

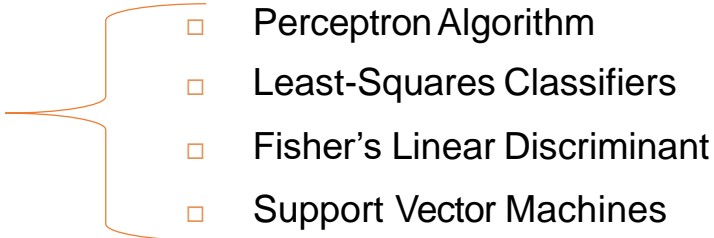
CS 4602

Introduction to Machine Learning

Linear Classifiers

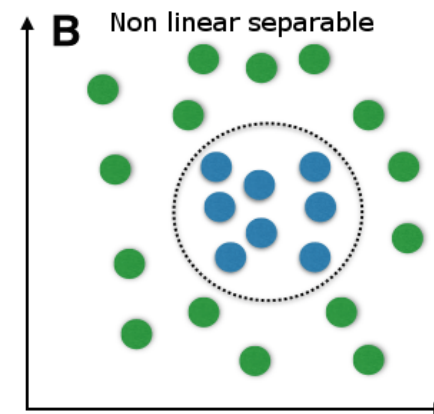
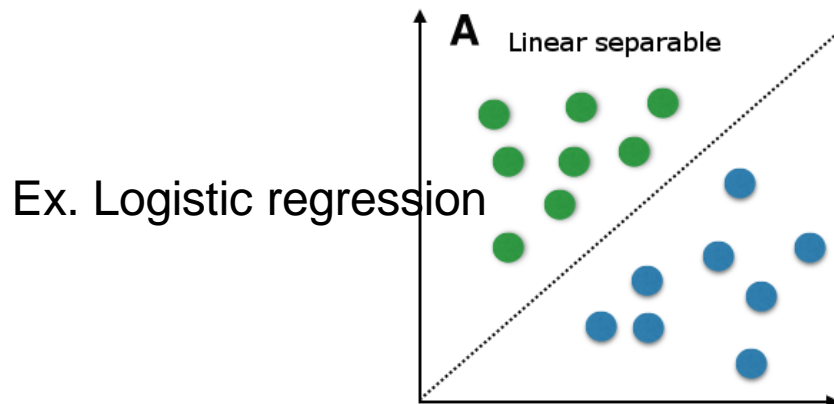
Instructor: Po-Chih Kuo

Roadmap

- Introduction and Basic Concepts
 - Regression
 - Bayesian Classifiers
 - Decision Trees
 - KNN
 - Linear Classifier
 - Neural Networks
 - Deep learning
 - Convolutional Neural Networks
 - RNN/Transformer
 - Reinforcement Learning
 - Model Selection and Evaluation
 - Clustering
 - Data Exploration & Dimensionality reduction
- 
- Perceptron Algorithm
 - Least-Squares Classifiers
 - Fisher's Linear Discriminant
 - Support Vector Machines

Recall

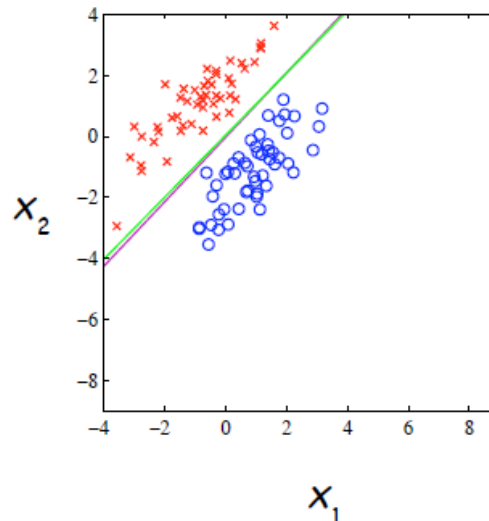
- In regression, we are modeling the relationship between a continuous input variable \mathbf{x} and a continuous target variable y .
- In classification, the input variable \mathbf{x} may still be continuous, but the target variable is discrete.



Ex. KNN

Linear Models for Classification

- Linear models for classification separate input vectors into classes using linear (hyperplane) *decision boundaries*.
- Example:
 - 2D Input vector $\mathbf{x} = (x_1, x_2)$
 - Two discrete classes C1 (red) and C2 (blue)



Two Class Discriminant Function

$$y(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

$y(\mathbf{x}) \geq 0$ \mathbf{x} assigned to C_1

$y(\mathbf{x}) < 0$ \mathbf{x} assigned to C_2

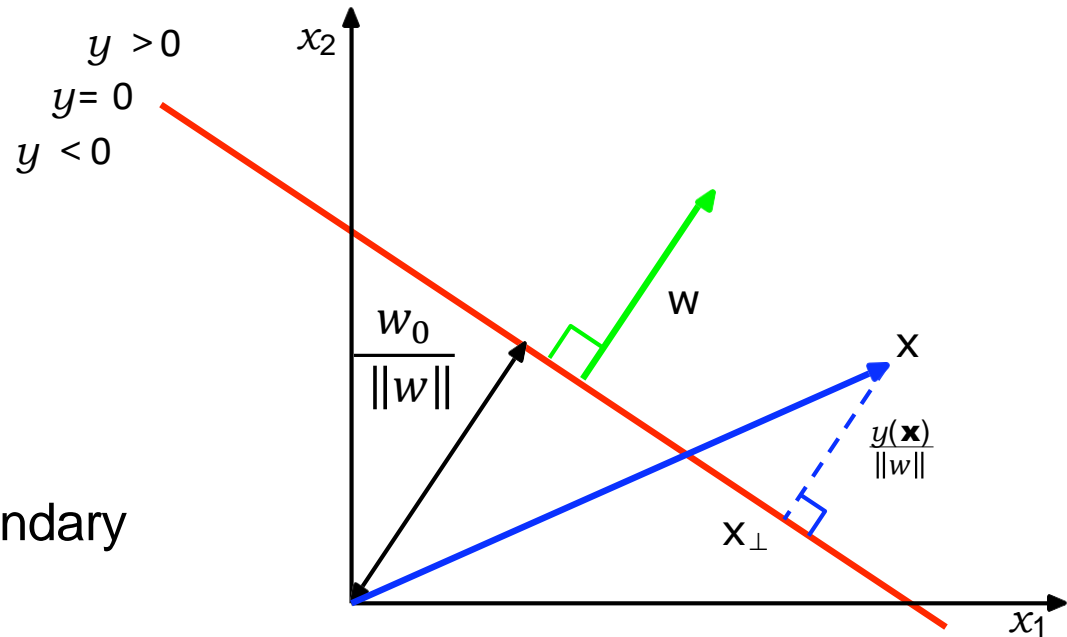
$y(\mathbf{x}) = 0$ defines the decision boundary

