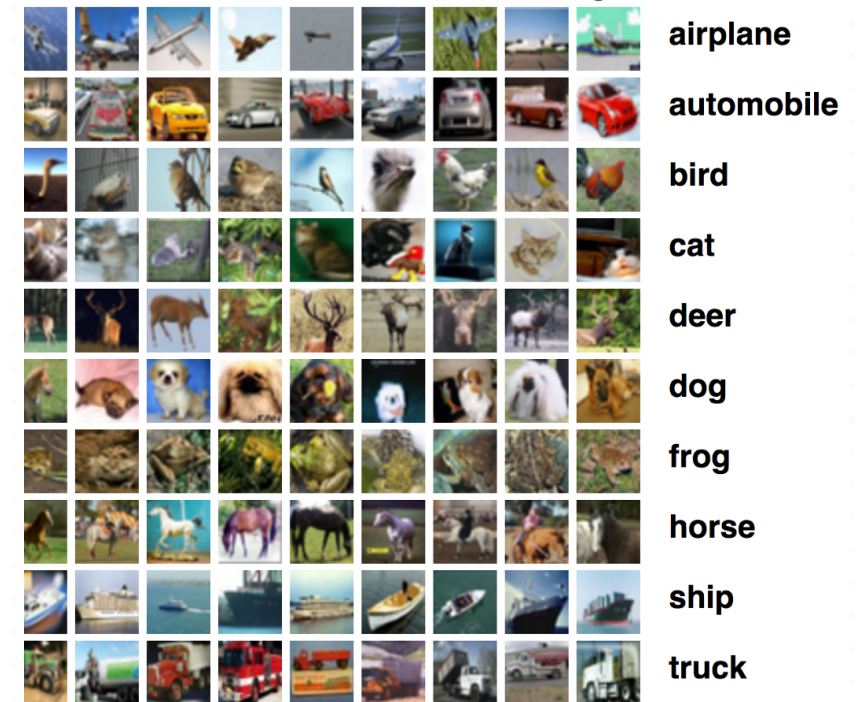


ECE 4803/8803: Fundamentals of Machine Learning (FunML)

Spring 2024

Lecture 4: Classifiers



Logistic Regression

Summary

- The modeling process we just went through is called logistic regression
- Logistic Regression is a **parametric model**
- Given feature vector $\mathbf{x} = [x_1, \dots, x_p]^T$, logistic regression classifier estimates the **posterior** $P(y|\mathbf{x})$ with *sigmoid function* $\sigma(\mathbf{x})$:

$$P(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b)$$

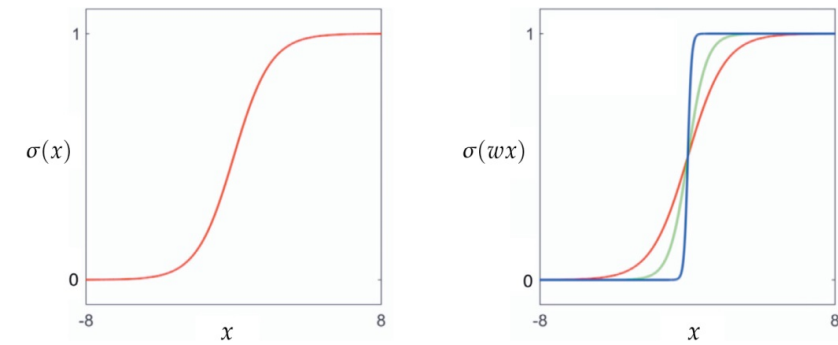
$$P(y = 0|\mathbf{x}) = 1 - \sigma(\mathbf{w}^T \mathbf{x} + b)$$

- Where $\mathbf{w} = [w_1, \dots, w_p]^T$ is **weight**, vector b is **bias**
- Classification rule:

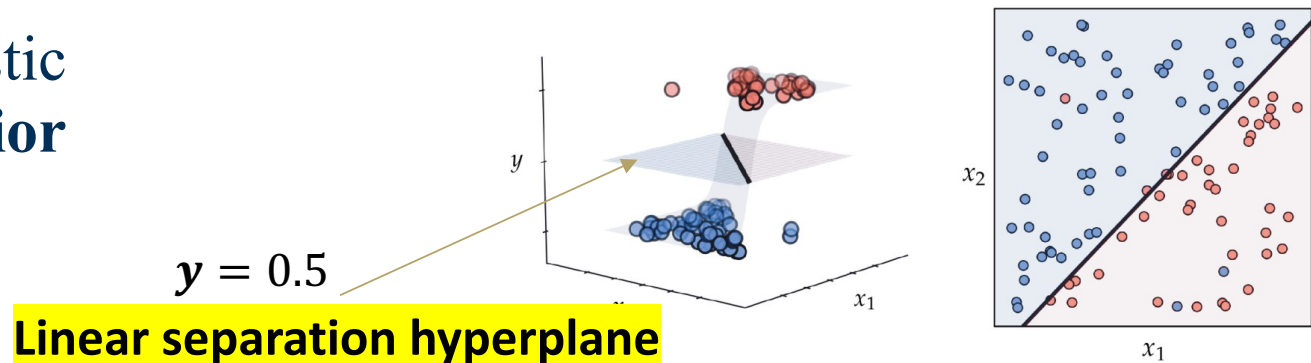
$$y = \begin{cases} 1 & P(y = 1|\mathbf{x}) > 0.5 \\ 0 & \text{otherwise} \end{cases}$$

$$\boldsymbol{\theta} = \{\mathbf{w}, b\}$$

$$\boldsymbol{\theta}(t+1) = \boldsymbol{\theta}(t) - \alpha \frac{\partial LL(\mathbf{x}, y, \boldsymbol{\theta})}{\partial \boldsymbol{\theta}}$$



Sigmoid function $\sigma(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}}}$



The overall negative log-likelihood cost function:

$$LL(\mathbf{x}, y) = -\frac{1}{N} \sum_{i=1}^N y_i \log(\sigma(\mathbf{x}_i)) + (1 - y_i) \log(1 - \sigma(\mathbf{x}_i))$$

N : number of samples, y_i : desired output, z_i : predicted output

$LL(\mathbf{x}, y)$ is referred to as the **Cross Entropy Cost** for logistic regression. In binary classification, this is also referred to as the **Binary Cross Entropy Cost**

Logistic Regression

Example: Sentiment Classification

Let's assume we want to predict the binary sentiment for the following movie review:

Input: “It's hokey. There are virtually no surprises, and the writing is second-rate. So why was it so enjoyable? For one thing, the cast is great. Another nice touch is the music. I was overcome with the urge to get off the couch and start dancing. It sucked me in, and it'll do the same to you.”

Output: positive (1) or negative (0)

Logistic Regression

Example: Sentiment Classification

Input: “It's hokey . There are virtually no surprises , and the writing is second-rate. So why was it so enjoyable? For one thing, the cast is great. Another nice touch is the music. I was overcome with the urge to get off the couch and start dancing. It sucked me in, and it'll do the same to you.”

Output: positive (1) or negative (0)

Let's assume the sentiment is represented by 6 features

$$\mathbf{x} = [x_1, \dots x_6]$$

Thus, we want to estimate

$$P(\text{sentiment} = 1|\mathbf{x})$$

$$P(\text{sentiment} = 0|\mathbf{x})$$

feature	description	value
x_1	Count positive words	
x_2	Count negative words	
x_3	$\begin{cases} 1 & \text{if "no" in text} \\ 0 & \text{otherwise} \end{cases}$	
x_4	Count 1 st and 2 nd pronouns	
x_5	$\begin{cases} 1 & \text{if "!" in text} \\ 0 & \text{otherwise} \end{cases}$	
x_6	log(word count)	

Logistic Regression

Example: Sentiment Classification

“It's **hokey**. There are virtually **no** surprises, and the writing is **second-rate**. So why was it so **enjoyable**? For one thing, the cast is **great**. Another **nice** touch is the music. **I** was overcome with the urge to get off the couch and start dancing. It sucked **me** in, and it'll do the same to **you**.”

$x_3 = 1$ $x_2 = 2$

$x_4 = 3$ $x_5 = 0$ $x_1 = 3$

feature	description	value
x_1	Count positive words	3
x_2	Count negative words	2
x_3	$\begin{cases} 1 & \text{if "no" in text} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	Count 1 st and 2 nd pronouns	3
x_5	$\begin{cases} 1 & \text{if "!" in text} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	$\ln(\text{word count})$	$\ln(56)=4.025$

Logistic Regression

Example: Sentiment Classification

Input: “It's hokey . There are virtually no surprises , and the writing is second-rate . So why was it so enjoyable ? For one thing , the cast is great . Another nice touch is the music . I was overcome with the urge to get off the couch and start dancing . It sucked me in , and it'll do the same to you.”

Output: positive (1) or negative (0)

Let's assume the weights and bias are:

$$\mathbf{w} = [w_1, \dots, w_6]^T = [2.5, -5.0, -1.2, 0.5, 2.0, 0.7]^T$$
$$b = 0.1$$

Thus

$$\begin{aligned} P(\text{sentiment} = 1 | \mathbf{x}) &= \sigma(\mathbf{w}^T \mathbf{x} + b) \\ &= \sigma([2.5, -5.0, -1.2, 0.5, 2.0, 0.7][3, 2, 1, 3, 0, 4.025]^T + 0.1) \\ &= 0.67 \\ P(\text{sentiment} = 0 | \mathbf{x}) &= 1 - \sigma(\mathbf{w}^T \mathbf{x} + b) = 0.33 \end{aligned}$$

feature	description	value
x_1	Count positive words	3
x_2	Count negative words	2
x_3	$\begin{cases} 1 & \text{if "no" in text} \\ 0 & \text{otherwise} \end{cases}$	1
x_4	Count 1 st and 2 nd pronouns	3
x_5	$\begin{cases} 1 & \text{if "!" in text} \\ 0 & \text{otherwise} \end{cases}$	0
x_6	ln(word count)	4.025

Classifier Comparison

Methods	Assumptions on Feat. Dist.	Feat. Normalization	Cost Function	Regularization	Linear Classifier	Prob. View of Prediction	Generative/Discriminative	Parametric/Non-parametric	Overfitting
Logistic Regression	No	Required	BCE (convex)	Additional term	Linear	Yes	Discriminative	Parametric	Not often
K Nearest Neighbors	No	Required	N/A	N/A	Non-linear	N/A	Discriminative	Non-parametric	when k is too small
Decision Trees	No	Not Required	N/A	N/A	Non-linear	N/A	Discriminative	Non-parametric	with large depth
Support Vector Machines	No	Required	Hinge (convex)	C (control robustness)	Linear/Non-linear(kernel)	N/A	Discriminative	Parametric	Not often
Naïve Bayes	Conditional independent	Not Required	N/A	N/A	Non-linear/Linear (Gaussian)	Yes	Generative	Parametric	Not often
Artificial Neural Networks	No	Required	Non-convex	Additional term	Non-linear	Yes	Discriminative/Generative	Parametric	with many layers

Overview

In this Course..

