



Samueli
Computer Science



CS145 Discussion: Week 3

Decision Tree & SVM

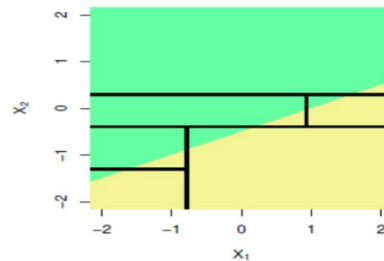
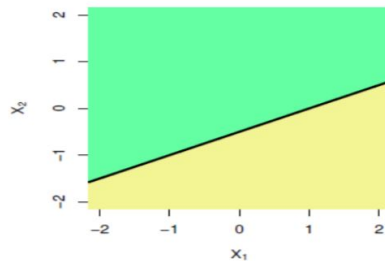
Junheng Hao
Friday, 10/23/2020

- Announcement
- Decision Tree
- SVM (Part I)

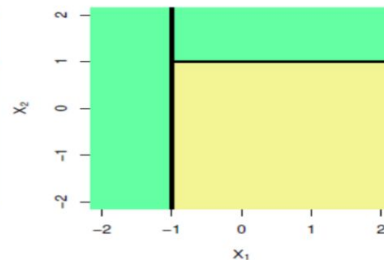
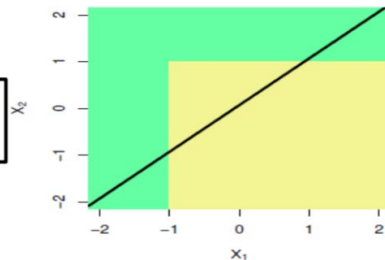
- Homework 1 due on **Oct 30 (Friday) 11:59 PT**
 - Submit through GradeScope of 1 PDF (2 python file and 1 jupyter notebook into 1 PDF file)
 - Assign pages to the questions on GradeScope
- Group formation
 - Please email the TA whose session you're enrolled in for help if you cannot find a group with 4-5 members.
 - You may also find 1 or 2 additional team members if your group has someone who has dropped the class (before the end of Week 3)

- Comparison: Logistic Regression vs Decision Tree

**Ground Truth:
Linear Boundary**



**Ground Truth:
Non-Linear Boundary**



**Fitted Model:
Linear Model**

**Fitted Model:
Trees**

One more question on logistics regression:

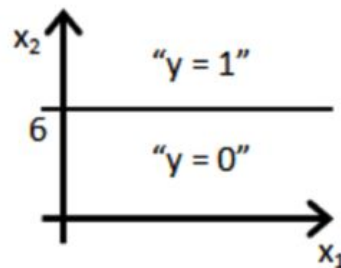
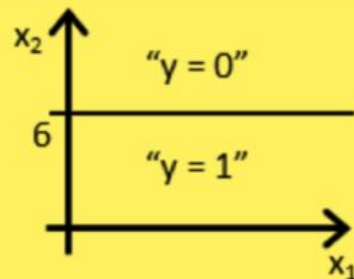
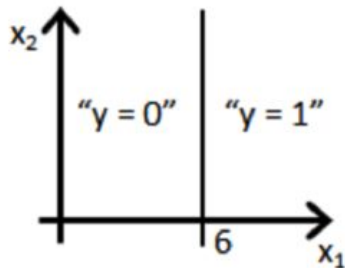
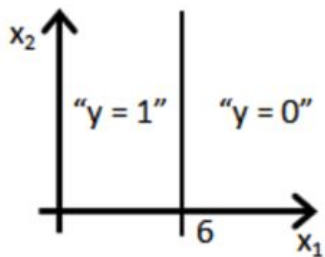
Suppose you train a logistic classifier $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$. Suppose $\theta_0 = 6, \theta_1 = 0, \theta_2 = -1$. Which of the following figures represents the decision boundary found by your classifier?

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

Decision Boundary: Exercise

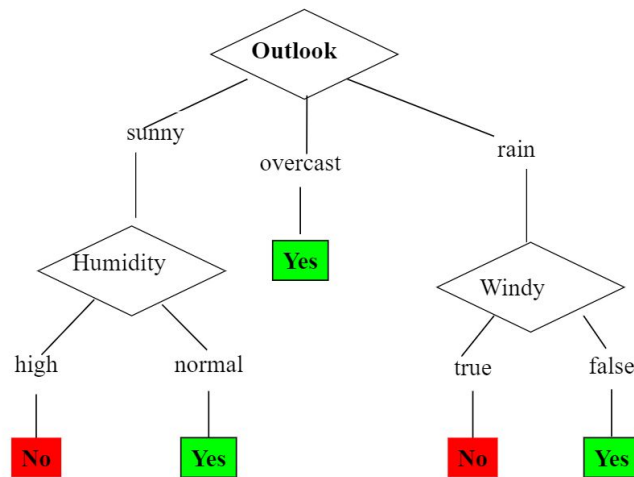


Suppose you train a logistic classifier $h_{\theta}(x) = g(\theta_0 + \theta_1 x_1 + \theta_2 x_2)$. Suppose $\theta_0 = 6, \theta_1 = 0, \theta_2 = -1$. Which of the following figures represents the decision boundary found by your classifier?



