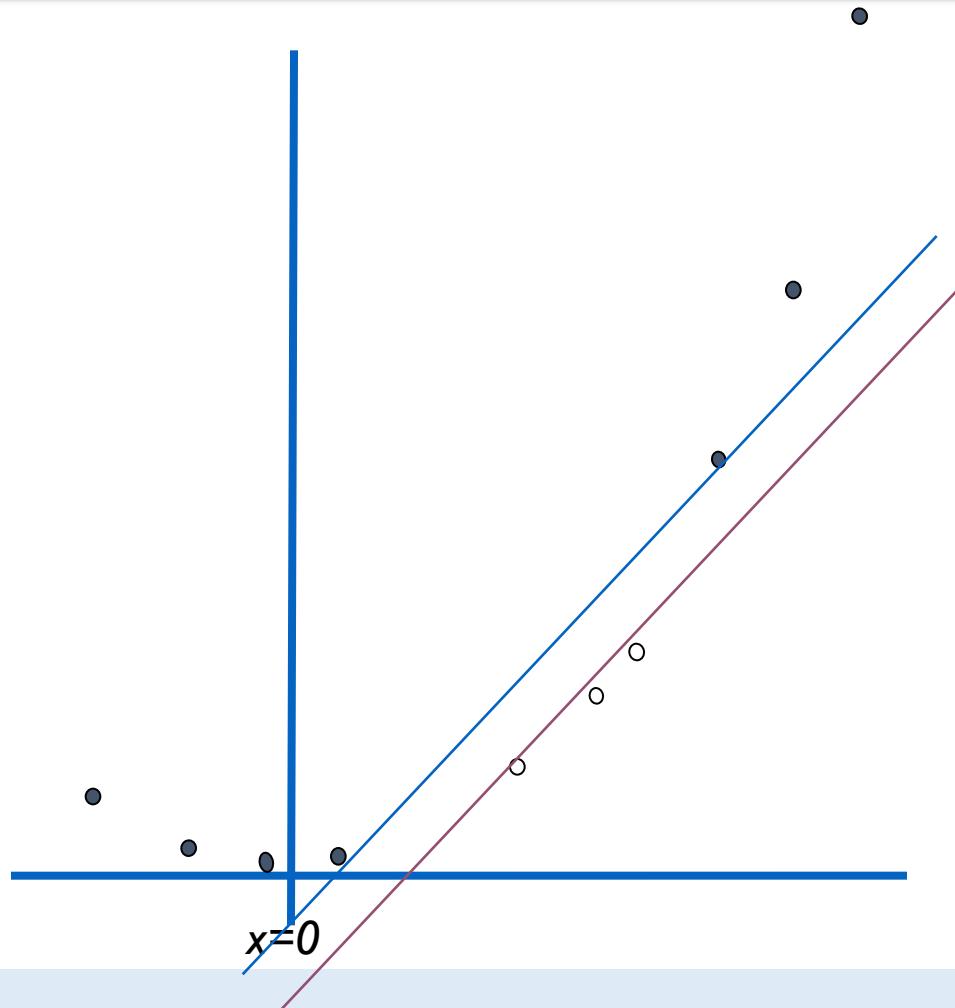
Harder 1-dimensional Dataset



Apply the following map

$$\mathbf{z}_k = (x_k, x_k^2)$$

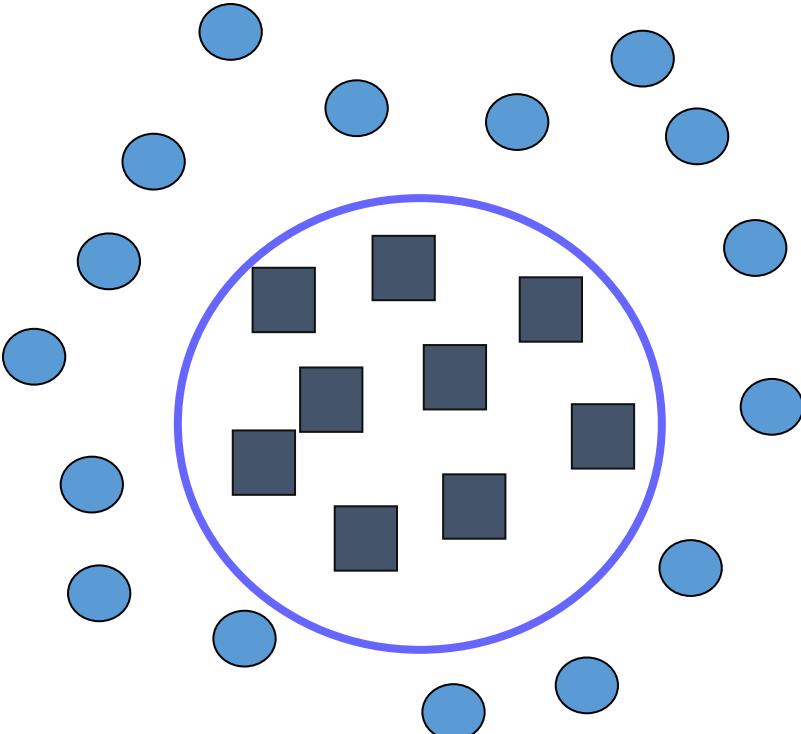
Harder 1-dimensional Dataset



Apply the following map

$$\mathbf{z}_k = (x_k, x_k^2)$$

Harder 2-dimensional Dataset



Apply the following map

$$\mathbf{z}_k = (x_k, y_k, x_k^2, y_k^2, x_k y_k)$$

Common SVM Basis Functions

\mathbf{z}_k = (polynomial terms of \mathbf{x}_k of degree 1 to q)

\mathbf{z}_k = (radial basis functions of \mathbf{x}_k)

$$\mathbf{z}_k[j] = \varphi_j(\mathbf{x}_k) = \exp\left(-\frac{|\mathbf{x}_k - \mathbf{c}_j|^2}{\sigma^2}\right)$$

\mathbf{z}_k = (sigmoid functions of \mathbf{x}_k)

Recall: Dual Form

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\mathbf{x}_k \cdot \mathbf{x}_l)$

Subject to these constraints:

$$0 \leq \alpha_k \leq c \quad \forall k \quad \sum_{k=1}^R \alpha_k y_k = 0$$

Then define:

$$\mathbf{w} = \frac{1}{2} \sum_{k=1}^R \alpha_k y_k \mathbf{x}_k$$

$$y_i(\mathbf{x}_i \bullet \mathbf{w} + b) - 1 = 0 \\ b = -y_i(\mathbf{x}_i \bullet \mathbf{w}) - 1$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, b) = sign(\mathbf{w} \cdot \mathbf{x} + b)$$

SVM QP with Basis Functions

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_l))$

Subject to these constraints:

$$0 \leq \alpha_k \leq C \quad \forall k$$

$$\sum_{k=1}^R \alpha_k y_k = 0$$

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \Phi(\mathbf{x}_k)$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \cdot \mathbf{f}(\mathbf{x}) - b)$$

Most important change:

$$\mathbf{x} \rightarrow \Phi(\mathbf{x})$$

SVM QP with Basis Functions

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_l))$

Subject to these constraints:

$$0 \leq \alpha_k \leq C$$

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \Phi(\mathbf{x}_k)$$

We must do $R^2/2$ dot products to get this matrix ready

Assuming a quadratic polynomial kernel, each dot product requires $m^2/2$ additions and multiplications

The whole thing costs $R^2 m^2 / 4$.