Nearest Neighbor

Naïve Bayes

Logistic Regression

Decision Trees (know it exists but will not study it)

- Overview
- ID3 Algorithm
- Entropy & Information Gain
- Example

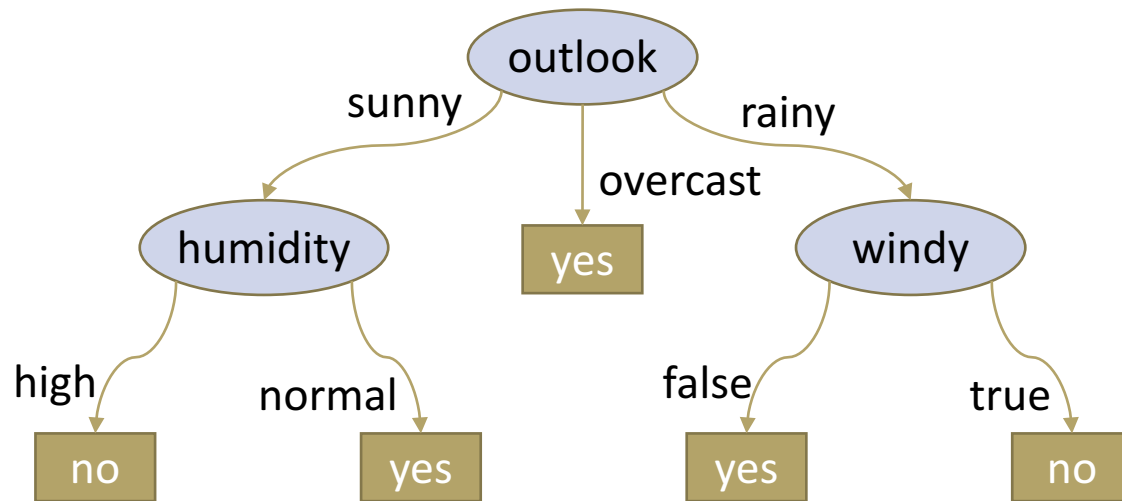
Support Vector Machines

Artificial Neural Networks

Decision Trees

Overview

- Can be used for both Classification and Regression
- The learning process builds classification model in the form of a tree structure that mimics rule-based human reasoning



Final decision tree computed based on data in table

outlook	temp.	humidity	windy	play
sunny	hot	high	false	no
sunny	hot	high	true	no
overcast	hot	high	false	yes
rainy	mild	high	false	yes
rainy	cool	normal	false	yes
rainy	cool	normal	true	no
Overcast	cool	normal	true	yes
sunny	mild	high	false	no
sunny	cool	normal	false	yes
rainy	mild	normal	false	yes
sunny	mild	normal	true	yes
overcast	mild	high	true	yes
overcast	hot	normal	false	yes
rainy	mild	high	true	no

Overview

In this Course..

Nearest Neighbor

Naïve Bayes

Logistic Regression

Decision Trees

Support Vector Machines

- Linear Separators
- Maximum Margin Classification
- Soft Margin Classification
- Multi-class SVM
- Non-linear SVM and kernel tricks

Artificial Neural Networks

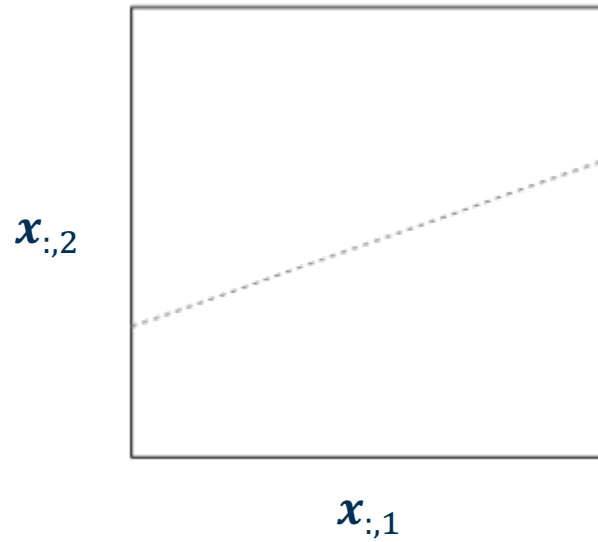
Support Vector Machines

Overview

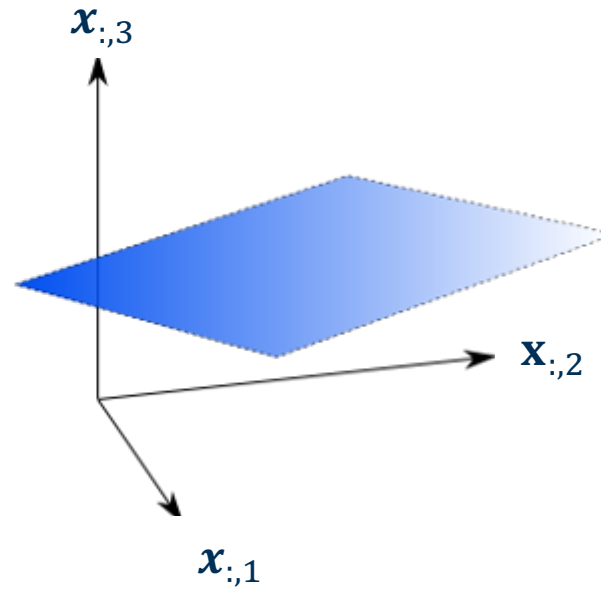
- Linear Separating Hyperplanes
- Maximal Margin Classifier
- Support Vector Classifiers
- Support Vector Machines
- Advantages and Disadvantages

Support Vector Machines

Linear Hyperplanes



One-dimensional hyperplane



Two-dimensional hyperplane

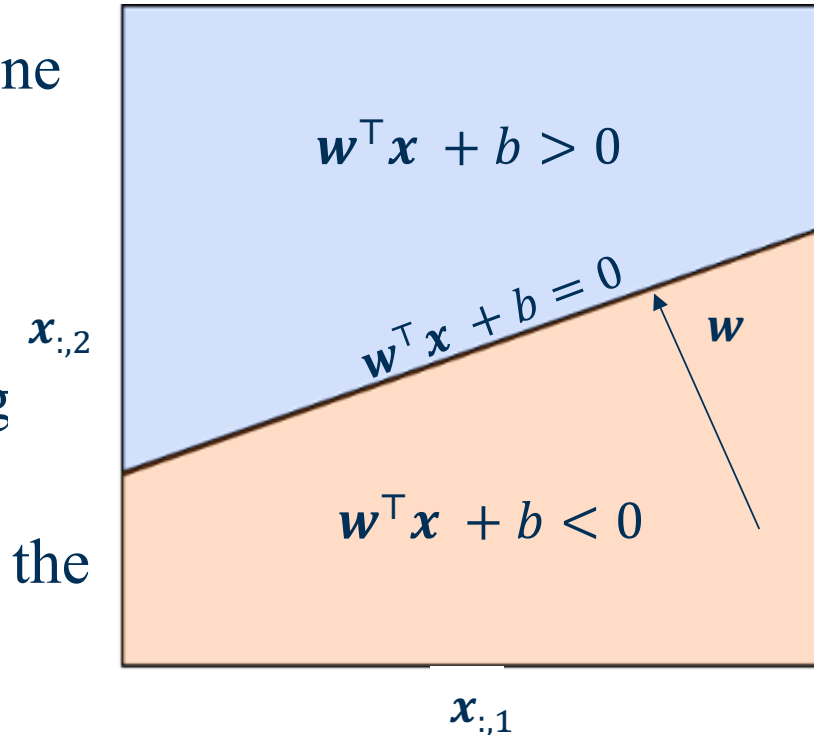
Support Vector Machines

Linear Hyperplanes

- Separating P -dimensional feature space by a hyperplane
- Given a single sample $\mathbf{x} = [x_1, x_2, \dots, x_P] \in \mathbb{R}^P$, we define a hyperplane by the following equation:

$$b + w_1 x_1 + \dots + w_P x_P = 0$$
$$\text{or } \mathbf{w}^\top \mathbf{x} + b = 0$$

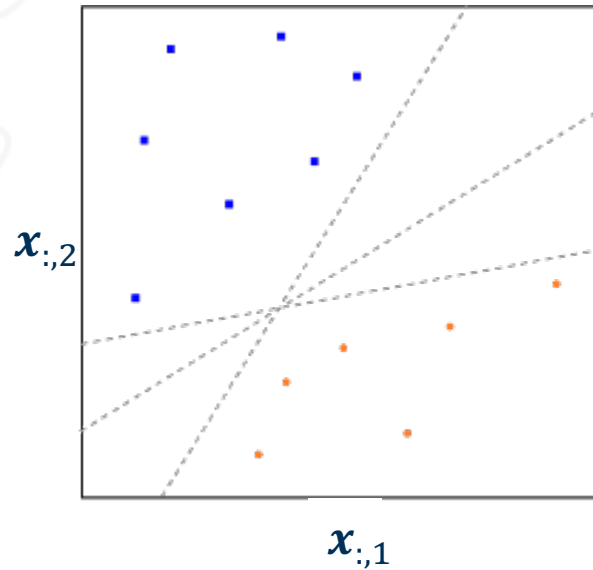
- $b \neq 0$ gives us an affine plane (*i.e.* it does not pass through the origin)
- Note: the vector \mathbf{w} is ANY vector that is perpendicular to the hyperplane. However, for simplicity and without loss of generality, we will assume \mathbf{w} is a unit norm vector.