- Decision Tree Classification: From data to model

Outlook	Temperature	Humidity	Windy	Play?
sunny	hot	high	false	No
sunny	hot	high	true	No
overcast	hot	high	false	Yes
rain	mild	high	false	Yes
rain	cool	normal	false	Yes
rain	cool	normal	true	No
overcast	cool	normal	true	Yes
sunny	mild	high	false	No
sunny	cool	normal	false	Yes
rain	mild	normal	false	Yes
sunny	mild	normal	true	Yes
overcast	mild	high	true	Yes
overcast	hot	normal	false	Yes
rain	mild	high	true	No

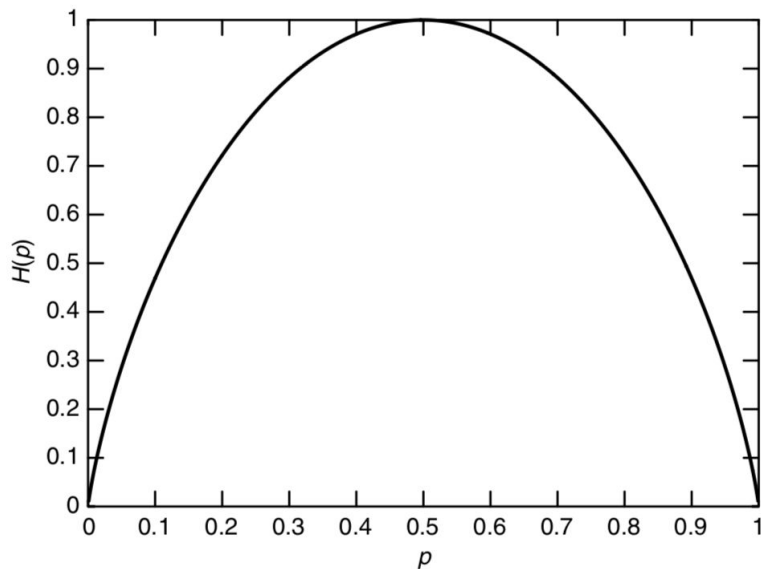


- Choosing the Splitting Attribute
- At each node, available attributes are evaluated on the basis of separating the classes of the training examples.
- A goodness function (information measurement) is used for this purpose:
 - **Information Gain**
 - **Gain Ratio**
 - Gini Index*

- Which is the best attribute?
 - The one which will result in the smallest tree
 - Heuristic: choose the attribute that produces the “purest” nodes
- Popular *impurity criterion: information gain*
 - Information gain increases with the average purity of the subsets that an attribute produces
- Strategy: choose attribute that results in greatest information gain

$$X = \begin{cases} 1 & \text{with probability } p \\ 0 & \text{with probability } 1 - p \end{cases}$$

$$H(X) = -p \log p - (1 - p) \log(1 - p) \stackrel{\text{def}}{=} H(p)$$



- Information in a split with x items of one class, y items of the second class

$$\begin{aligned}\text{info}([x, y]) &= \text{entropy}\left(\frac{x}{x+y}, \frac{y}{x+y}\right) \\ &= -\frac{x}{x+y} \log\left(\frac{x}{x+y}\right) - \frac{y}{x+y} \log\left(\frac{y}{x+y}\right)\end{aligned}$$

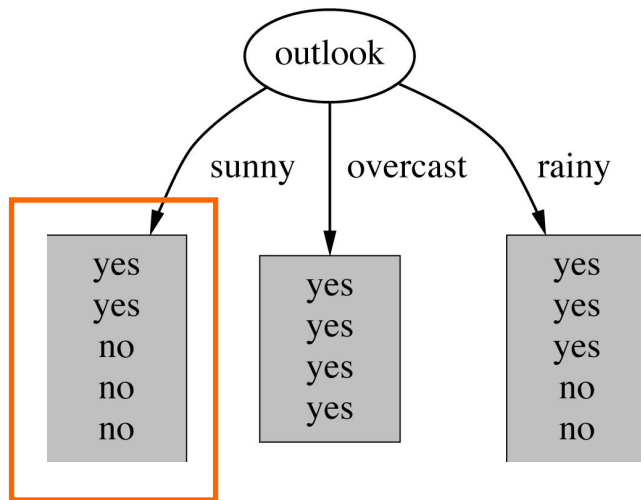
Decision Tree: Example for Practice

Attribute: “Outlook” = “Sunny”



- “Outlook” = “Sunny”: 2 and 3 split

$$\text{info}([2,3]) = \text{entropy}(2/5, 3/5) = -\frac{2}{5}\log\left(\frac{2}{5}\right) - \frac{3}{5}\log\left(\frac{3}{5}\right) = 0.971 \text{ bits}$$



Decision Tree: Example for Practice

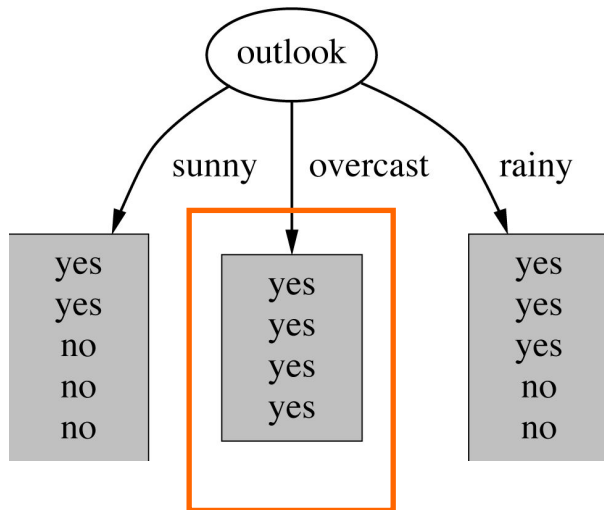
Attribute: “Outlook” = “Overcast”



- “Outlook” = “Overcast”: 4/0 split

$$\text{info}([4,0]) = \text{entropy}(1,0) = -1\log(1) - 0\log(0) = 0 \text{ bits}$$

Note: $\log(0)$ is not defined, but we evaluate $0 \cdot \log(0)$ as zero.