Let

$$\mathbf{w} = [w_1 \dots w_M]^T \Rightarrow [w_0 \ w_1 \dots w_M]^T$$

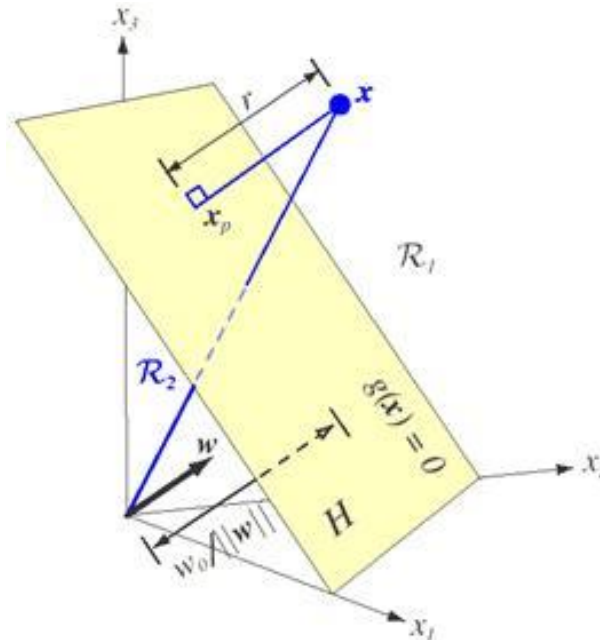
$$\mathbf{x} = [x_1 \dots x_M]^T \Rightarrow [1 \ x_1 \dots x_M]^T$$

$$y(\mathbf{X}) = \mathbf{W}^T \mathbf{X}$$



Hyperplane

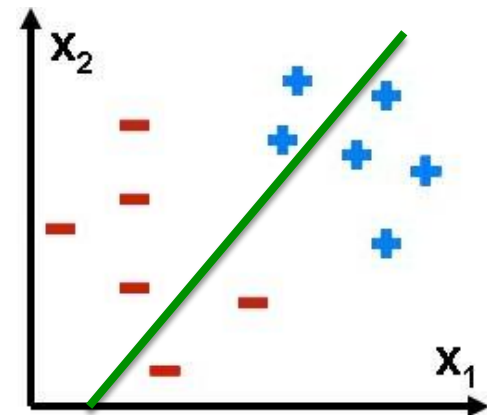
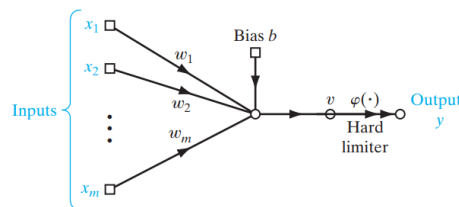
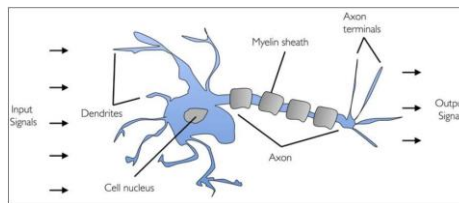
- Decision boundary $g(x)$ is a hyperplane
- A hyperplane is
 - a point in 1D
 - a line in 2D
 - a plane in 3D



- **Perceptron Algorithm**
- Least-Squares Classifiers
- Fisher's Linear Discriminant
- Support Vector Machines

The Perceptron Algorithm

- The perceptron algorithm was invented by Frank Rosenblatt (1962).
- The strategy is to start with a random **guess** at the weights \mathbf{w} , and to then iteratively change the weights to move the hyperplane in a direction that lowers the classification error.



The Perceptron Algorithm

- As we change the weights continuously, the classification error changes in discontinuous, piecewise constant fashion.
- Thus we cannot use the classification error as our objective function to minimize.
- What would be a better objective function?

The Perceptron Algorithm

□ we seek \mathbf{w} such that

- $\mathbf{w}^T \mathbf{x} \geq 0$ when $t = +1$
- $\mathbf{w}^T \mathbf{x} < 0$ when $t = -1$

□ In other words, we would like

$$\mathbf{w}^T \mathbf{x}_n t_n \geq 0, \forall n$$

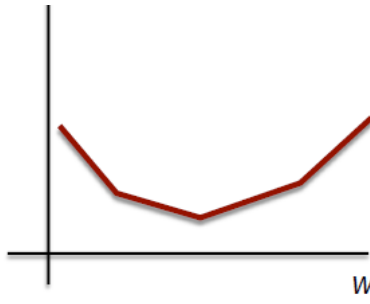
□ Thus we seek to minimize

$$E(\mathbf{w}) = - \sum_{n \in M} \mathbf{w}^T \mathbf{x}_n t_n$$

M is the set of misclassified inputs

The Perceptron Algorithm

- $E(\mathbf{w})$ is always non-negative.
- $E(\mathbf{w})$ is continuous and piecewise linear, and thus easier to minimize.



- Again, we can use gradient descent!

$$\mathbf{w}^{\theta+1} = \mathbf{w}^{\theta} - \eta \nabla E(\mathbf{w}) = \mathbf{w}^{\theta} + \eta \sum_{n \in M} \mathbf{x}_n t_n$$

If an input from $C_1(t = +1)$ is misclassified, we need to make its projection on \mathbf{w} more positive.

Not Monotonic

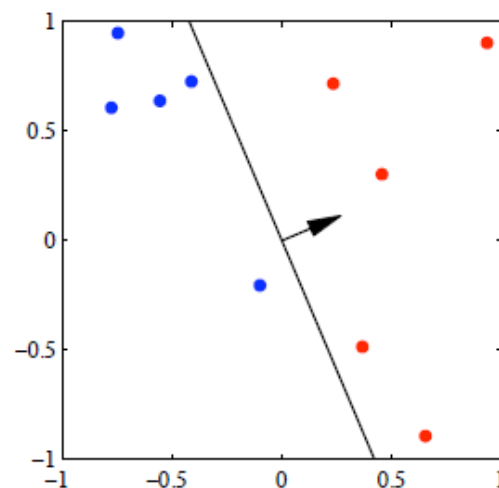
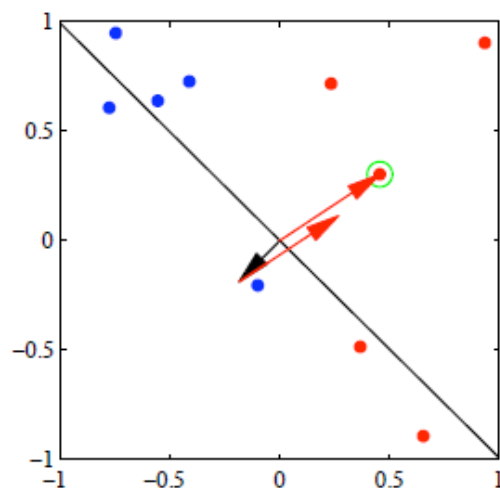
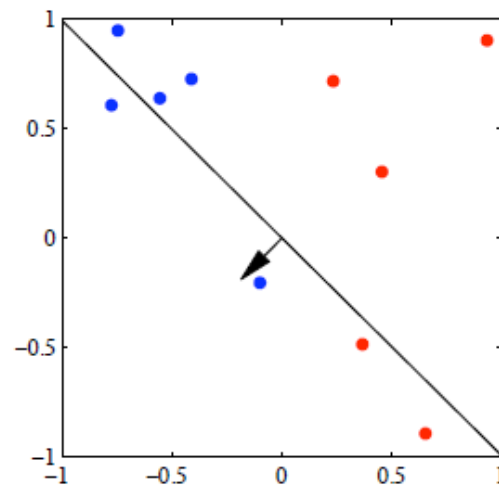
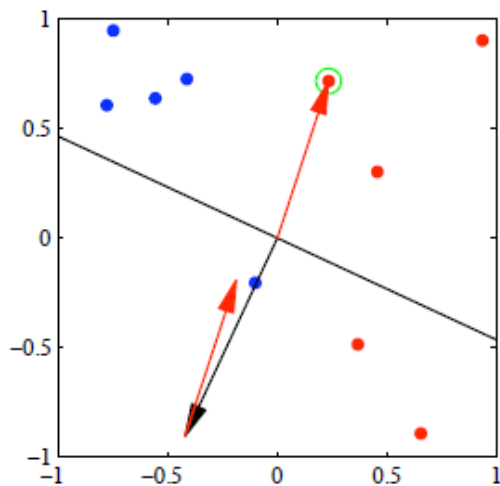
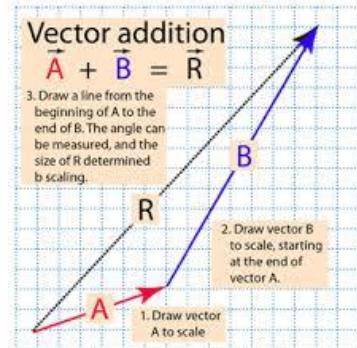
- While updating with respect to a misclassified input n will lower the error for that input, the error for other misclassified inputs may increase.
- Also, new inputs that had been classified correctly may now be misclassified.
- The result is that the perceptron algorithm is not guaranteed to reduce the total error monotonically at each stage.