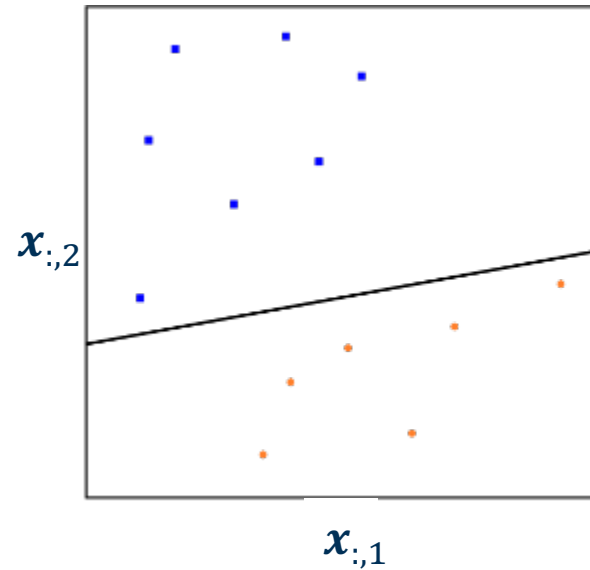
Separation of two-dimensional feature space by a hyperplane

Support Vector Machines

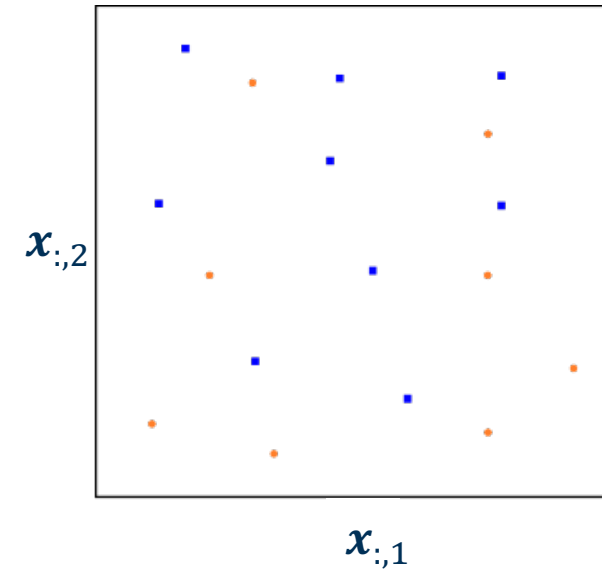
Linear Hyperplanes



Multiple separating hyperplanes



Perfect separation of class data



No possibility of a true separating hyperplane

Support Vector Machines

Maximal Margin Classifiers

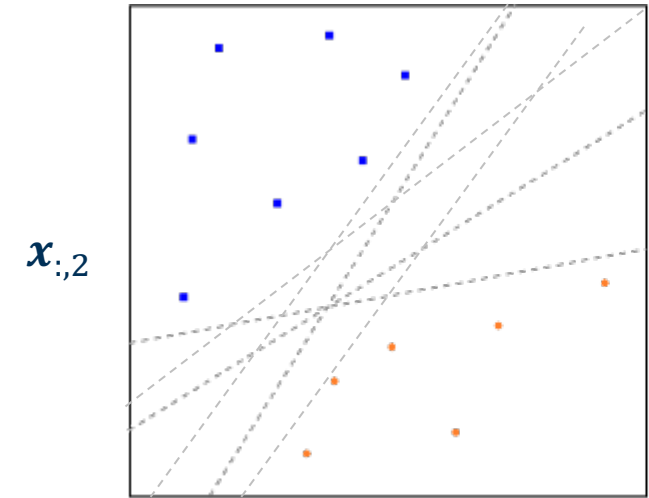
- Our goal is to develop a classifier based on training data that will correctly classify test data using only their feature values.
- We will initially suppose that it is possible, via a means yet to be determined, to construct a hyperplane that separates training data *perfectly* according to their class labels.
- The following diagram is only showing $P = 2$. However, for real examples (like spam email filtering) $P > 10^6$.

$$\mathbf{w}^\top \mathbf{x} + b > 0, \text{ if } y = 1 \text{ (class 1)}$$

and

$$\mathbf{w}^\top \mathbf{x} + b < 0, \text{ if } y = -1 \text{ (class 2)}$$

- However, based on the training data, it's possible to find many hyperplanes that satisfy the above. How do we choose the “**best**” ?



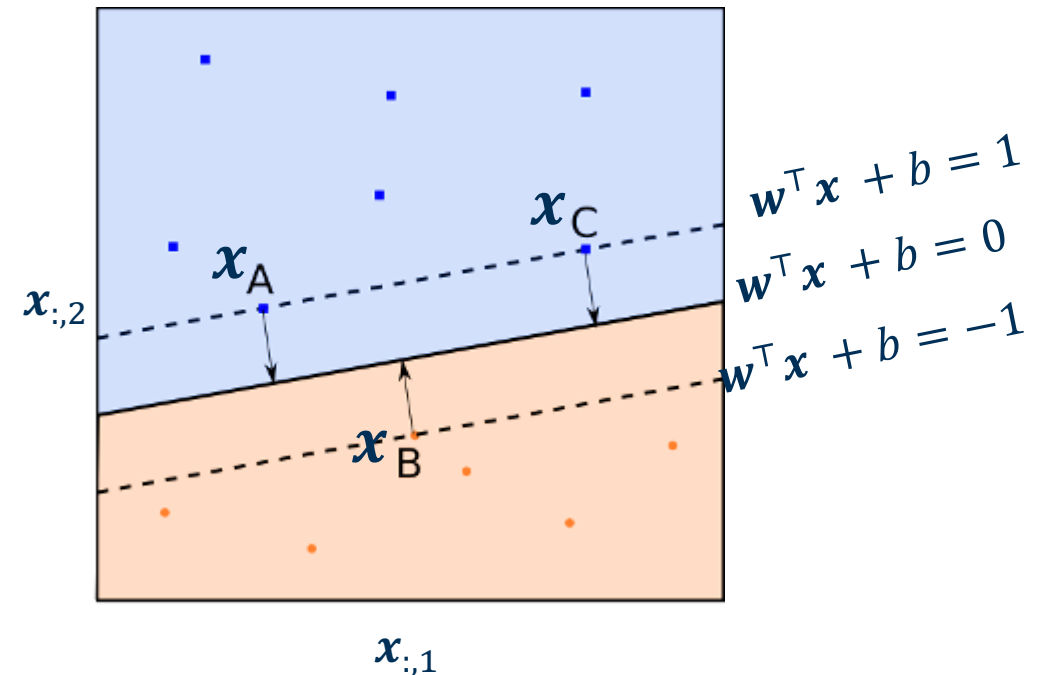
Support Vector Machines

Maximal Margin Classifiers

- The **maximal margin hyperplane** (MMH) is the separating hyperplane that is farthest from any training data.
- The location of the MMH only depends on the training observations that lie directly on the margin boundary (\mathbf{x}_A , \mathbf{x}_B and \mathbf{x}_C)
- (\mathbf{x}_A , \mathbf{x}_B and \mathbf{x}_C) are known as the **support vectors**.
- Thus, the goal is:

finding an algorithm that can produce the best MMH.

Best = Highest width and highest accuracy



Support Vector Machines

Width of Maximal Margin Classifiers

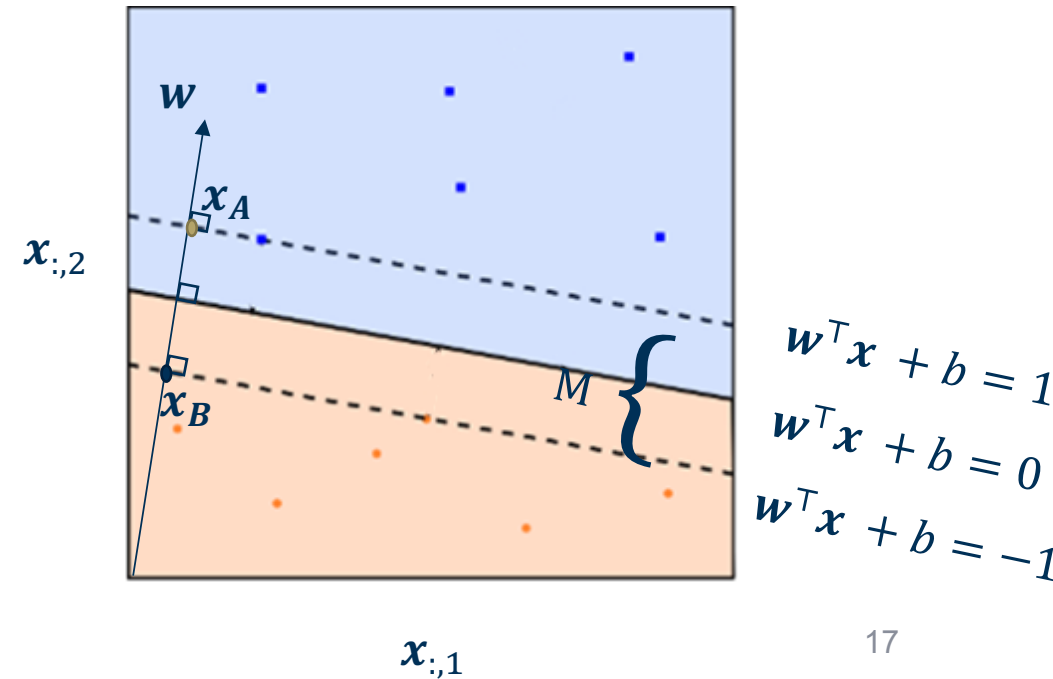
- Let's try to write the width of the margin as a function of something we know.
- Let's change the equation as follows, to ensure no samples are within the margin M
 - $\mathbf{w}^\top \mathbf{x}^+ + b \geq 1$ \mathbf{x}^+ is a class 1 example
 - $\mathbf{w}^\top \mathbf{x}^- + b \leq -1$ \mathbf{x}^- is a class 2 example
- Note the following from the figure:

$$(\mathbf{x}_A - \mathbf{x}_B)^\top \mathbf{w} = 2$$

But the left-hand side is the inner product and can

Be written as: $\|\mathbf{x}_A - \mathbf{x}_B\| \|\mathbf{w}\| \cos(0) = 2$

$$\text{Thus, } M = \frac{2}{\|\mathbf{w}\|} = \|\mathbf{x}_A - \mathbf{x}_B\|$$



Support Vector Machines

Width of Maximal Margin Classifiers

- Let's try to write the width of the margin as a function of something we know.
- Let's change the equation as follows, to ensure no samples are within the margin M
 - $\mathbf{w}^\top \mathbf{x}^+ + b \geq 1$ \mathbf{x}^+ is a class 1 example
 - $\mathbf{w}^\top \mathbf{x}^- + b \leq -1$ \mathbf{x}^- is a class 2 example
- Note the following from the figure:

$$(\mathbf{x}_A - \mathbf{x}_B)^\top \mathbf{w} = 2$$

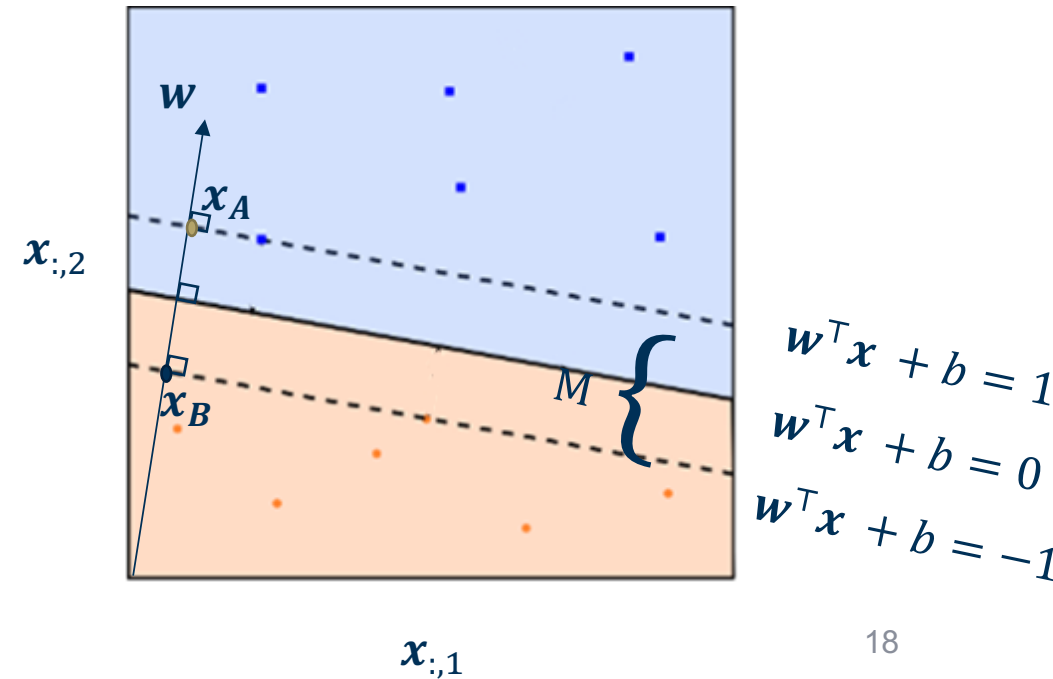
But the left-hand side is the inner product and can