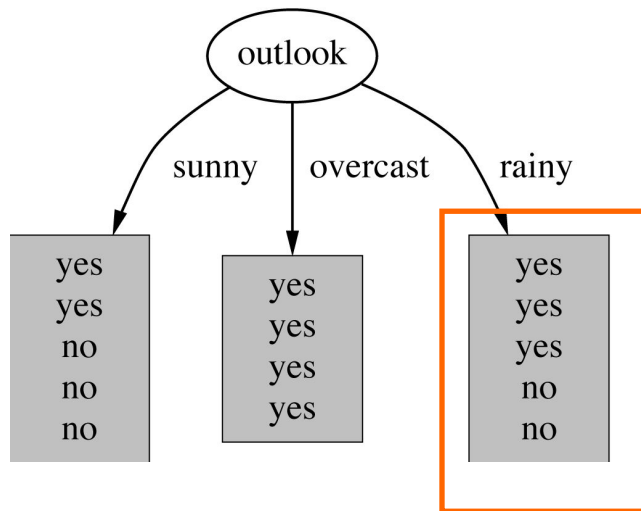
Decision Tree: Example for Practice

Attribute: “Outlook” = “Rainy”



- “Outlook” = “Rainy”:

$$\text{info}([3,2]) = \text{entropy}(3/5, 2/5) = -\frac{3}{5} \log\left(\frac{3}{5}\right) - \frac{2}{5} \log\left(\frac{2}{5}\right) = 0.971 \text{ bits}$$



Expected Information of Attribute “Outlook”

Expected information for attribute:

$$\begin{aligned}\text{info}([3,2],[4,0],[3,2]) &= (5/14) \times 0.971 + (4/14) \times 0 + (5/14) \times 0.971 \\ &= 0.693 \text{ bits}\end{aligned}$$

Information gain:

(information before split) – (information after split)

$$\begin{aligned}\text{gain("Outlook")} &= \text{info}([9,5]) - \text{info}([2,3],[4,0],[3,2]) = 0.940 - 0.693 \\ &= 0.247 \text{ bits}\end{aligned}$$

Information gain for attributes from all weather data:

$$\text{gain("Outlook")} = 0.247 \text{ bits}$$

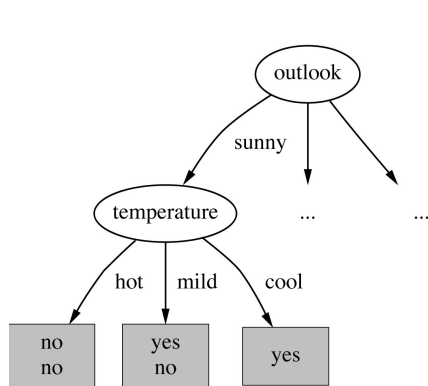
$$\text{gain("Temperature")} = 0.029 \text{ bits}$$

$$\text{gain("Humidity")} = 0.152 \text{ bits}$$

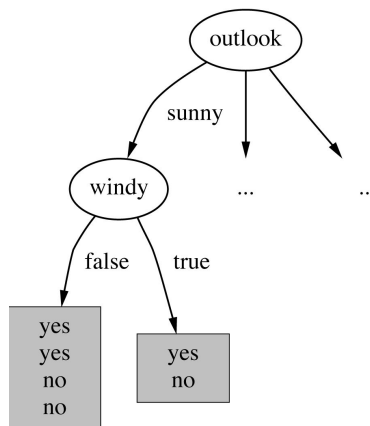
$$\text{gain("Windy")} = 0.048 \text{ bits}$$

Decision Tree: Example for Practice

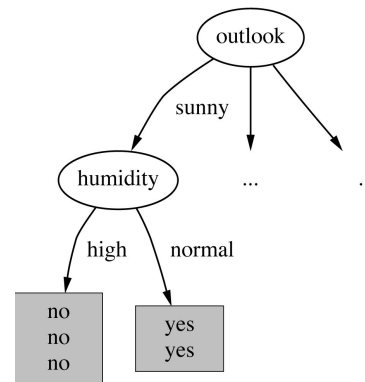
Continue to Split



gain("Temperature") = 0.571 bits



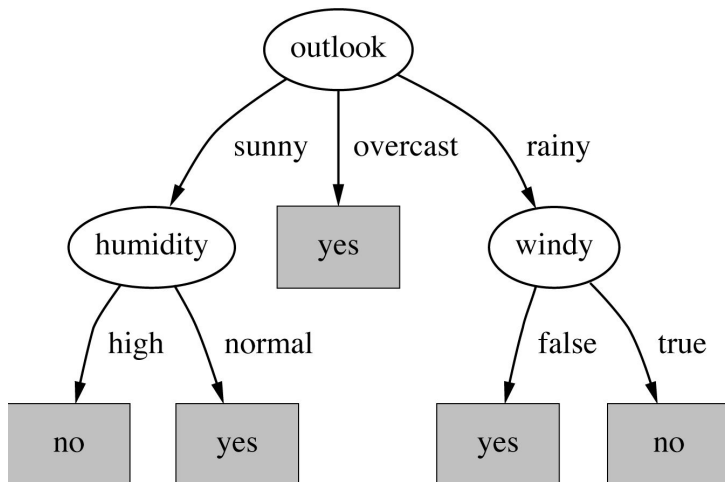
gain("Windy") = 0.020 bits



gain("Humidity") = 0.971 bits

Decision Tree: Example for Practice

Final Tree



- Note: Not all leaves need to be pure. Sometimes identical instances have different classes. → Splitting can stop when data can't be split any further