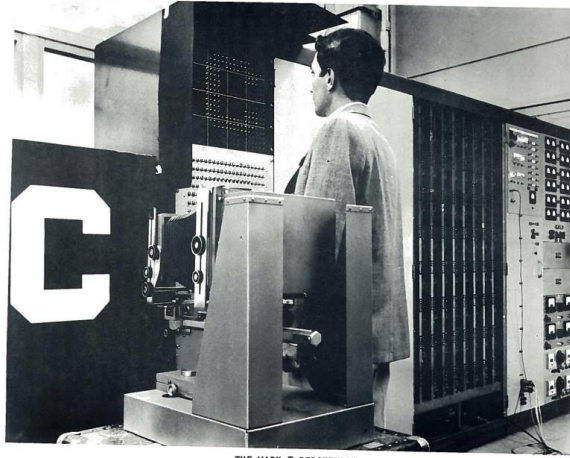
Rosenblatt

if in fact the data are linearly separable, then the algorithm is guaranteed to find an exact solution in a finite number of steps

Example



Mark 1 Perceptron Hardware (1960)



THE MARK 1 PERCEPTRON

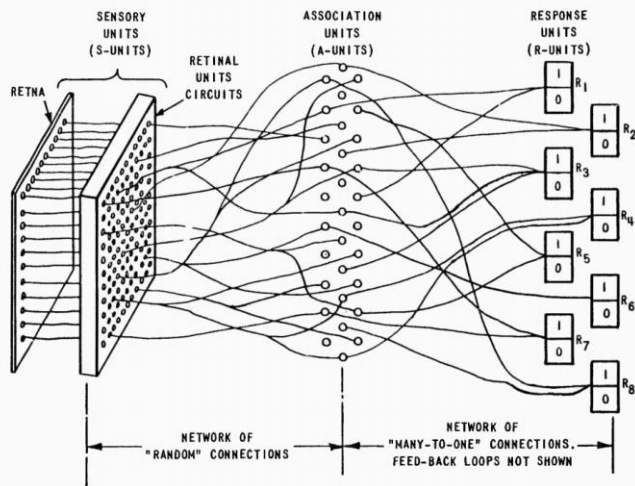
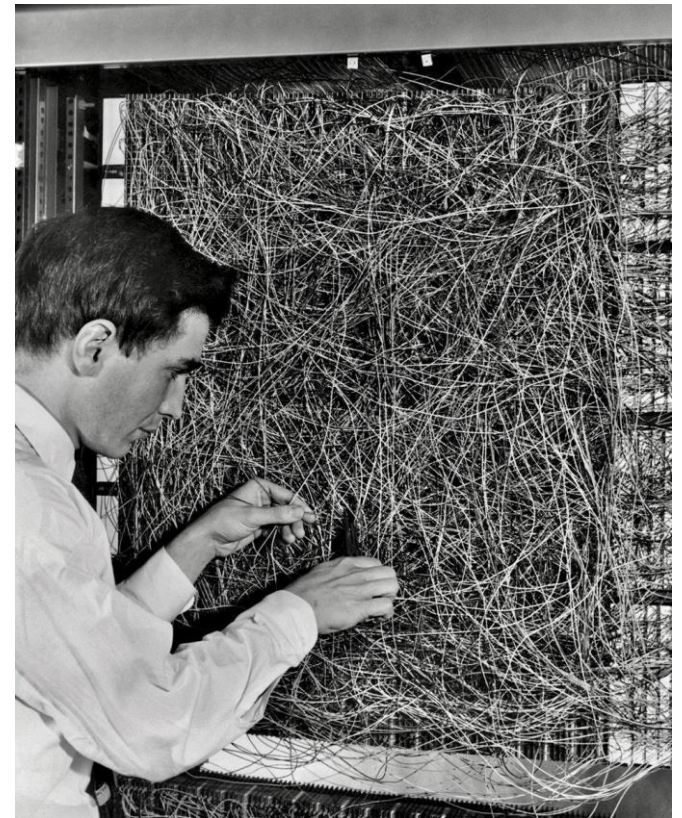


Figure 1 ORGANIZATION OF THE MARK 1 PERCEPTRON



Limitation

- The Perceptron Convergence Theorem is an important result. However, there are practical limitations:
 - Convergence may be slow
 - If the data are not separable, the algorithm will not converge.
 - We will only know that the data are separable once the algorithm converges.
 - The solution is in general not unique, and will depend upon initialization, scheduling of input vectors, and the learning rate η .

- Perceptron Algorithm
- **Least-Squares Classifiers**
- Fisher's Linear Discriminant
- Support Vector Machines

Dealing with Non-Linearly Separable Inputs

- The perceptron algorithm fails when the training data are not perfectly linearly separable.
- Let's now turn to methods for learning the parameter vector \mathbf{w} of a perceptron (linear classifier) even when the training data are not linearly separable.

Learning the Parameters

$$y(\mathbf{X}) = \mathbf{W}^T \mathbf{X}$$

Training dataset $(\mathbf{x}_n, \mathbf{t}_n)$, $n=1, \dots, N$

where we use the 1-of- K coding scheme for \mathbf{t}_n

Let \mathbf{T} be the $N \times K$ matrix whose n^{th} row is \mathbf{t}_n

Let \mathbf{X} be the $N \times (D+1)$ matrix whose n^{th} row is \mathbf{x}_n

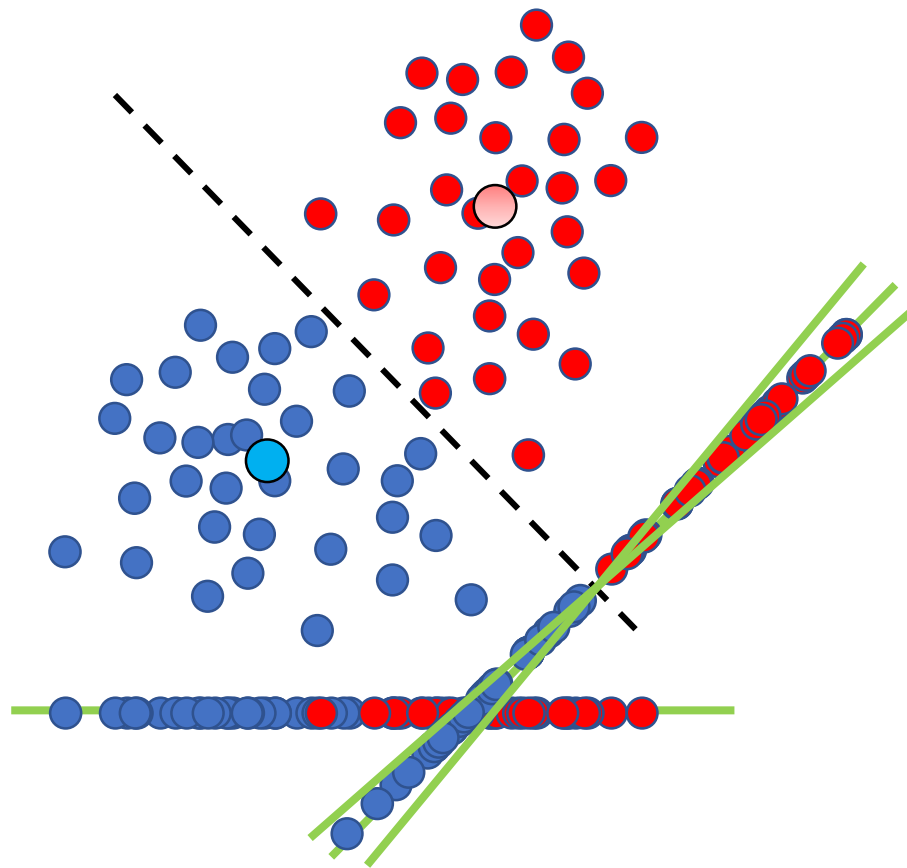
Let $\mathbf{R}(\mathbf{W}) = \mathbf{XW} - \mathbf{T}$