Be written as: $\|\mathbf{x}_A - \mathbf{x}_B\| \|\mathbf{w}\| \cos(0) = 2$

$$\text{Thus, } M = \frac{2}{\|\mathbf{w}\|} = \|\mathbf{x}_A - \mathbf{x}_B\|$$

Therefore, we want to:

- $\max \frac{1}{\|\mathbf{w}\|}$ or $\min \|\mathbf{w}\|$ or $\min \frac{1}{2} \|\mathbf{w}\|^2$

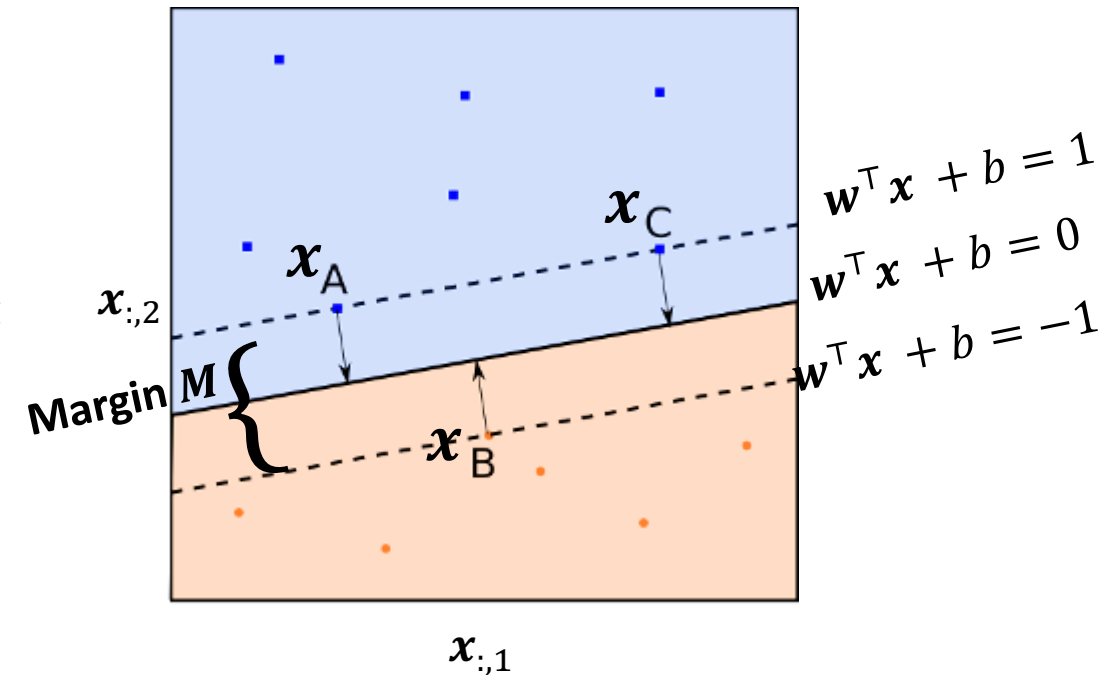


Support Vector Machines

Accuracy of Maximal Margin Classifiers

- Perceptron with a linear boundary $\mathbf{w}^T \mathbf{x} + \mathbf{b} = 0$:
$$\begin{cases} \mathbf{w}^T \mathbf{x} + \mathbf{b} > 0, & \text{if } y = +1 \\ \mathbf{w}^T \mathbf{x} + \mathbf{b} < 0, & \text{if } y = -1 \end{cases}$$
- Suppose that all samples could be separated further by $\mathbf{w}^T \mathbf{x} + \mathbf{b} = \pm 1$, the above could be written as following:
$$\begin{cases} \mathbf{w}^T \mathbf{x} + \mathbf{b} \geq 1, & \text{if } y = +1 \\ \mathbf{w}^T \mathbf{x} + \mathbf{b} \leq -1, & \text{if } y = -1 \end{cases}$$
- The above formula could be written as:
$$y(\mathbf{w}^T \mathbf{x} + \mathbf{b}) - 1 \geq 0$$
- Or equivalently as **Margin-Perceptron cost**,
$$L(\mathbf{x}_i, y_i, \boldsymbol{\theta}) = \max(0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + \mathbf{b})) = 0$$

This additional 1 prevents the trivial zero solution with the original perceptron



Support Vector Machines

Accuracy of Maximal Margin Classifiers

- We want our classifier to be correct. Hence, we need to satisfy:
 - $\mathbf{w}^\top \mathbf{x}^+ + b \geq 1$
 - $\mathbf{w}^\top \mathbf{x}^- + b \leq -1$
- Let's multiply both inequalities by y_i .
- Since $y_i = 1$ for \mathbf{x}_i^+ and $y_i = -1$ for \mathbf{x}_i^- , we get
 - $y_i(\mathbf{w}^\top \mathbf{x}_i^+ + b) \geq y_i \rightarrow y_i(\mathbf{w}^\top \mathbf{x}_i^+ + b) \geq 1$
 - $y_i(\mathbf{w}^\top \mathbf{x}_i^- + b) \geq -y_i \rightarrow y_i(\mathbf{w}^\top \mathbf{x}_i^- + b) \geq 1$
- So, the two inequalities become identical, and our optimization objectives become:

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 \\ & \text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0 \end{aligned}$$

Support Vector Machines

Summary

- $y(\mathbf{w}^\top \mathbf{x} + b) - 1 \geq 0$
- The optimization objectives become:

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0$$

- Or equivalently,

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{Subject to } \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b)) = 0$$

Margin-Perceptron cost

Support Vector Machines

Optimization

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0$$

- Since we have a constrained optimization objective, we use Lagrange multiplier to find our cost function.

$$L = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1]$$

- We can solve this by taking the derivative and equating it to 0

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0 \rightarrow \mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^N \alpha_i y_i = 0 \rightarrow \sum_{i=1}^N \alpha_i y_i = 0$$

- Next, let's plug those values back in the cost function

Support Vector Machines

Optimization

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$
$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$\begin{aligned} L &= \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^\top \mathbf{x}_i + b) - 1] \\ &= \frac{1}{2} \left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \right)^\top \left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right) - \sum_{i=1}^N \alpha_i \left[y_i \left(\left(\sum_{j=1}^N \alpha_j y_j \mathbf{x}_j \right)^\top \mathbf{x}_i + b \right) - 1 \right] \\ &= \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j - b \sum_{i=1}^N \alpha_i y_i + \sum_{i=1}^N \alpha_i \\ L &= -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i \end{aligned}$$

Note that the loss (cost) function is only a sum of the dot products of pairs of samples ($\mathbf{x}_i^\top \mathbf{x}_j$)

Support Vector Machines

Optimization

- **Final Objective Function:**

$$\text{maximize}(L(\alpha_i) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_{i=1}^N \alpha_i)$$

Subject to: $\sum_{i=1}^N \alpha_i y_i = 0$ and $\alpha_i \geq 0$

Note that loss (cost) function is only a sum of the dot products of pairs of samples ($\mathbf{x}_i^T \mathbf{x}_j$)

- Can intuitively think of dot product as measure of *similarity*.
- **Case 1:** If two features x_i, x_j are completely dissimilar, their dot product is 0, and they don't contribute to L.
- **Case 2:** If two, features x_i, x_j are completely alike, their dot product is 1.
 - **Subcase 1:** x_i and x_j predict the same output value y_i (either +1 or -1). Then $y_i y_j x_i^T x_j$ is always 1, and the value of $\alpha_i \alpha_j y_i y_j x_i^T x_j$ will be positive. But this would decrease the value of L. So, the objective function downgrades similar feature vectors that make the same prediction.
 - **Subcase 2:** x_i and x_j make opposite predictions about the output value y_i (i.e., one is +1, the other -1), but are otherwise very closely similar: then the product $\alpha_i \alpha_j y_i y_j x_i^T x_j$ is negative and we are subtracting it, so this adds to the sum, maximizing it. This is precisely the examples we are looking for: the critical ones that tell the two classes apart.

Support Vector Machines

Optimization

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$

- This optimization problem can be solved using any solver to find the values of α 's and hence, the values of w_i and b
- $\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$
- After finding the optimal weights, when we get a new unseen sample (\mathbf{u}) to classify, we can simply compute:

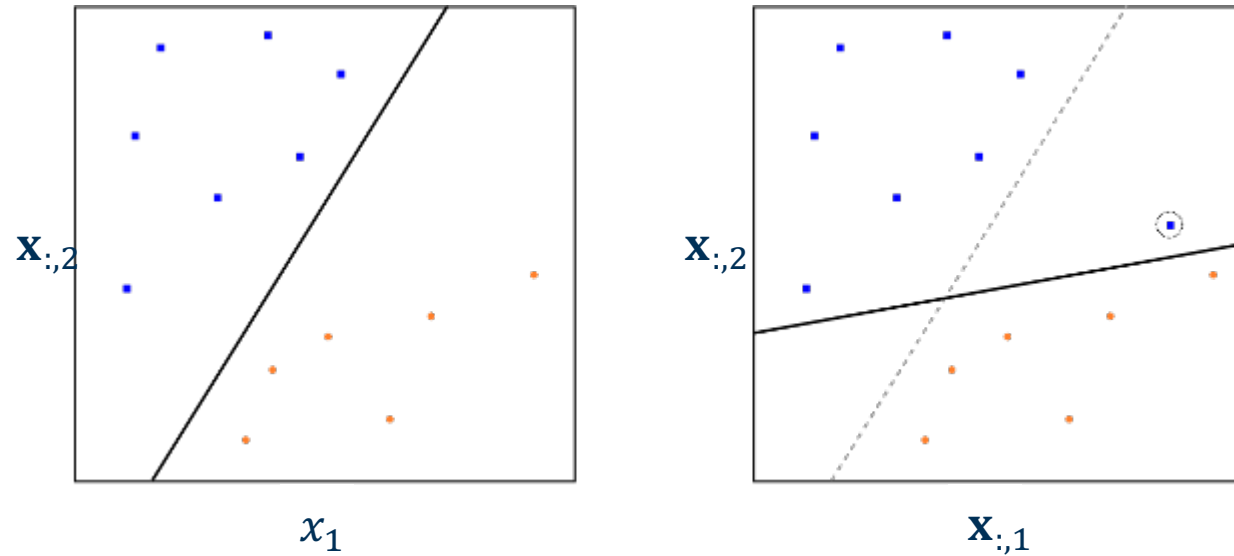
$$\mathbf{w}^\top \mathbf{u} + b$$
$$(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i) \cdot \mathbf{u} + b$$