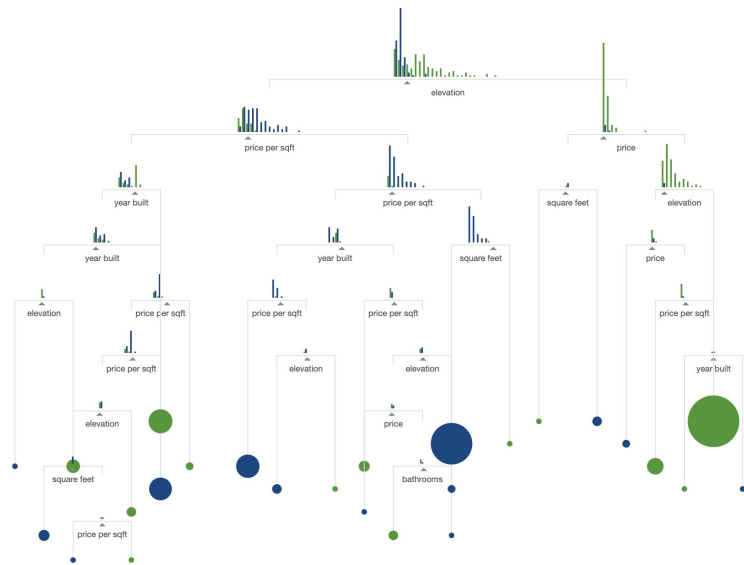
- SplitInfo and Gain Ratio

$$SplitInfo_A(D) = -\sum_{j=1}^v \frac{|D_j|}{|D|} \times \log_2\left(\frac{|D_j|}{|D|}\right)$$

$$GainRatio(A) = \text{Gain(A)} / \text{SplitInfo(A)}$$

- Why Gain Ratio?
 - Information gain: biased towards attributes with a large number of values
- Practice: What is the gain ratio for attribute “Outlook” in the previous example?

- Demo links
 - <http://www.r2d3.us/visual-intro-to-machine-learning-part-1/>
 - <http://explained.ai/decision-tree-viz/>
- Does decision tree also have the bias-variance trade-off?
 - A visual demo: <http://www.r2d3.us/visual-intro-to-machine-learning-part-2/>



- Links: <https://cs.stanford.edu/people/karpathy/svmjs/demo/>



- Hyperplane separating the data points

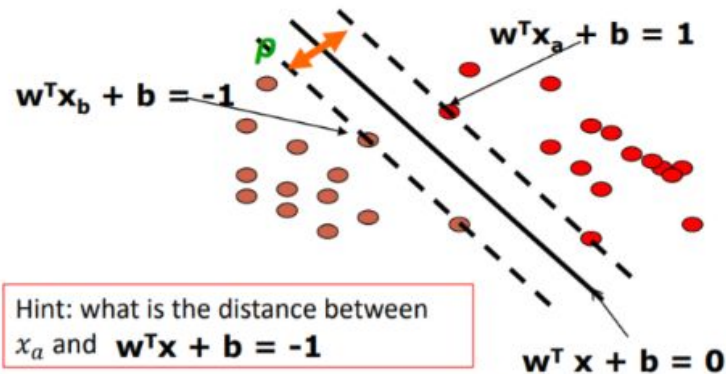
$$\mathbf{w}^T \mathbf{x} + b = 0$$

- Maximize margin

$$\rho = \frac{2}{\|\mathbf{w}\|}$$

- Solution by solving its dual problem

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = \sum_{k: \alpha_k \neq 0} (y_k - \mathbf{w}^T \mathbf{x}_k) / N_k$$



SVM: Start with the margin



Margin Lines

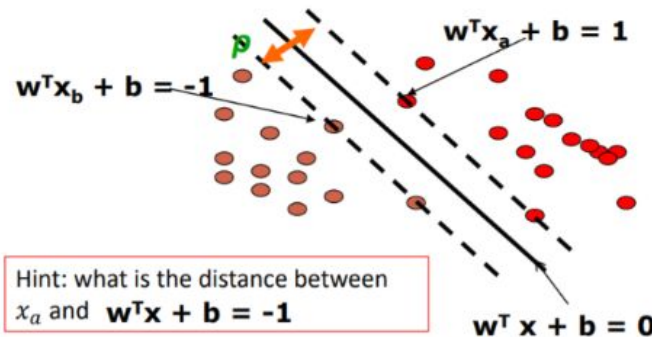
$$\mathbf{w}^T \mathbf{x}_a + b = 1 \quad \mathbf{w}^T \mathbf{x}_b + b = -1$$

Distance between parallel lines of $ax_1 + bx_2 = c_1/c_2$

$$d = \frac{|c_2 - c_1|}{\sqrt{a^2 + b^2}}$$

Margin

$$\rho = \frac{|(b+1) - (b-1)|}{\|\mathbf{w}\|} = \frac{2}{\|\mathbf{w}\|}$$



1. Formulation of the Linear SVM problem: maximizing margin
2. Formulation of Quadratic Programming (optimization with linear constraints) → Primal problem
3. Solving linear SVM problem with “great” math*
 - a. (Generalized) Lagrange function, lagrange multiplier
 - b. Identify primal and dual problem (duality) → KKT conditions
 - c. Solution to w and b regarding α
4. Support Vectors, SVM Classifier Inference
5. Non-linear SVM, Kernel tricks

- Slides: <http://people.csail.mit.edu/dsontag/courses/ml13/slides/lecture6.pdf>
- Notes: <https://see.stanford.edu/materials/aimlcs229/cs229-notes3.pdf>

**To show in hand notes*

Given labeled data and alpha values

- Positively labeled data points (1 to 4)

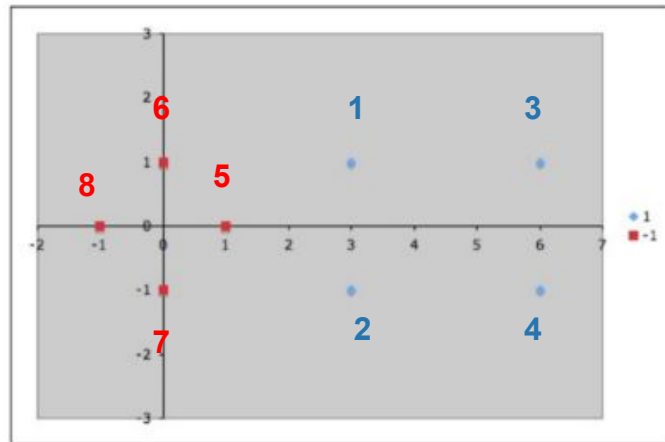
$$\left\{ \begin{pmatrix} 3 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ -1 \end{pmatrix}, \begin{pmatrix} 6 \\ 1 \end{pmatrix}, \begin{pmatrix} 6 \\ -1 \end{pmatrix} \right\}$$