We define the error as $E(\mathbf{W}) = \frac{1}{2} \sum_{i,j} R_{i,j}^2 = \frac{1}{2} \text{Tr}\{\mathbf{R}(\mathbf{W})^T \mathbf{R}(\mathbf{W})\}$

Setting derivative wrt \mathbf{W} to 0 yields:

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{T}$$

- Perceptron Algorithm
- Least-Squares Classifiers
- **Fisher's Linear Discriminant**
- Support Vector Machines



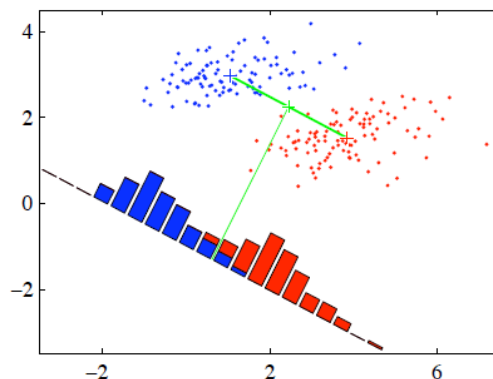
Fisher's Linear Discriminant

- Another way to view linear discriminants: find the 1D subspace that maximizes the separation between the two classes.

- Let

$$\mathbf{m}_1 = \frac{1}{N_1} \sum_{n \in C_1} \mathbf{x}_n, \quad \mathbf{m}_2 = \frac{1}{N_2} \sum_{n \in C_2} \mathbf{x}_n$$

- We might choose \mathbf{w} to maximize $\mathbf{w}^T(\mathbf{m}_2 - \mathbf{m}_1)$, subject to $\|\mathbf{w}\|=1$



$$\mathbf{w} \propto (\mathbf{m}_2 - \mathbf{m}_1)$$

Not optimal if conditional distributions are not isotropic!

Fisher's Linear Discriminant

Let $m_1 = \mathbf{w}^T \mathbf{m}_1$, $m_2 = \mathbf{w}^T \mathbf{m}_2$ be the conditional means on the 1D subspace.

Let $s_k^2 = \sum_{n \in C_k} (y_n - m_k)^2$ be the within-class variance on the subspace for class C_k

The Fisher criterion is then $J(\mathbf{w}) = \frac{(m_2 - m_1)^2}{s_1^2 + s_2^2}$

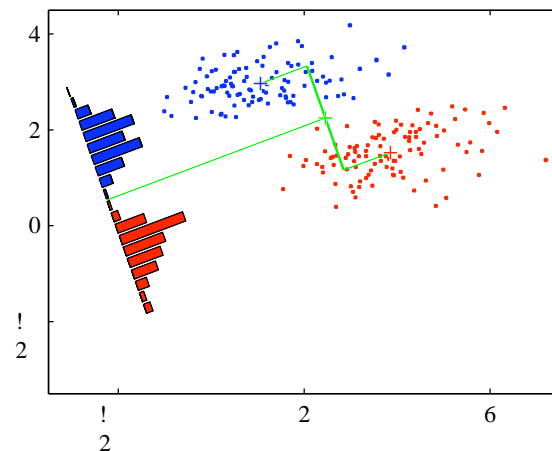
This can be rewritten as $J(\mathbf{w}) = \frac{\mathbf{w}^t \mathbf{S}_B \mathbf{w}}{\mathbf{w}^t \mathbf{S}_W \mathbf{w}}$

where

$\mathbf{S}_B = (\mathbf{m}_2 - \mathbf{m}_1)(\mathbf{m}_2 - \mathbf{m}_1)^t$ is the between-class variance

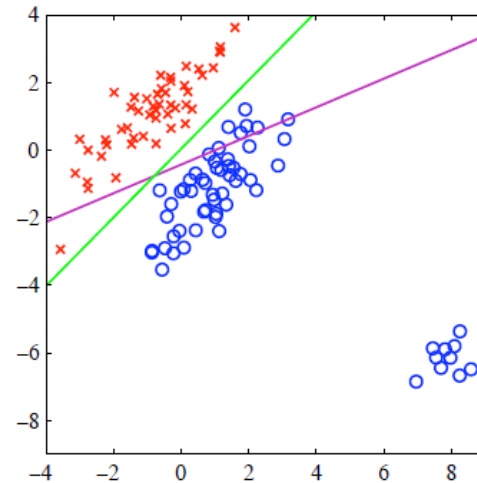
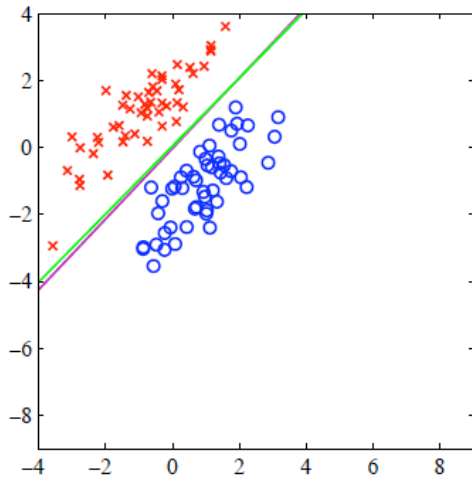
$\mathbf{S}_W = \sum_{n \in C_1} (\mathbf{x}_n - \mathbf{m}_1)(\mathbf{x}_n - \mathbf{m}_1)^t + \sum_{n \in C_2} (\mathbf{x}_n - \mathbf{m}_2)(\mathbf{x}_n - \mathbf{m}_2)^t$ is the within-class variance

$J(\mathbf{w})$ is maximized for $\mathbf{w} \propto \mathbf{S}_W^{-1} (\mathbf{m}_2 - \mathbf{m}_1)$