- If $\mathbf{w}^\top \mathbf{u} + b > 0$, then \mathbf{u} is classified as class 1 (+1)
 - If $\mathbf{w}^\top \mathbf{u} + b < 0$, then \mathbf{u} is classified as class 2 (-1)
- Or simply, class = $\text{sgn}(\mathbf{w}^\top \mathbf{u} + b)$

Support Vector Machines

Classifiers

- Addition of a single point dramatically changes the MMH line



Support Vector Machines

Classifiers

- Clearly, the case of perfect separability is an ideal one. Most "real world" datasets will not have such perfect separability via a linear hyperplane
- However, if there is no separability then we are unable to construct a MMC by the optimization procedure above. So, how do we create a form of separating hyperplane?
- Essentially, we must relax the requirement that a separating hyperplane will perfectly separate every training observation on the correct side of the line (*i.e.* guarantee that it is associated with its true class label), using what is called a **soft margin**. This motivates the concept of a **support vector classifier (SVC)**.

Support Vector Machines

Soft Margin Classifiers

- Maximum Margin Classifiers (MMC) can be extremely sensitive to the addition of new training observations.
- We could consider a classifier based on a separating hyperplane that:
 - Does have a greater robustness to the addition of *new* individual observations
 - Has a better classification on *most* of the training data.
- This comes at the expense of some misclassification of a few training data.

$$\begin{aligned} & \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \\ & \text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) \geq 1 - \xi_i \end{aligned}$$

Recall

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2$$

$$\text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) - 1 \geq 0$$

Support Vector Machines

Soft Margin Classifiers

- Maximum Margin Classifiers (MMC) can be extremely sensitive to the addition of new training observations.
- We could consider a soft margin classifier based on a separating hyperplane that:
 - Does have a greater robustness to the addition of *new* individual observations
 - Has a better classification on *most* of the training data.
- This comes at the expense of some misclassification of a few training data.

$$\text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i \quad \text{regularization parameter}$$

$$\begin{aligned} \text{Subject to } y_i(\mathbf{w}^\top \mathbf{x}_i + b) &\geq 1 - \xi_i \\ \xi_i &\geq 0 \end{aligned}$$

- The above objective could be written as following:

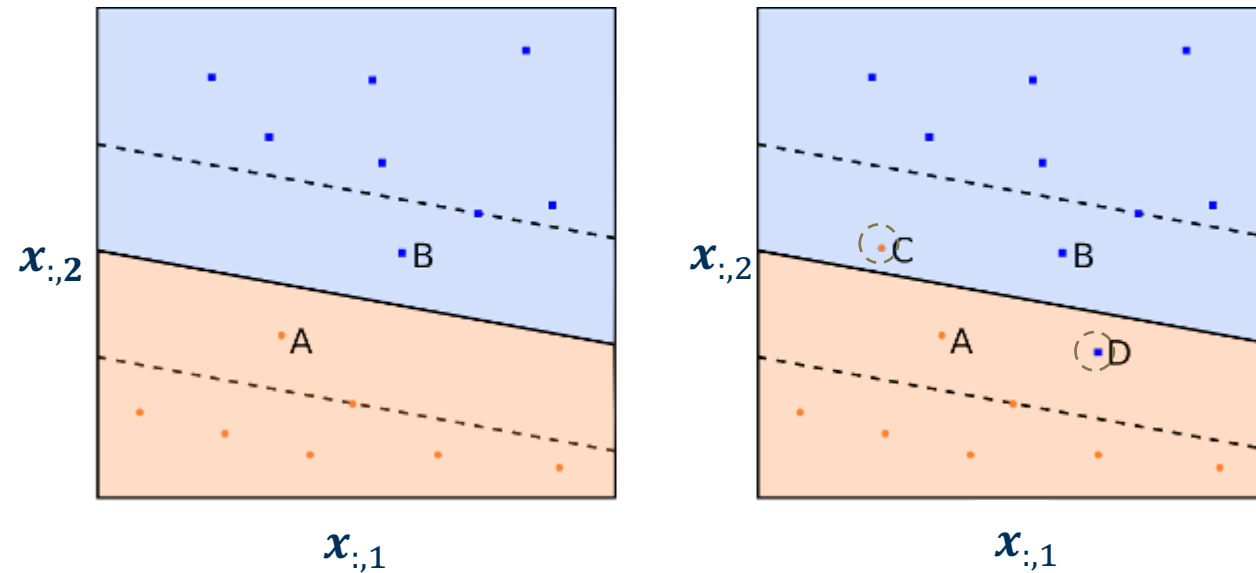
$$\text{Margin-Perceptron cost} \quad \text{minimize } \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^N \max(0, 1 - y_i(\mathbf{w}^\top \mathbf{x}_i + b))$$

L₂ norm regularization points to $\frac{1}{2} \|\mathbf{w}\|^2$

Support Vector Machines

Soft Margin Classifiers

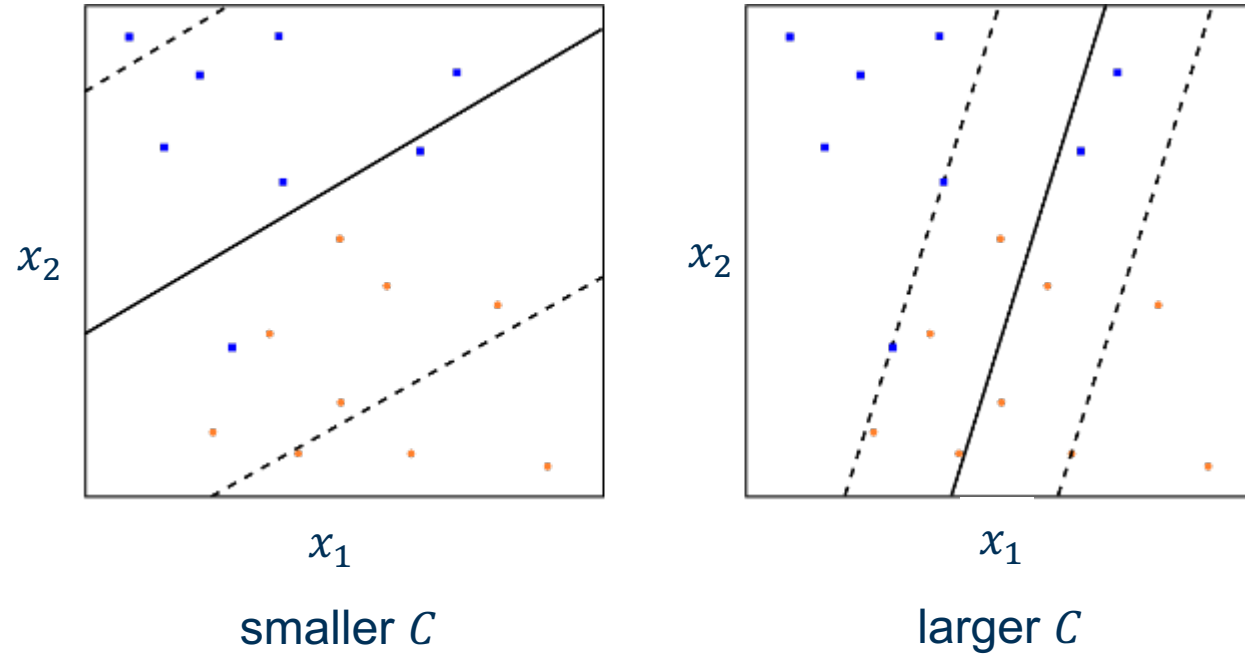
- A *Soft Margin Classifier* allows some observations to be on the incorrect side of the margin (or hyperplane), hence it provides a "soft" separation.



Support Vector Machines

Soft Margin Classifiers

- Different values of the tuning parameter C

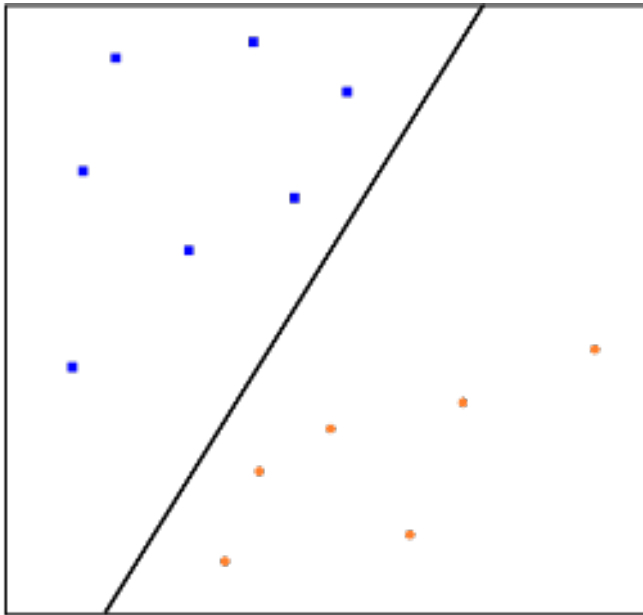


Larger C = Higher penalty for misclassification, results in narrower margin

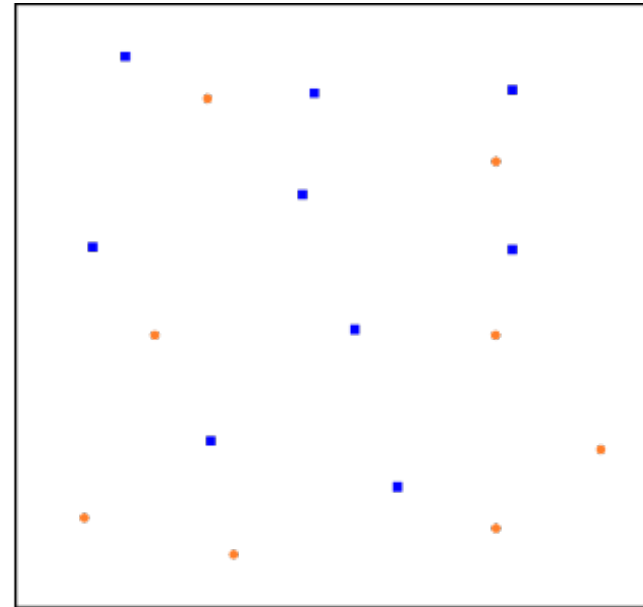
Support Vector Machines

Non-Linear Soft Margin Classifiers

- What if our data samples are not linearly separable?