- Negatively labeled data points (5 to 8)

$$\left\{ \begin{pmatrix} 1 \\ 0 \end{pmatrix}, \begin{pmatrix} 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 0 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 0 \end{pmatrix} \right\}$$

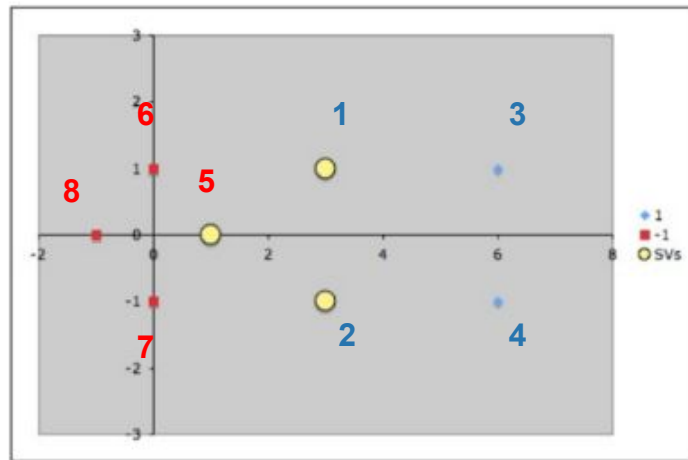
- Alpha values

- $\alpha_1 = 0.25$
- $\alpha_2 = 0.25$
- $\alpha_5 = 0.5$
- Others = 0



- Which points are support vectors?
- Calculate normal vector of hyperplane: \mathbf{w}
- Calculate the bias term
- What is the decision boundary?
- Predict class of new point (4, 1)

$$\mathbf{w} = \sum \alpha_i y_i \mathbf{x}_i \quad b = \sum_{k:\alpha_k \neq 0} (y_k - \mathbf{w}^T \mathbf{x}_k) / N_k$$



Linear SVM: Example for Practice

Predictions for new data



$$y \leftarrow \text{sign}(\vec{w} \cdot \vec{x} + b)$$



Using dual solution

$$y \leftarrow \text{sign} \left[\sum_i \alpha_i y_i (\underbrace{\vec{x}_i \cdot \vec{x}}_{\text{dot product}}) + b \right]$$

dot product of feature vectors of
new example with support vectors

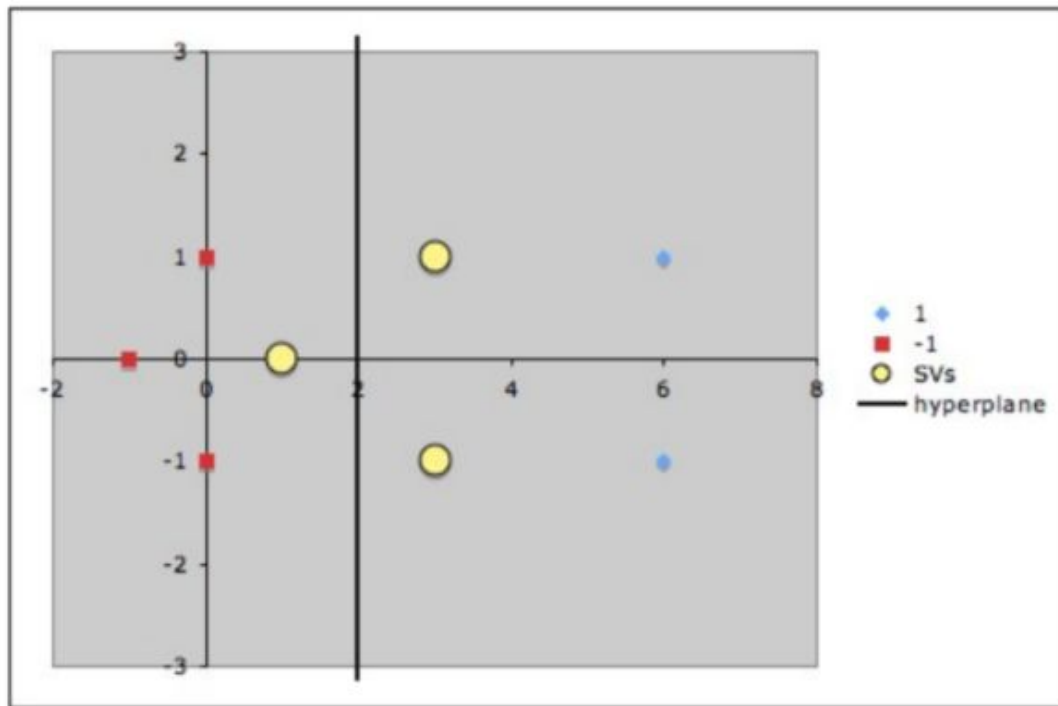
$$\mathbf{w} = \sum_i \alpha_i y_i \mathbf{x}_i$$