$$\mathbf{x} \rightarrow \Phi(\mathbf{x})$$

$$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) = \begin{pmatrix} 1 \\ \sqrt{2}a_1 \\ \sqrt{2}a_2 \\ \vdots \\ \sqrt{2}a_m \\ a_1^2 \\ a_2^2 \\ \vdots \\ a_m^2 \\ \sqrt{2}a_1a_2 \\ \sqrt{2}a_1a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \sqrt{2}a_2a_3 \\ \vdots \\ \sqrt{2}a_1a_m \\ \vdots \\ \sqrt{2}a_{m-1}a_m \end{pmatrix} \bullet \begin{pmatrix} 1 \\ \sqrt{2}b_1 \\ \sqrt{2}b_2 \\ \vdots \\ \sqrt{2}b_m \\ b_1^2 \\ b_2^2 \\ \vdots \\ b_m^2 \\ \sqrt{2}b_1b_2 \\ \sqrt{2}b_1b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \sqrt{2}b_2b_3 \\ \vdots \\ \sqrt{2}b_1b_m \\ \vdots \\ \sqrt{2}b_{m-1}b_m \end{pmatrix}$$

1
+
 $\sum_{i=1}^m 2a_i b_i$
+
 $\sum_{i=1}^m a_i^2 b_i^2$
+
 $\sum_{i=1}^m \sum_{j=i+1}^m 2a_i a_j b_i b_j$

Quadratic Dot Products

Quadratic Dot Products

Just out of interest, let's look at another function of \mathbf{a} and \mathbf{b} :

$$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) = \\ 1 + 2 \sum_{i=1}^m a_i b_i + \sum_{i=1}^m a_i^2 b_i^2 + \sum_{i=1}^m \sum_{j=i+1}^m 2 a_i a_j b_i b_j$$

$$\begin{aligned} & (\mathbf{a} \cdot \mathbf{b} + 1)^2 \\ &= (\mathbf{a} \cdot \mathbf{b})^2 + 2\mathbf{a} \cdot \mathbf{b} + 1 \\ &= \left(\sum_{i=1}^m a_i b_i \right)^2 + 2 \sum_{i=1}^m a_i b_i + 1 \\ &= \sum_{i=1}^m \sum_{j=1}^m a_i b_i a_j b_j + 2 \sum_{i=1}^m a_i b_i + 1 \\ &= \sum_{i=1}^m (a_i b_i)^2 + 2 \sum_{i=1}^m \sum_{j=i+1}^m a_i b_i a_j b_j + 2 \sum_{i=1}^m a_i b_i + 1 \end{aligned}$$

Quadratic Dot Products

They're the same!

And this is only $O(m)$ to compute!

$$\Phi(\mathbf{a}) \bullet \Phi(\mathbf{b}) =$$

$$1 + 2 \sum_{i=1}^m a_i b_i + \sum_{i=1}^m a_i^2 b_i^2 + \sum_{i=1}^m \sum_{j=i+1}^m 2 a_i a_j b_i b_j$$

Just out of interest, let's look at another function of \mathbf{a} and \mathbf{b} :

$$(\mathbf{a} \cdot \mathbf{b} + 1)^2$$

$$= (\mathbf{a} \cdot \mathbf{b})^2 + 2\mathbf{a} \cdot \mathbf{b} + 1$$

$$= \left(\sum_{i=1}^m a_i b_i \right)^2 + 2 \sum_{i=1}^m a_i b_i + 1$$

$$= \sum_{i=1}^m \sum_{j=1}^m a_i b_i a_j b_j + 2 \sum_{i=1}^m a_i b_i + 1$$

$$= \sum_{i=1}^m (a_i b_i)^2 + 2 \sum_{i=1}^m \sum_{j=i+1}^m a_i b_i a_j b_j + 2 \sum_{i=1}^m a_i b_i + 1$$

SVM QP with Basis Functions

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l (\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_l))$