Linearly Separable

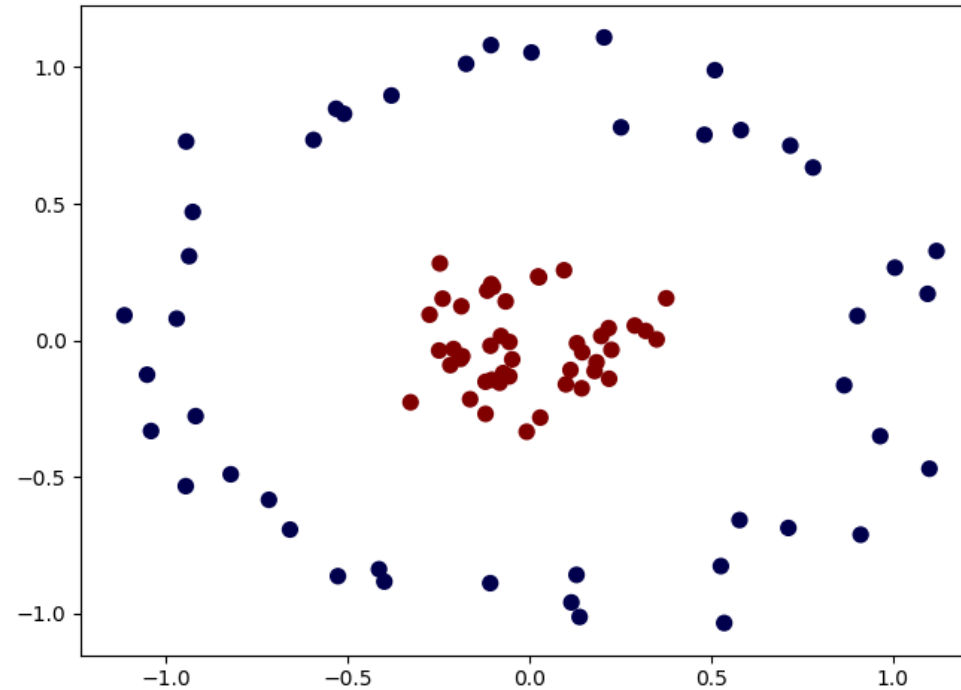


NOT Linearly Separable

Support Vector Machines

Non-Linear Soft Margin Classifiers

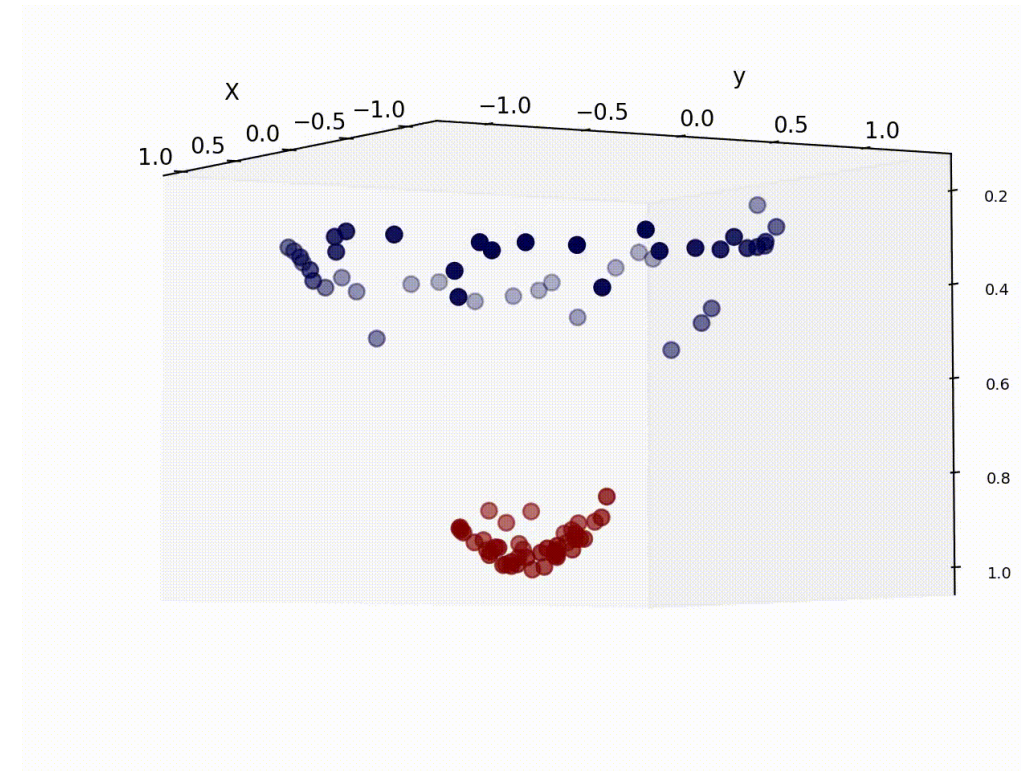
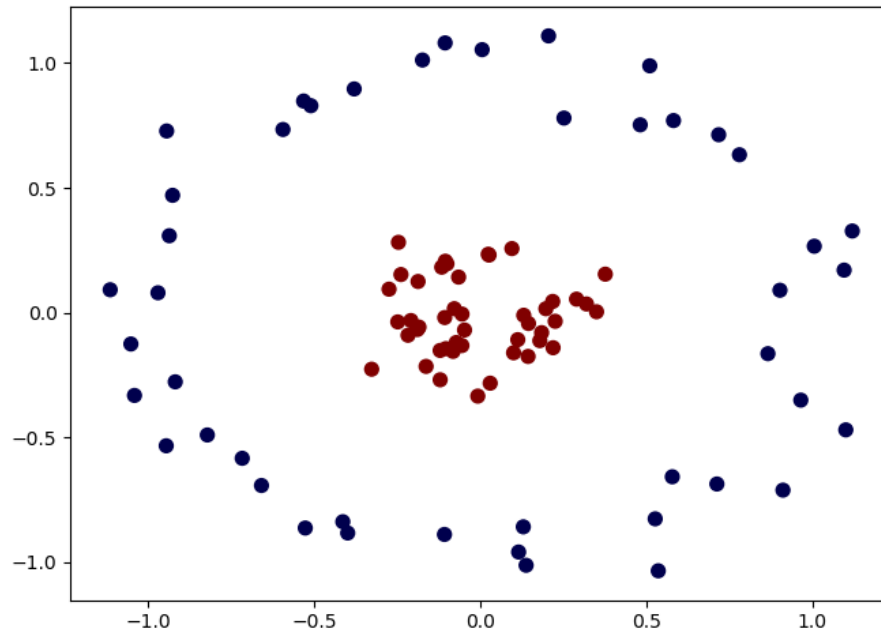
- There is no single hyperplane that can correctly separate the classes



Support Vector Machines

Non-Linear Soft Margin Classifiers

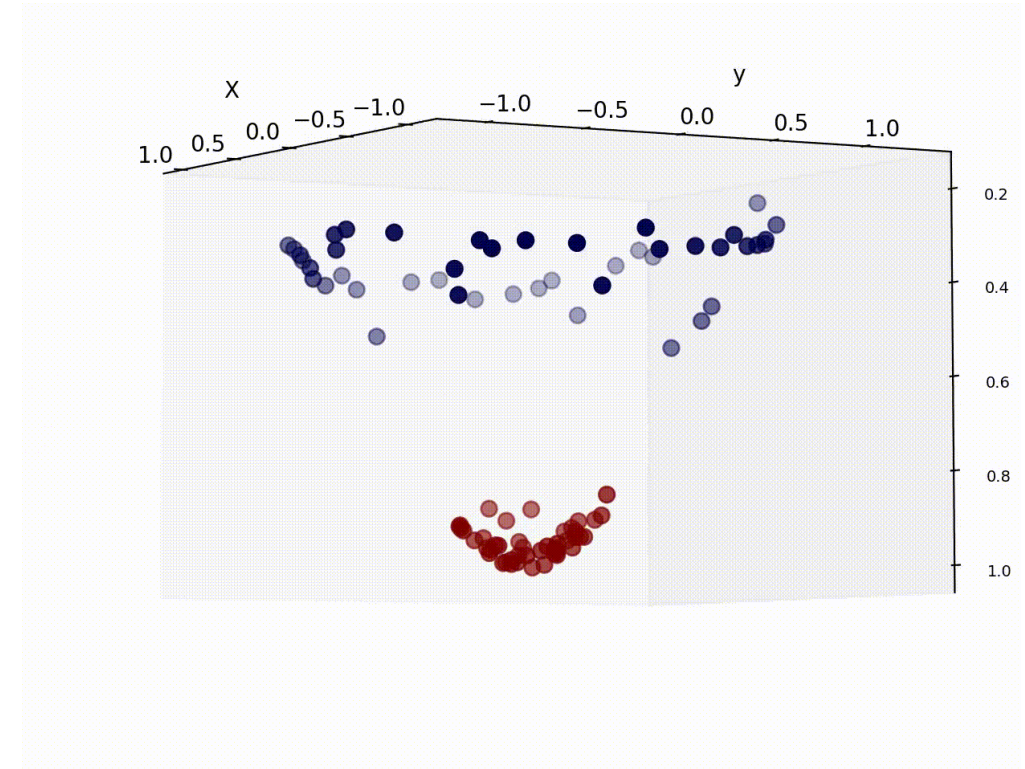
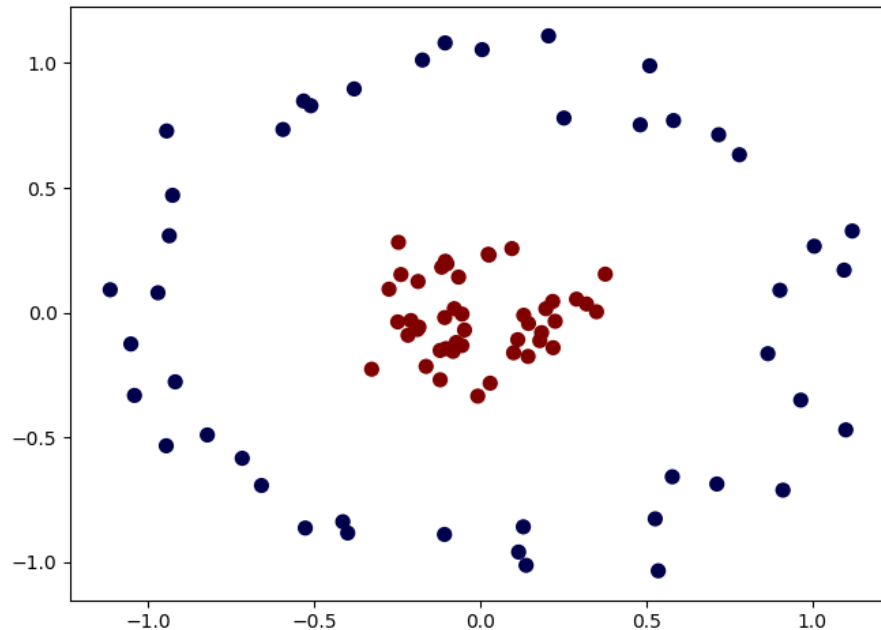
- But we can transform to another subspace where the data is linearly separable.