$$b = y_k - \mathbf{w} \cdot \mathbf{x}_k$$

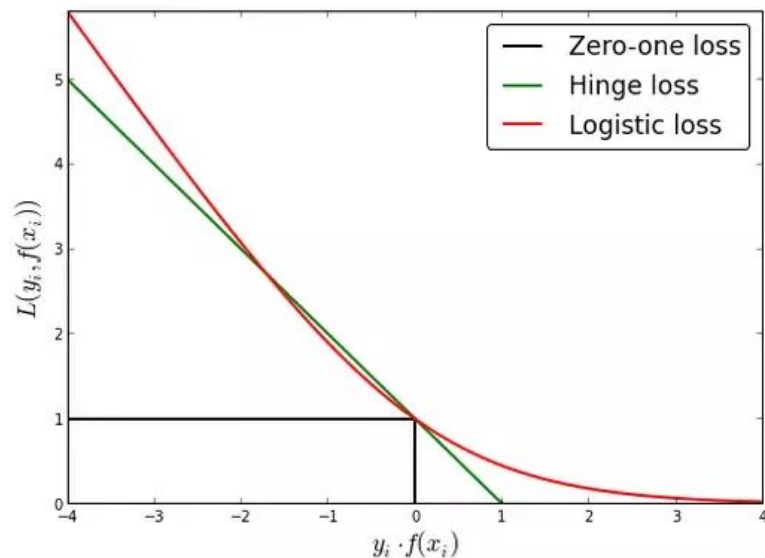
for any k where $C > \alpha_k > 0$

Linear SVM: Example for Practice

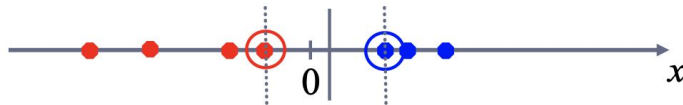
Plot



- Decision boundaries?
- Loss functions?



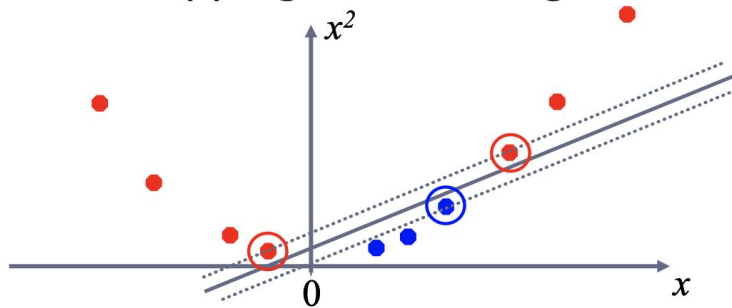
- Datasets that are linearly separable (with some noise) work out great:



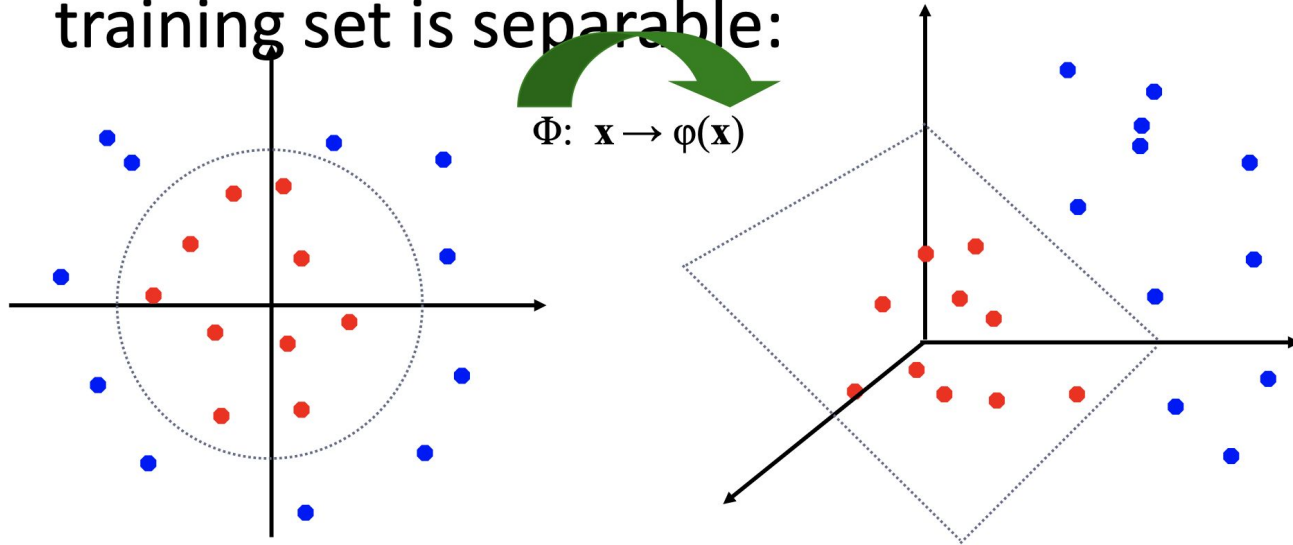
- But what are we going to do if the dataset is just too hard?



- How about ... mapping data to a higher-dimensional space:



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



$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

$$\text{maximize}_{\alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j)$$

$$K(\mathbf{x}_i, \mathbf{x}_j) = \Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)$$

$$\sum_i \alpha_i y_i = 0$$

$$C \geq \alpha_i \geq 0$$

- The linear SVM relies on an inner product between data vectors,

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$$

- If every data point is mapped into high-dimensional space via transformation, the inner product becomes,

$$K(\mathbf{x}_i, \mathbf{x}_j) = \phi^T(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$$

- Do we need to compute $\phi(\mathbf{x})$ explicitly for each data sample? → **Directly compute kernel function $K(\mathbf{x}_i, \mathbf{x}_j)$**

$$\begin{aligned}
 k(\mathbf{x}, \mathbf{z}) &= (\mathbf{x}^T \mathbf{z} + c)^2 = \left(\sum_{j=1}^n x^{(j)} z^{(j)} + c \right) \left(\sum_{\ell=1}^n x^{(\ell)} z^{(\ell)} + c \right) \\
 &= \sum_{j=1}^n \sum_{\ell=1}^n x^{(j)} x^{(\ell)} z^{(j)} z^{(\ell)} + 2c \sum_{j=1}^n x^{(j)} z^{(j)} + c^2 \\
 &= \sum_{j,\ell=1}^n (x^{(j)} x^{(\ell)}) (z^{(j)} z^{(\ell)}) + \sum_{j=1}^n (\sqrt{2c} x^{(j)}) (\sqrt{2c} z^{(j)}) + c^2,
 \end{aligned}$$