Subject to these constraints:

$$0 \leq \alpha_k \leq C$$

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \Phi(\mathbf{x}_k)$$

We must do $R^2/2$ dot products to get this matrix ready

Now, each dot product now only requires m additions and multiplications

प्राप्ति विकास
विद्यालय

Most important change:

$$\mathbf{x} \rightarrow \Phi(\mathbf{x})$$

Higher-Order Polynomials

Poly-nomial	$f(x)$	Cost to build Q_{kl} matrix traditionally	Cost if 100 dimensions	$f(a).f(b)$	Cost to build Q_{kl} matrix sneakily	Cost if 100 dimensions
Quadratic	All $m^2/2$ terms up to degree 2	$m^2 R^2 / 4$	2,500 R^2	$(a.b+I)^2$	$m R^2 / 2$	50 R^2
Cubic	All $m^3/6$ terms up to degree 3	$m^3 R^2 / 12$	83,000 R^2	$(a.b+I)^3$	$m R^2 / 2$	50 R^2
Quartic	All $m^4/24$ terms up to degree 4	$m^4 R^2 / 48$	1,960,000 R^2	$(a.b+I)^4$	$m R^2 / 2$	50 R^2

SVM QP with Basis Functions

Maximize $\sum_{k=1}^R \alpha_k - \frac{1}{2} \sum_{k=1}^R \sum_{l=1}^R \alpha_k \alpha_l Q_{kl}$ where $Q_{kl} = y_k y_l K(\mathbf{x}_k, \mathbf{x}_l)$

Subject to these constraints:

$$0 \leq \alpha_k \leq C \quad \forall k$$

$$\sum_{k=1}^R \alpha_k y_k$$

Kernel gram matrix

Then define:

$$\mathbf{w} = \sum_{k \text{ s.t. } \alpha_k > 0} \alpha_k y_k \Phi(\mathbf{x}_k)$$

Then classify with:

$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(K(\mathbf{w}, \mathbf{x}) - b)$$

Most important change:

$$\Phi(\mathbf{x}_k) \cdot \Phi(\mathbf{x}_l) \rightarrow K(\mathbf{x}_k, \mathbf{x}_l)$$

SVM Kernel Functions

- $K(\mathbf{a}, \mathbf{b}) = (\mathbf{a} \cdot \mathbf{b} + 1)^d$ is an example of an SVM Kernel Function
- Beyond polynomials there are other very high dimensional basis functions that can be made practical by finding the right Kernel Function
 - Radial-Basis-style Kernel Function:

$$K(\mathbf{a}, \mathbf{b}) = \exp\left(-\frac{(\mathbf{a} - \mathbf{b})^2}{2\sigma^2}\right)$$

- Sigmoidal function

Kernel Tricks

- Replacing dot product with a kernel function
- Not all functions are kernel functions
 - Need to be decomposable
 - $K(a,b) = \phi(a) \cdot \phi(b)$
- **Mercer's condition**

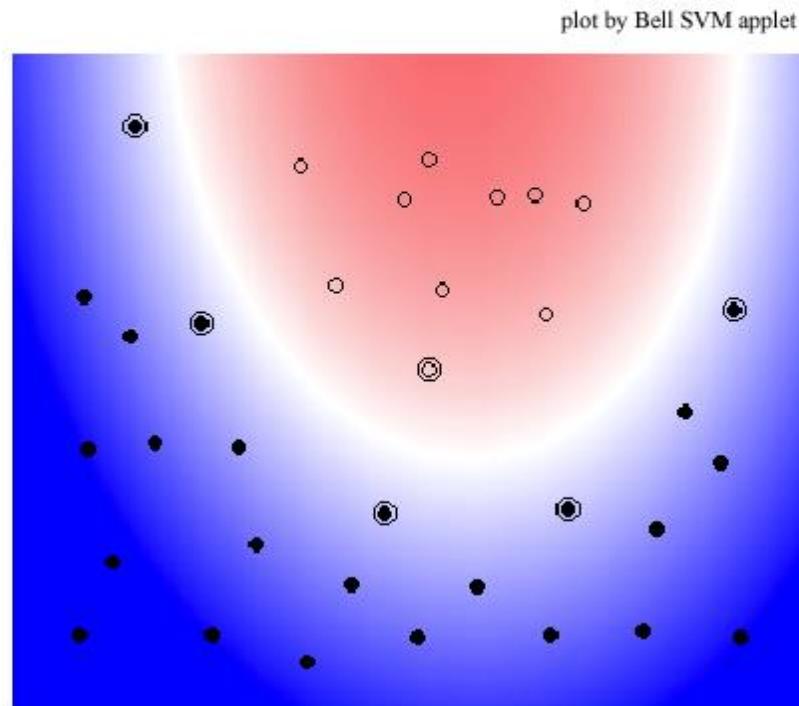
To expand Kernel function $K(x,y)$ into a dot product, i.e.
 $K(x,y)=\Phi(x)\cdot\Phi(y)$, $K(x, y)$ has to be positive semi-definite function,
i.e., for any function $f(x)$ whose $\int f^2(x)dx$ is finite, the following
inequality holds