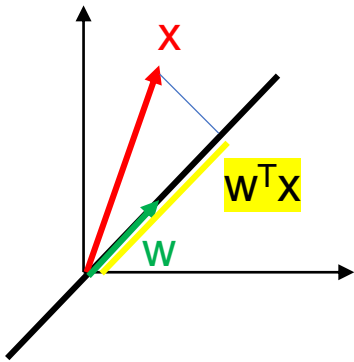


Limitation

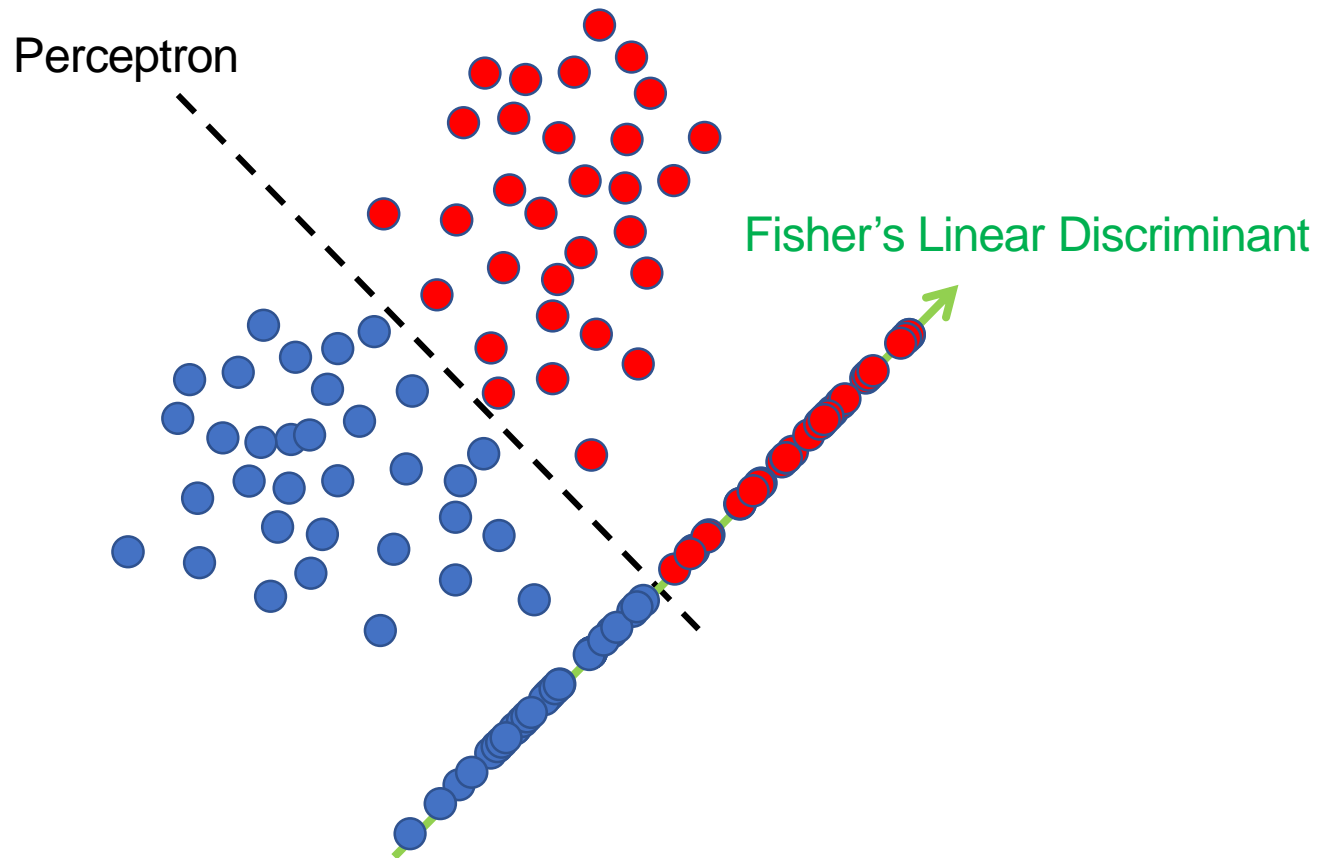
- Sensitivity to outliers



Perceptron vs LD



$w^T x$ is the projection of x into a 1-dimensional subspace spanned by w

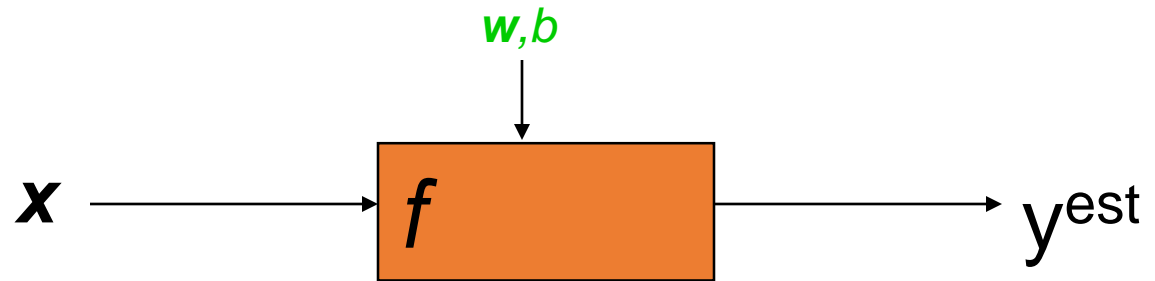


- Perceptron Algorithm
- Least-Squares Classifiers
- Fisher's Linear Discriminant
- **Support Vector Machines**

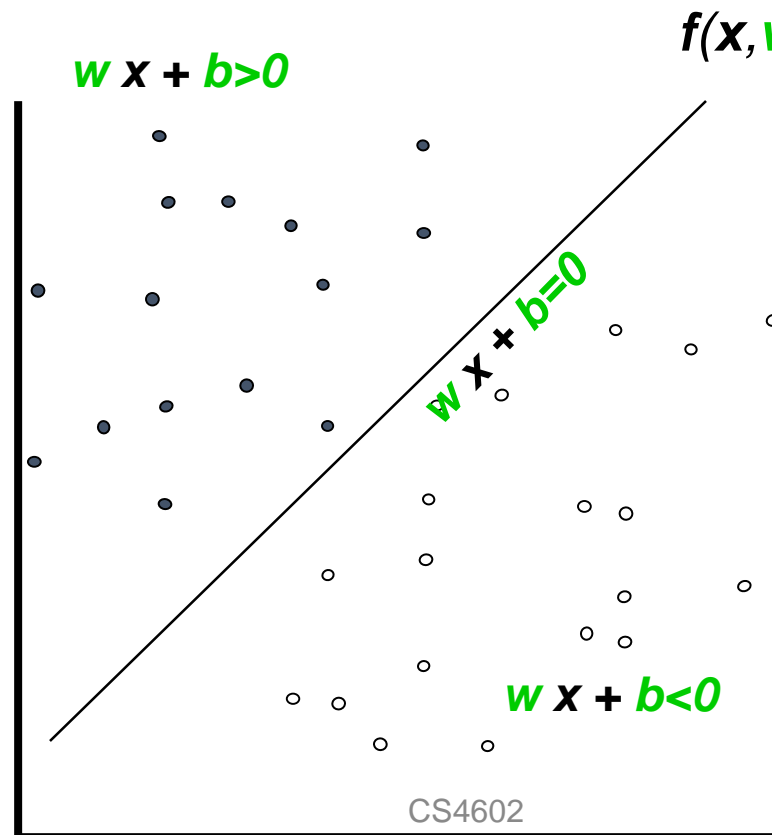
SVM: Motivation

- The perceptron algorithm is guaranteed to provide a linear decision surface that separates the training data, if one exists
- What if it doesn't exist? (back to this later)
- There are an infinite number of solutions, and the solution returned by the perceptron algorithm depends on the initial conditions, the learning rate and the order in which training data are processed.
- While all solutions achieve a perfect score on the training data, they won't all necessarily generalize as well to new inputs.