Feature mapping given by:

$$\Phi(\mathbf{x}) = [x^{(1)2}, x^{(1)} x^{(2)}, \dots, x^{(3)2}, \sqrt{2c} x^{(1)}, \sqrt{2c} x^{(2)}, \sqrt{2c} x^{(3)}, c]$$

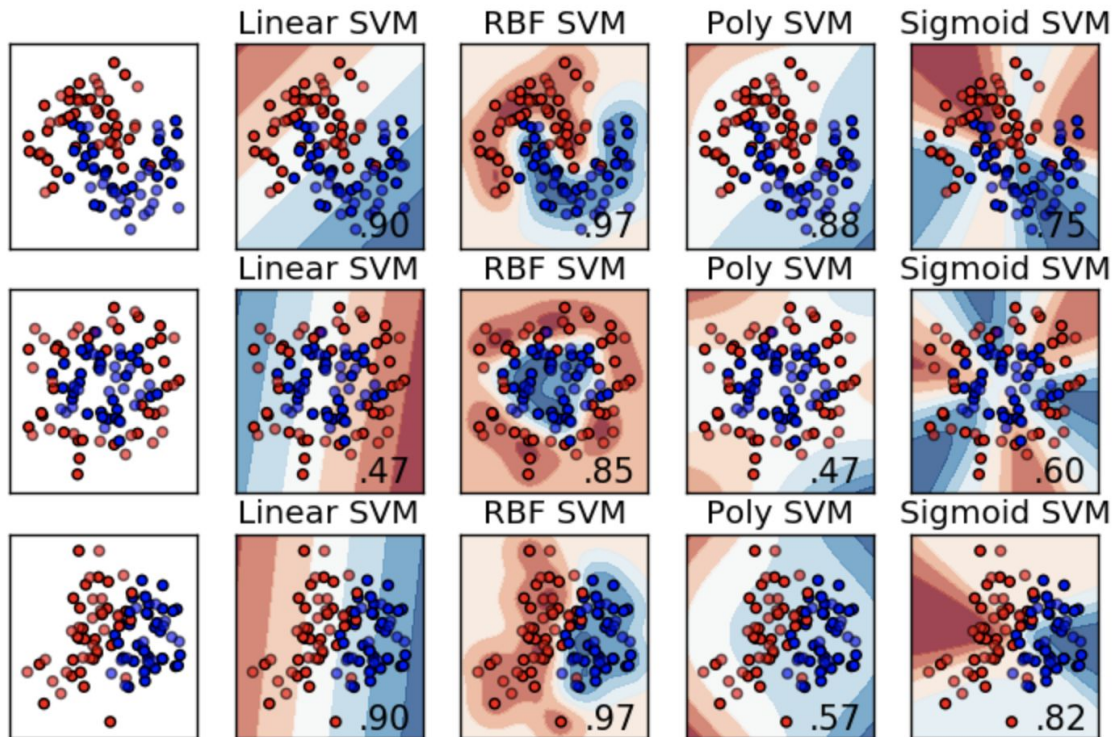
Polynomial kernel of degree h : $K(\mathbf{X}_i, \mathbf{X}_j) = (\mathbf{X}_i \cdot \mathbf{X}_j + 1)^h$

Gaussian radial basis function kernel : $K(\mathbf{X}_i, \mathbf{X}_j) = e^{-\|\mathbf{X}_i - \mathbf{X}_j\|^2 / 2\sigma^2}$

Sigmoid kernel : $K(\mathbf{X}_i, \mathbf{X}_j) = \tanh(\kappa \mathbf{X}_i \cdot \mathbf{X}_j - \delta)$

- Given the same data samples, what is the difference between linear kernel and non-linear kernel? Is the decision boundary linear (in original feature space)?

SVM: Demo of different kernels



- Positively labeled data points (1 to 4)

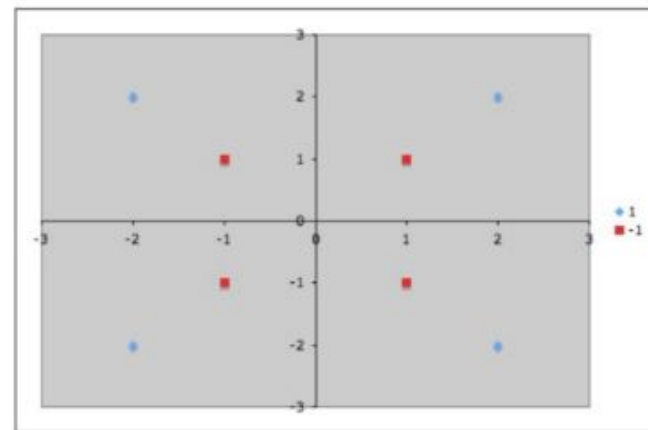
$$\left\{ \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 2 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ -2 \end{pmatrix}, \begin{pmatrix} -2 \\ 2 \end{pmatrix} \right\}$$

- Negatively labeled data points (5 to 8)

$$\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\}$$

- Non-linear mapping

$$\Phi_1 \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} = \begin{cases} \begin{pmatrix} 4 - x_2 \\ 4 - x_1 \\ x_1 \\ x_2 \end{pmatrix} & \text{if } \sqrt{x_1^2 + x_2^2} > 2 \\ \text{otherwise} \end{cases}$$



- New positively labeled data points (1 to 4)