$$\int dx dy f(x) K(x, y) f(y) \geq 0$$

Non-Linear Kernel: Visualization

Example: SVM with Polynomial of Degree 2

Kernel: $K(\vec{x}_i, \vec{x}_j) = [\vec{x}_i \cdot \vec{x}_j + 1]^2$



How to choose a kernel function?

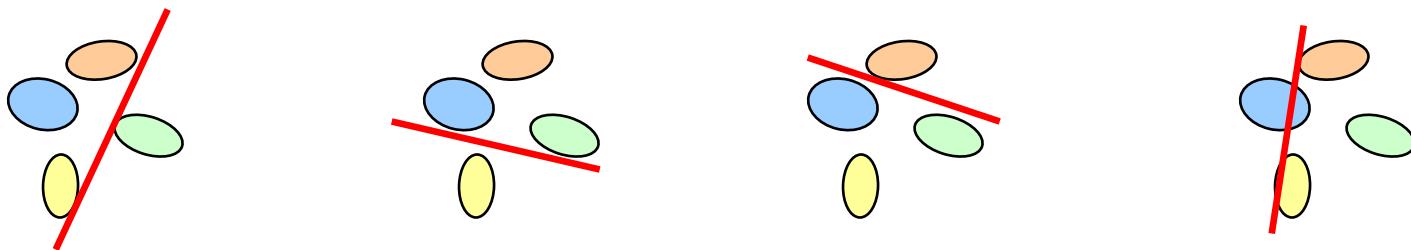
- Not easy! Remember – this depends on your data geometry
- If linear works, go with it
- RBF kernels are considered good in general, especially for images (and other smooth functions/data)
- For discrete data, chi-square kernel preferred of late (especially for histogram data)
- You can also do Multiple Kernel Learning
 - Beyond scope of current course – you can read more at:
https://en.wikipedia.org/wiki/Multiple_kernel_learning
- Still not sure? Use cross-validation to select a kernel function from some basic options