Linear Classifiers



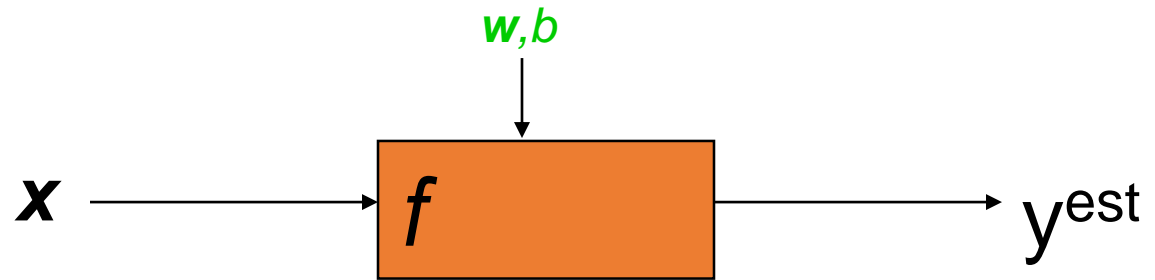
- +1
- -1



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

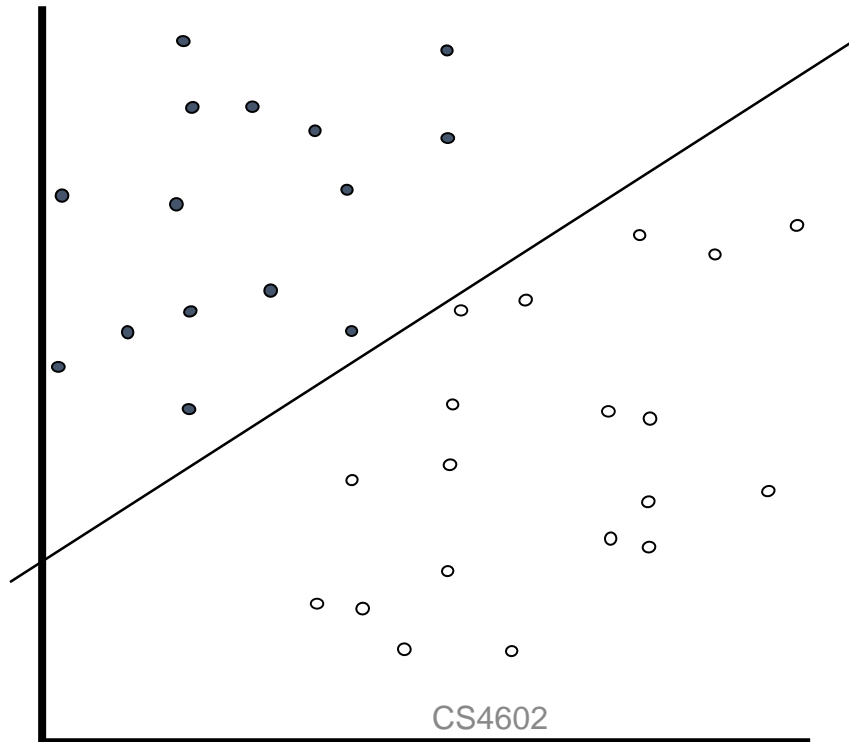
How would you classify this data?

Linear Classifiers



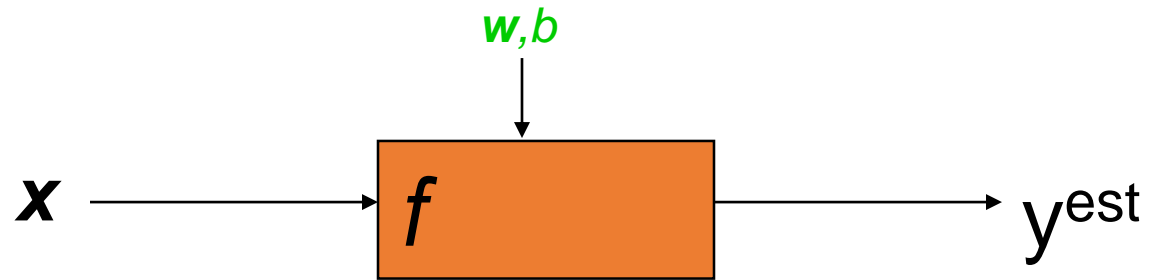
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

- +1
- -1



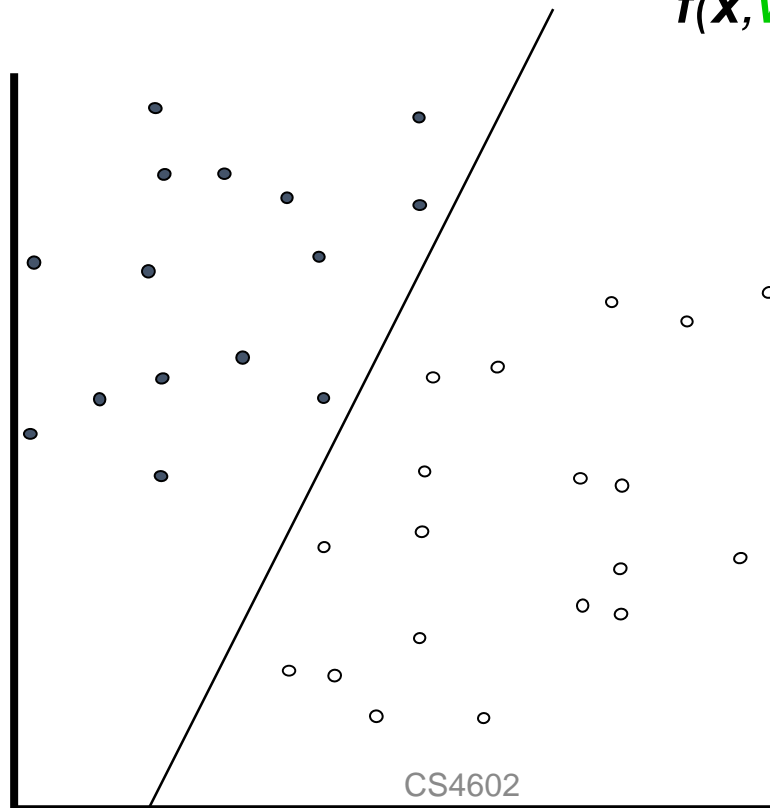
How would you classify this data?

Linear Classifiers



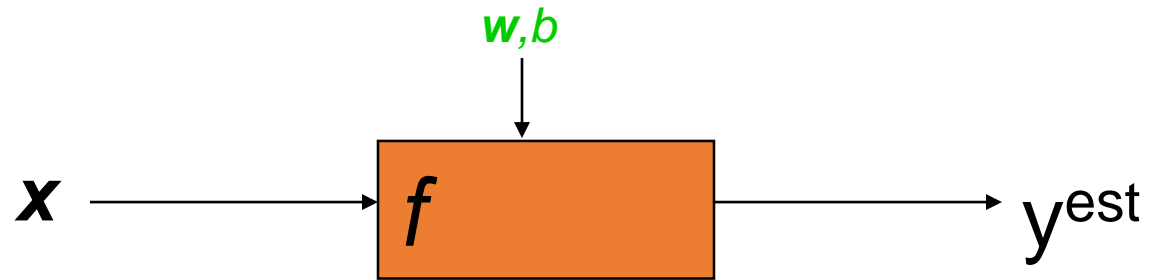
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