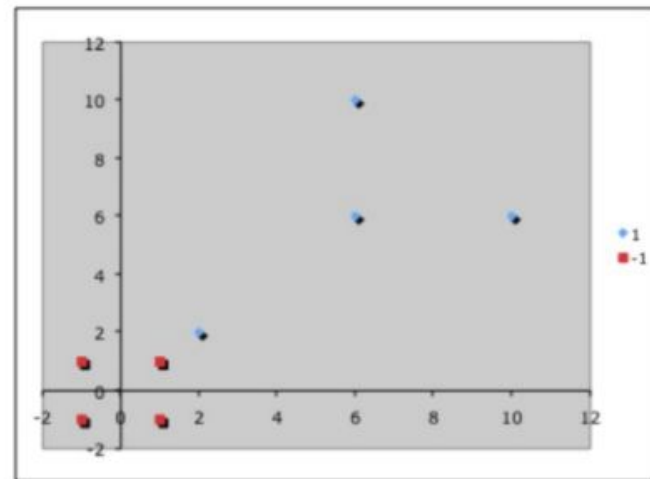
$$\left\{ \begin{pmatrix} 2 \\ 2 \end{pmatrix}, \begin{pmatrix} 6 \\ 2 \end{pmatrix}, \begin{pmatrix} 6 \\ 6 \end{pmatrix}, \begin{pmatrix} 2 \\ 6 \end{pmatrix} \right\}$$

- New negatively labeled data points (5 to 8)

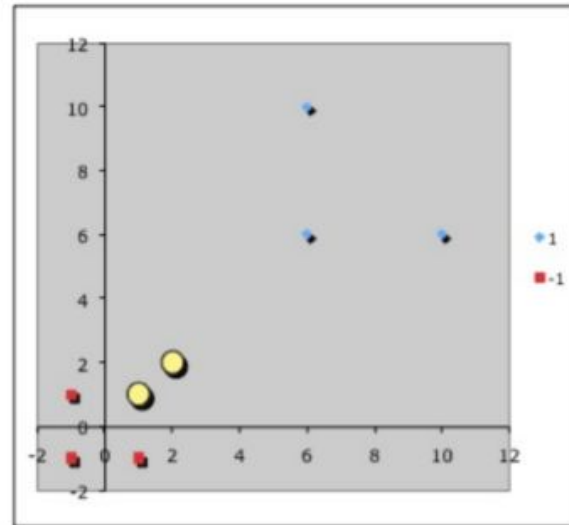
$$\left\{ \begin{pmatrix} 1 \\ 1 \end{pmatrix}, \begin{pmatrix} 1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ -1 \end{pmatrix}, \begin{pmatrix} -1 \\ 1 \end{pmatrix} \right\}$$

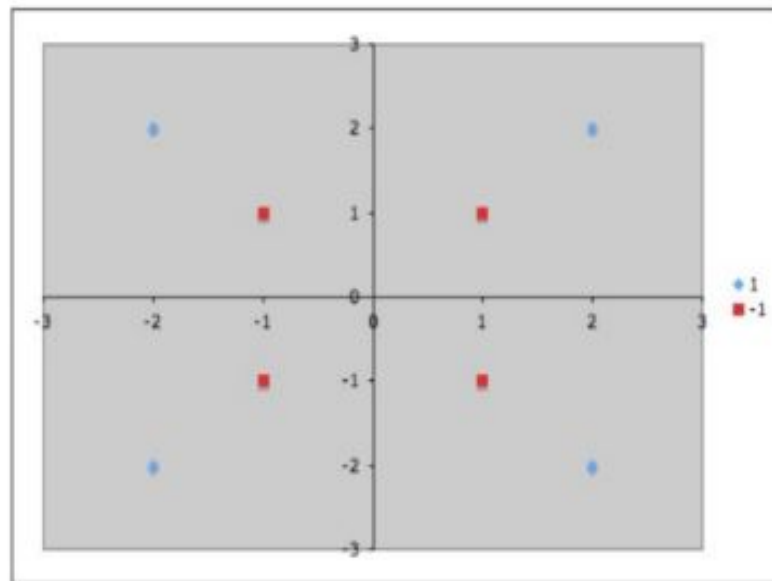
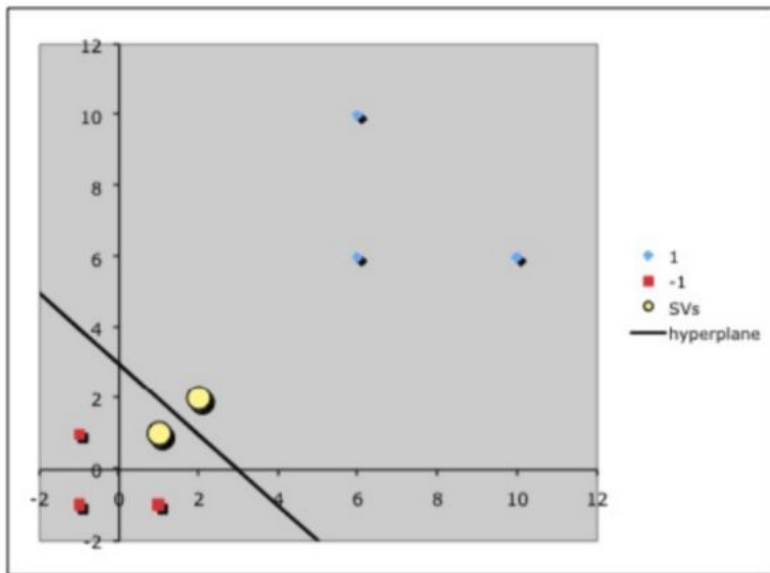
- Alpha values

- $\alpha_1 = 1.0$
- $\alpha_5 = 1.0$
- Others = 0



- Which points are support vectors?
- Calculate normal vector of hyperplane: \mathbf{w}
- Calculate the bias term
- What is the decision boundary?
- Predict class of new point (4, 5)







Samueli
Computer Science



Thank you!

Q & A