From Binary Classification to Multi-class Classification

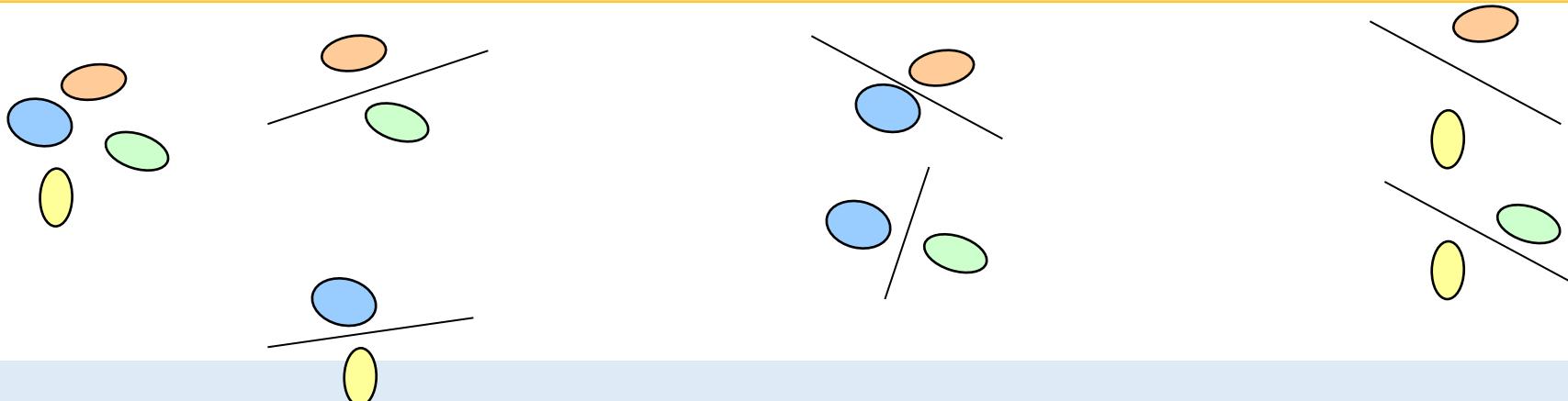
- SVMs can only handle two-class outputs.
- What can be done?
- Answer: with output arity N, learn N SVM's
 - SVM 1 learns "Output==1" vs "Output != 1"
 - SVM 2 learns "Output==2" vs "Output != 2"
 - :
 - SVM N learns "Output==N" vs "Output != N"
- Then to predict the output for a new input, just predict with each SVM and find out which one puts the prediction the furthest into the positive region.
- Other approaches
 - Pair-wise SVM, Tree-structured SVM

Multi-class Classification using SVM

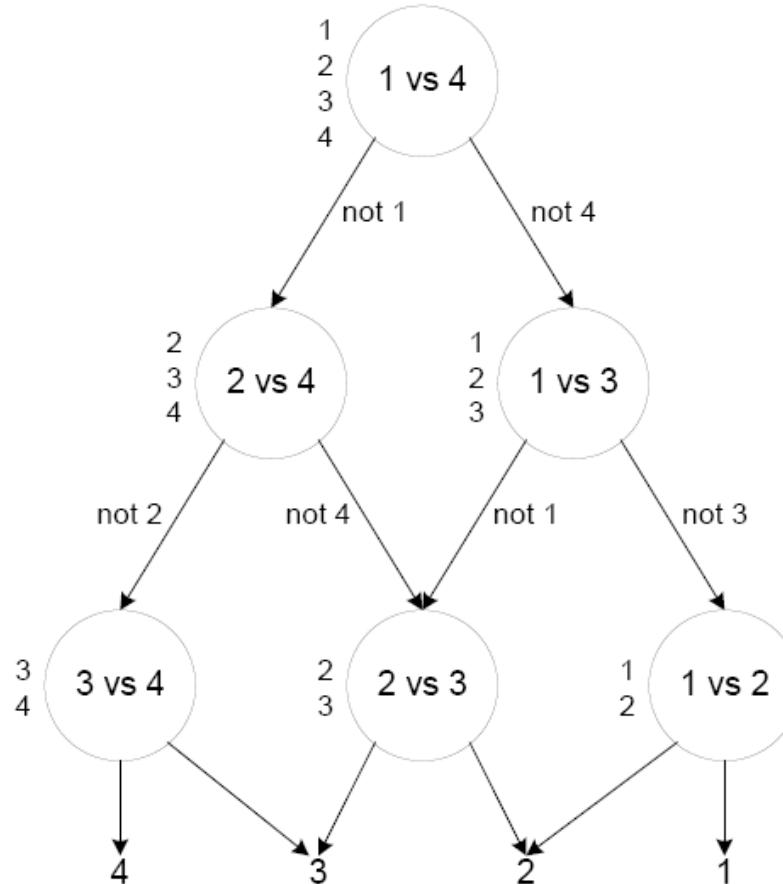
One- versus-all



One- versus-one



Tree-Structured SVM



Readings

- “Introduction to Machine Learning” by Ethem Alpaydin, Chapters 8, 9, 13, 16.1-16.3