- +1
- -1

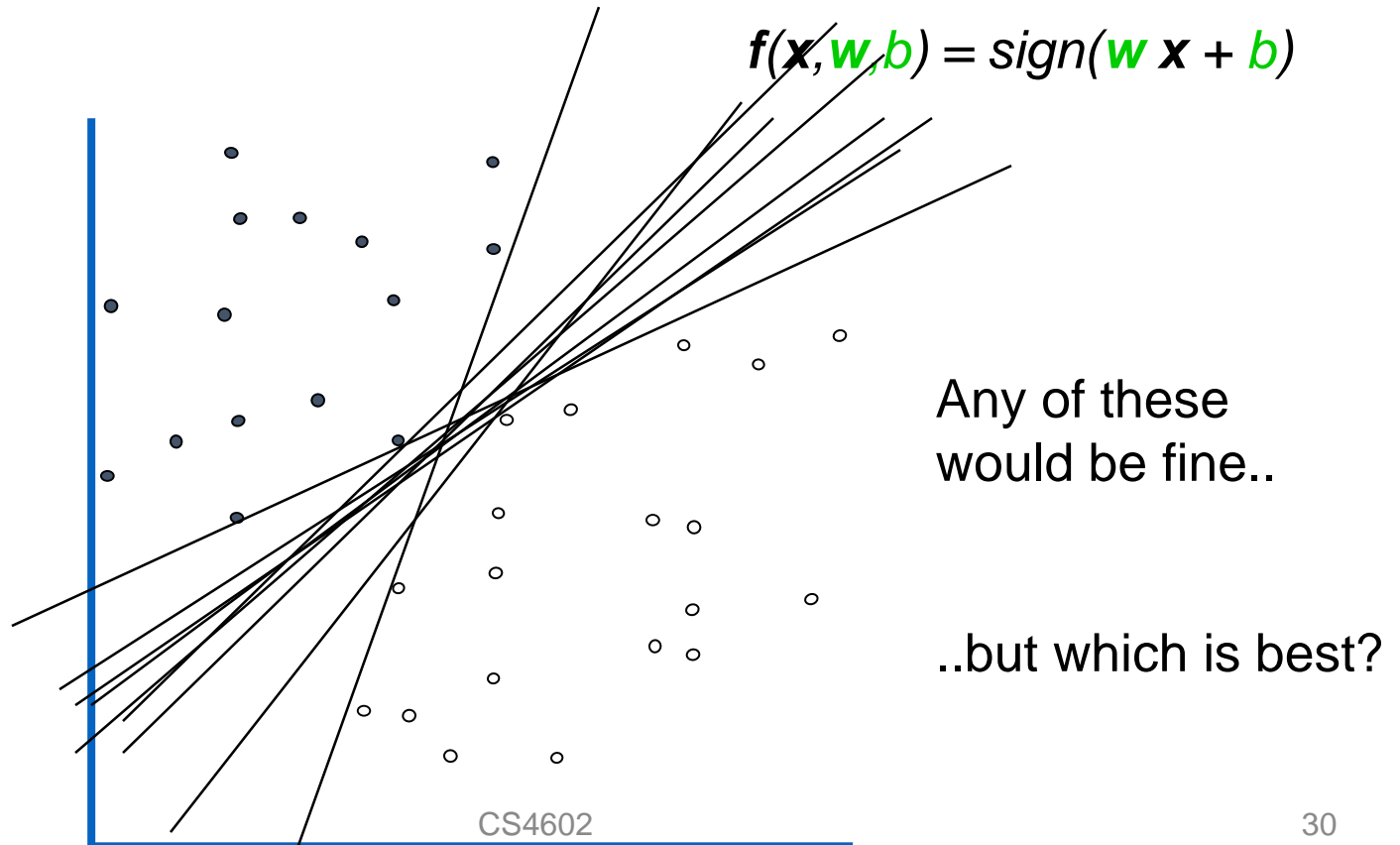


How would you classify this data?

Linear Classifiers



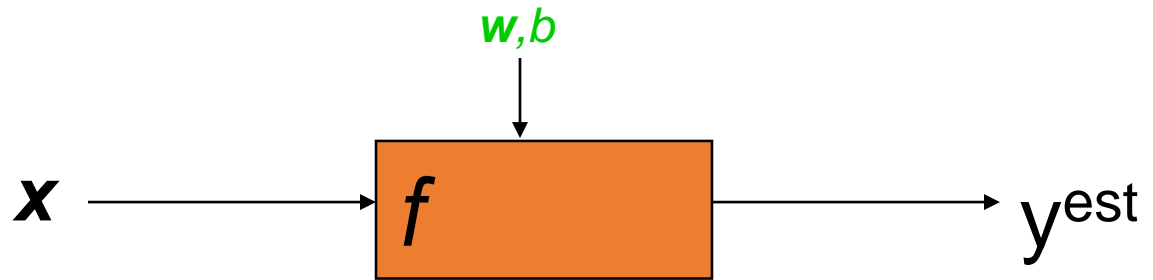
- +1
- -1



Support Vector Machine

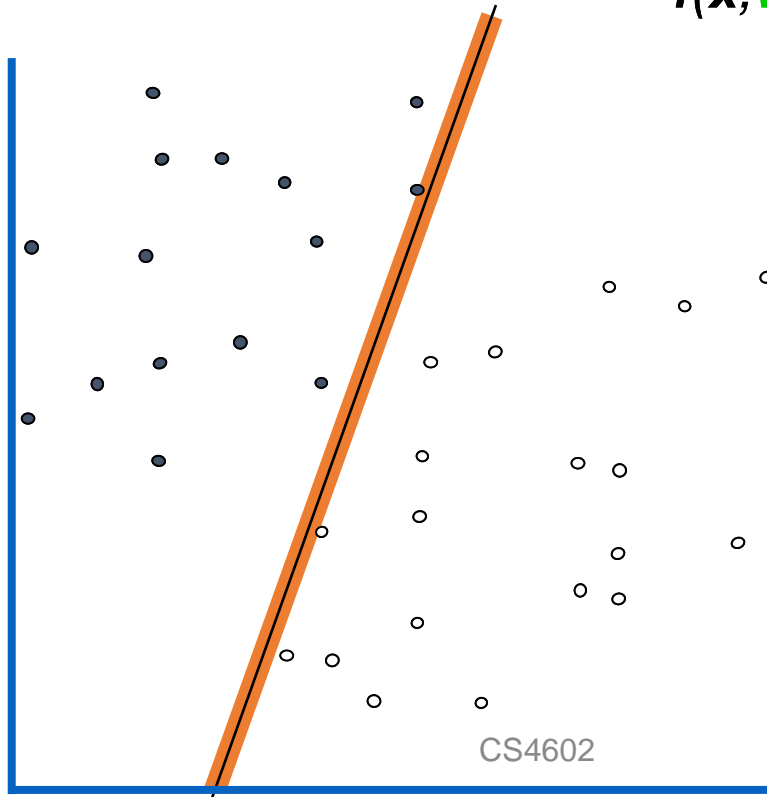
- Which hyperplane to “best” separate data?
- What if it is impossible to separate the data by a hyperplane?