Support Vector Machines

Kernel Trick

- We can use kernel tricks $K(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j)$ to transform linearly non-separable data to another subspace where the data is linearly separable.



Support Vector Machines

Kernel Trick

- Recall the loss function of the MMC

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \mathbf{x}_i^\top \mathbf{x}_j + \sum_{i=1}^N \alpha_i$$

- Let us define the transform $\phi(\mathbf{x}_i)$

$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \phi(\mathbf{x}_i)^\top \phi(\mathbf{x}_j) + \sum_{i=1}^N \alpha_i$$

- Let us define the kernel:

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

- Now, we don't need to know the transform ϕ , we can simply use a kernel K which is much easier.

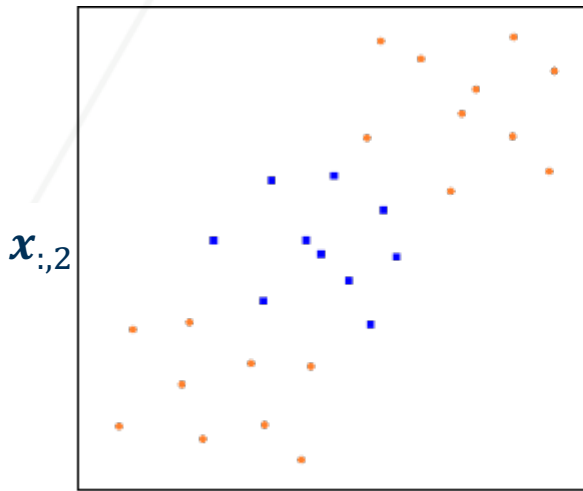
$$L = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j K(\mathbf{x}_i, \mathbf{x}_j) + \sum_{i=1}^N \alpha_i$$

Support Vector Machines

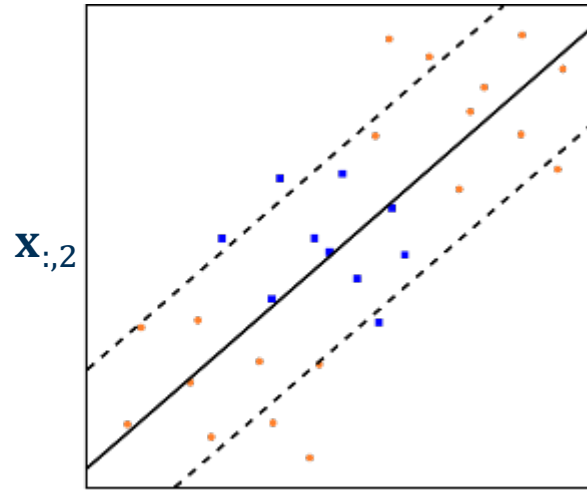
Kernel Trick

- Example kernels

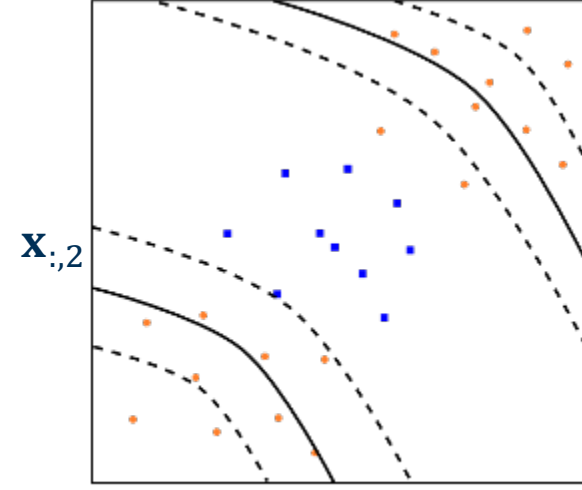
- $K(\mathbf{x}_i, \mathbf{x}_j) = (\mathbf{x}_i^\top \mathbf{x}_j + 1)^p$; polynomial kernel
- $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{1}{2\sigma^2}(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)}$; Gaussian kernel
- $K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\gamma(\mathbf{x}_i - \mathbf{x}_j)^\top (\mathbf{x}_i - \mathbf{x}_j)}$; Radian Basis Function (RBF)



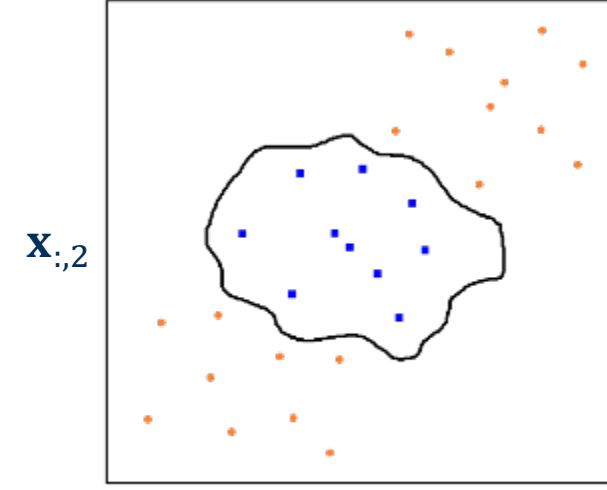
$\mathbf{x}_{:,1}$
data



$\mathbf{x}_{:,1}$
Linear SVM



$\mathbf{x}_{:,1}$
Polynomial Kernel

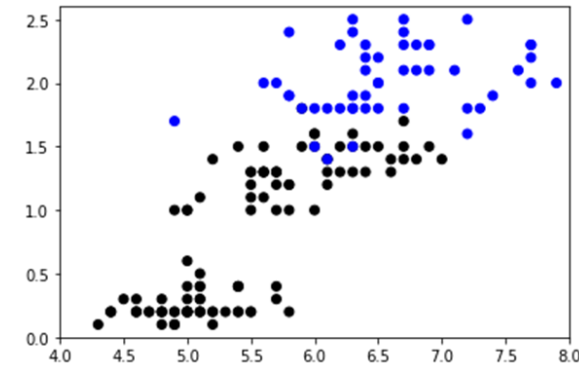
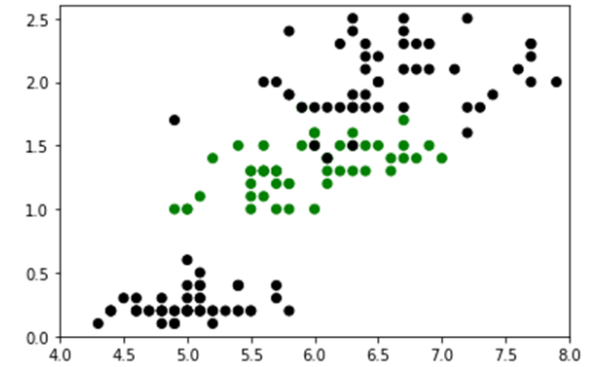
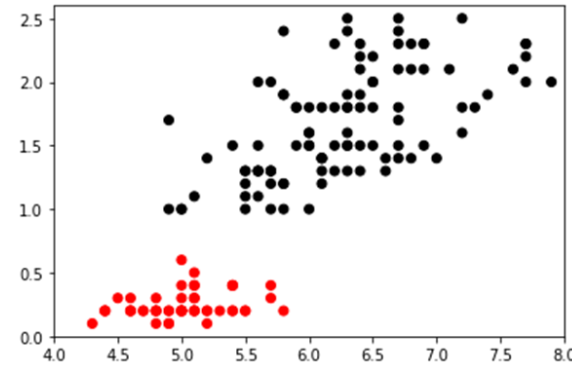
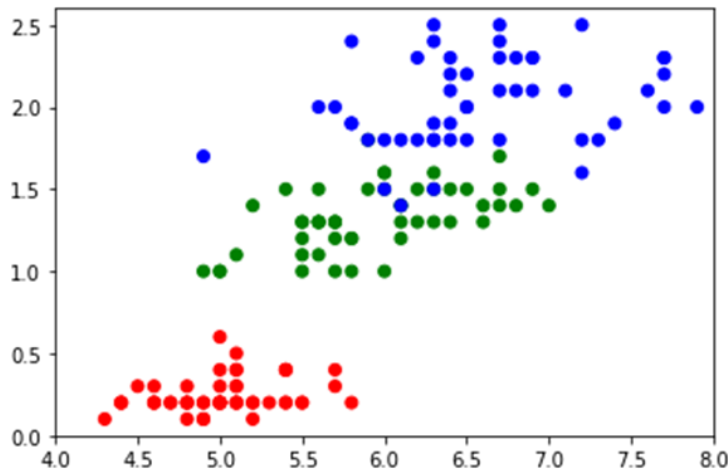


$\mathbf{x}_{:,1}$
RBF Kernel

Support Vector Machines

Multi-Class SVM

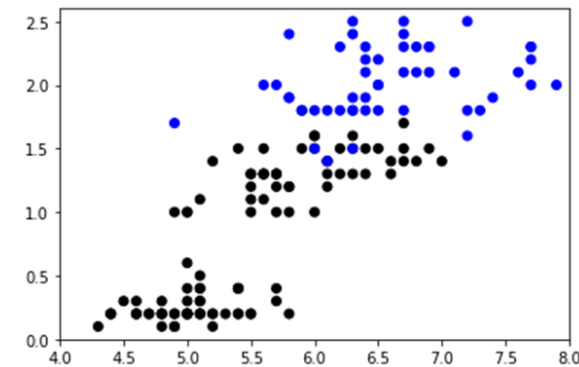
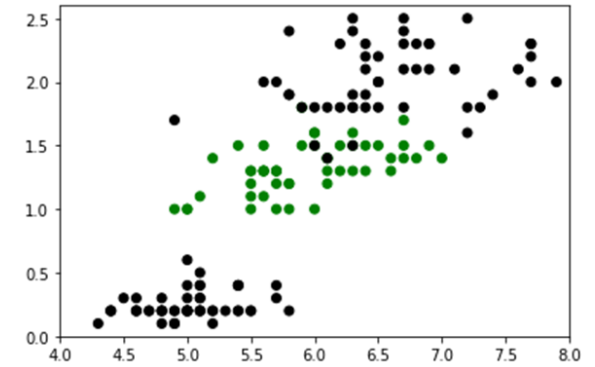
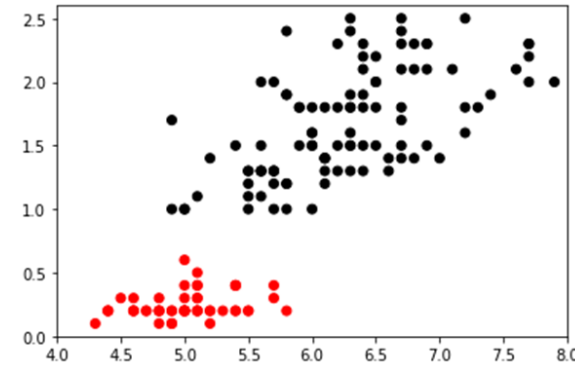
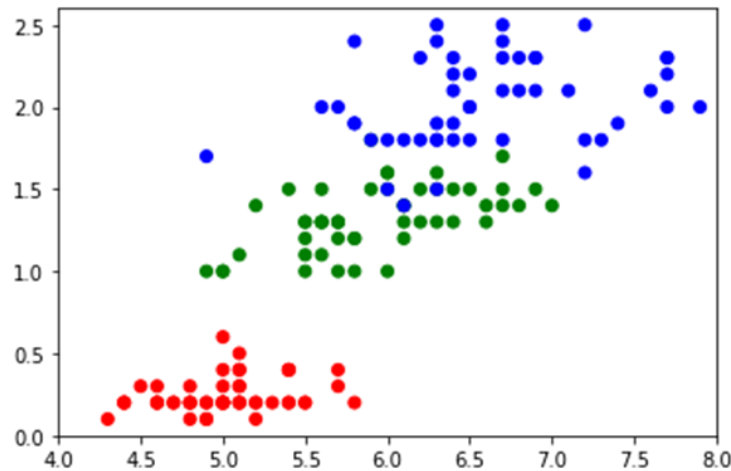
- SVMs are by default binary. However, they can be used to perform multi-class classification.
- Strategy 1:
 - One-Versus-All
 - K classes $\rightarrow K$ classifiers
 - Class with largest confidence is chosen



Support Vector Machines

Multi-Class SVM

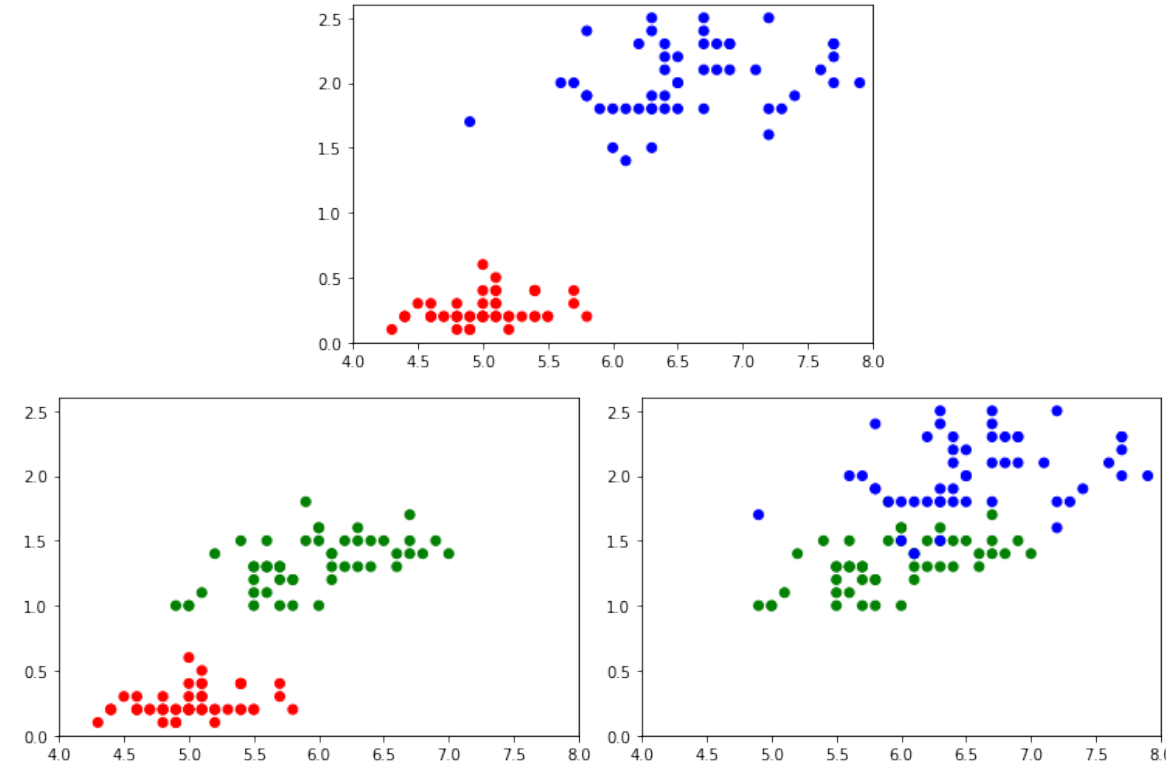
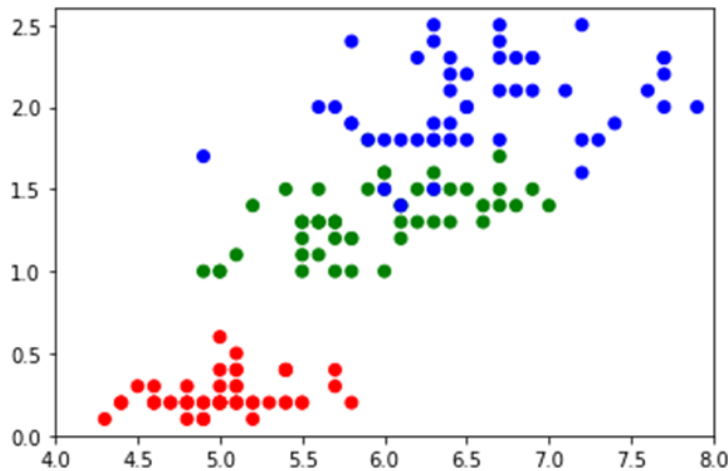
- One-Versus-All
- K classes $\rightarrow K$ classifiers
- Issue:
 - Unbalanced data



Support Vector Machines

Multi-Class SVM

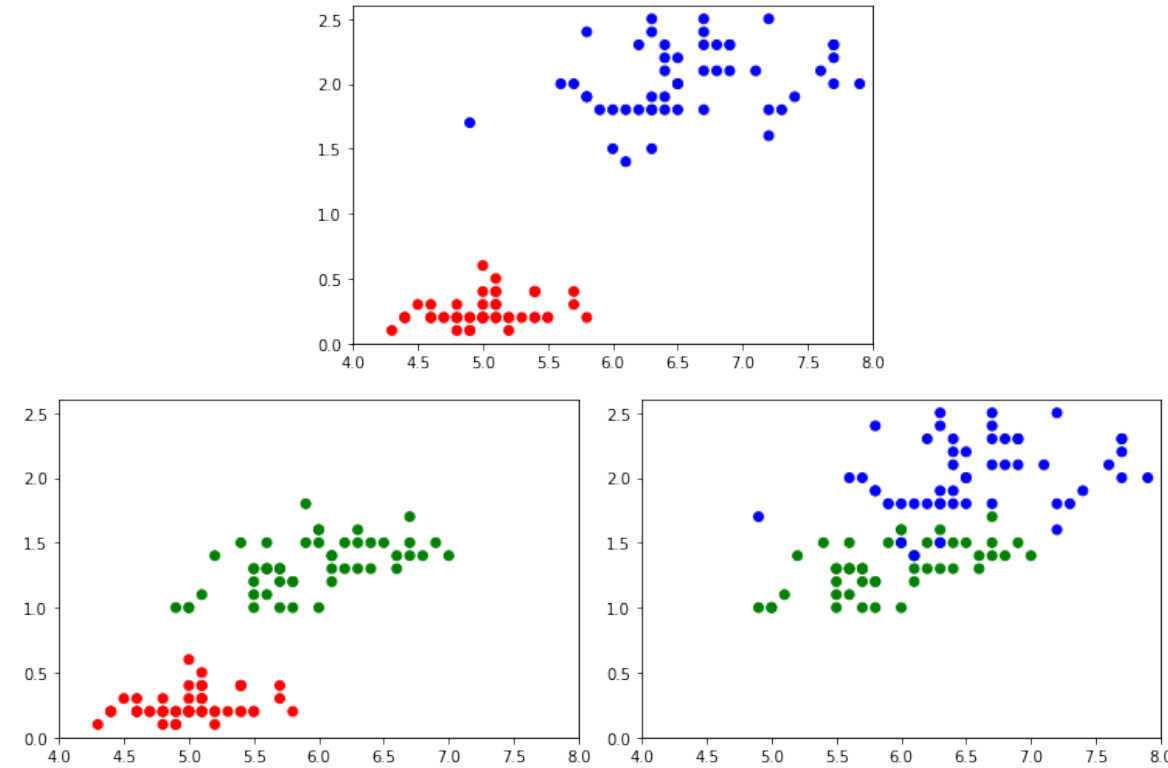
- Strategy 2:
 - One-Versus-One
 - K classes $\rightarrow \frac{K(K-1)}{2}$ classifiers
 - Then class is chosen using majority vote



Support Vector Machines

Multi-Class SVM

- One-Versus-One
- K classes $\rightarrow \frac{K(K-1)}{2}$ classifiers
- Issue:
 - Computationally expensive



Support Vector Machines

Pros and Cons

- Advantages
 - High-Dimensionality
 - Memory Efficiency
 - Versatility
- Disadvantages
 - Feature dimension $P \gg$ the number of training samples N
 - Non-probabilistic

Classifier Comparison

Methods	Assumptions on Feat. Dist.	Feat. Normalization	Cost Function	Regularization	Linear Classifier	Prob. View of Prediction	Generative/Discriminative	Parametric/Non-parametric	Overfitting
Logistic Regression	No	Required	BCE (convex)	Additional term	Linear	Yes	Discriminative	Parametric	Not often
K Nearest Neighbors	No	Required	N/A	N/A	Non-linear	N/A	Discriminative	Non-parametric	when k is too small
Decision Trees	No	Not Required	N/A	N/A	Non-linear	N/A	Discriminative	Non-parametric	with large depth
Support Vector Machines	No	Required	Hinge (convex)	C (control robustness)	Linear/Non-linear(kernel)	N/A	Discriminative	Parametric	Not often
Naïve Bayes	Conditional independent	Not Required	N/A	N/A	Non-linear/Linear (Gaussian)	Yes	Generative	Parametric	Not often
Artificial Neural Networks	No	Required	Non-convex	Additional term	Non-linear	Yes	Discriminative/Generative	Parametric	with many layers

Appendix A: Notations

- x_i : a single feature
- \mathbf{x}_i : feature vector (data sample)
- \mathbf{X} : matrix of feature vectors (dataset)
- N : number of data samples
- m : degree of polynomial
- P : number of features in a feature vector
- θ_i : a single model coefficient (parameter)
- $\boldsymbol{\theta}$: coefficient vector
- ε : error margin
- α : learning rate
- γ : bias factor
- Bold letter/symbol: vector
- Bold capital letters/symbol: matrix