Classifier Margin



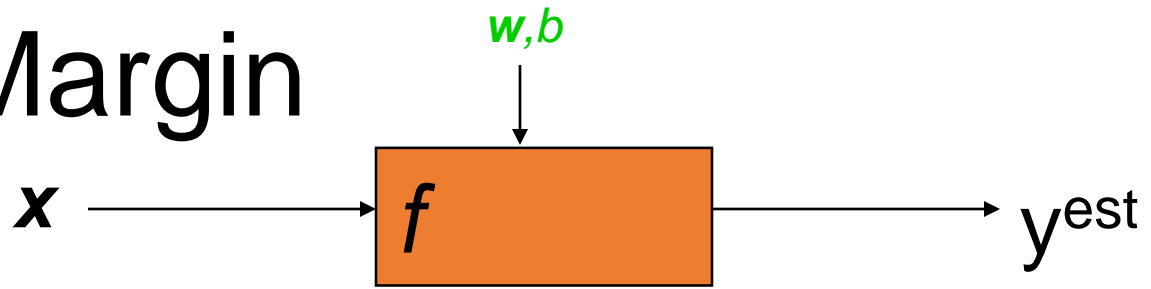
$$f(\mathbf{x}, \mathbf{w}, b) = \text{sign}(\mathbf{w} \mathbf{x} + b)$$

- +1
- -1



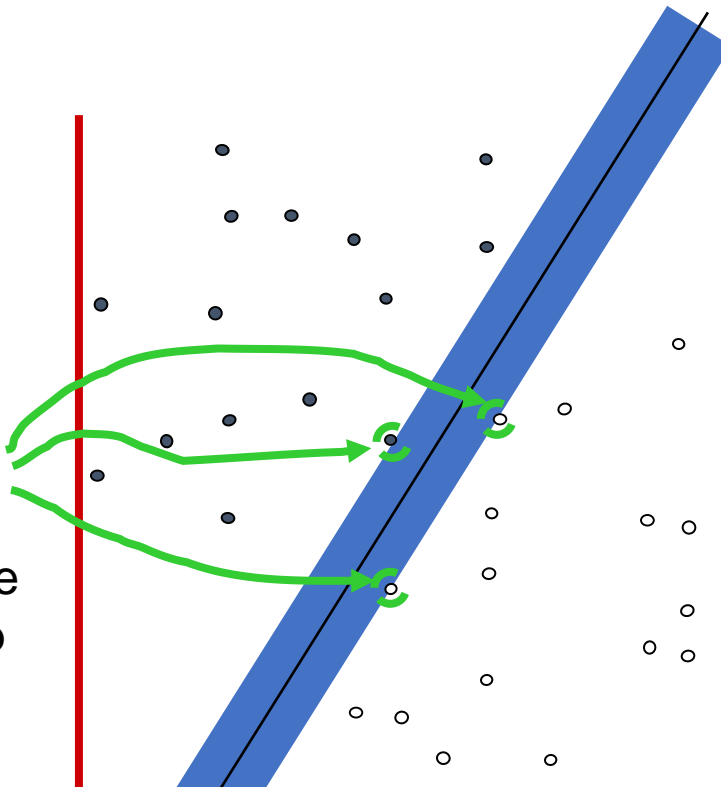
Define the **margin** of a linear classifier as the width that the boundary could be increased by before hitting a datapoint.

Maximum Margin



- +1
- -1

Support Vectors
are those
datapoints that the
margin pushes up
against



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

The **maximum margin linear classifier** is the linear classifier with the maximum margin.

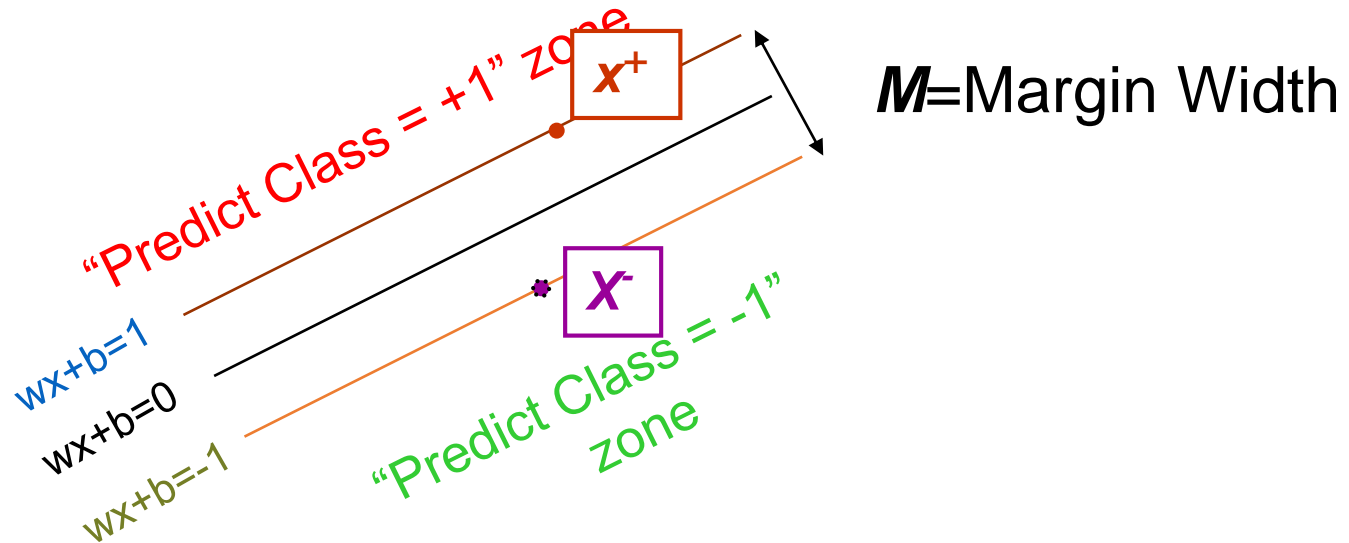
This is the simplest kind of SVM (Called an LSVM)

Implies that only support vectors are important; other training examples are ignorable.

CS4602

Linear

Linear SVM Mathematically



What we know:

- $\mathbf{w} \mathbf{x}^+ + b = +1$
- $\mathbf{w} \mathbf{x}^- + b = -1$
- $\mathbf{w} (\mathbf{x}^+ - \mathbf{x}^-) = 2$

$$M = \frac{(\mathbf{x}^+ - \mathbf{x}^-) \cdot \mathbf{w}}{|\mathbf{w}|} = \frac{2}{|\mathbf{w}|}$$

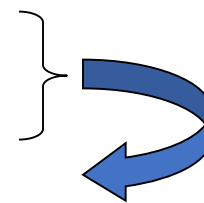
Linear SVM Mathematically

- Goal: 1) **Correctly classify all training data**

$$wx_i + b \geq 1 \quad \text{if } y_i = +1$$

$$wx_i + b \leq -1 \quad \text{if } y_i = -1$$

$$y_i(wx_i + b) \geq 1 \quad \text{for all } i$$



- 2) **Maximize the Margin**
same as minimize

$$M = \frac{2}{|w|}$$

$$\frac{1}{2} w^t w$$

- We can formulate a **Quadratic Optimization Problem** and solve for w and b

- Minimize $\Phi(w) = \frac{1}{2} w^t w$

subject to $y_i(wx_i + b) \geq 1 \quad \forall i$

The Dual Problem

$$\begin{aligned} \max. \quad W(\boldsymbol{\alpha}) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{subject to } \alpha_i &\geq 0, \sum_{i=1}^n \alpha_i y_i = 0 \end{aligned}$$

- This is a quadratic programming (QP) problem
 - A global maximum of α_i can always be found

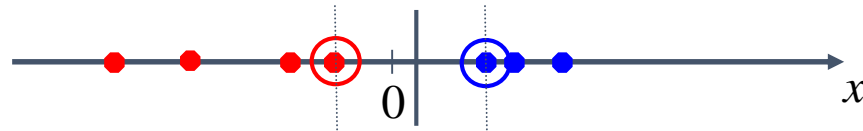
- \mathbf{w} can be recovered by
$$\mathbf{w} = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i$$

Support Vector Machine

- Which hyperplane to separate data?
- What if it is impossible to separate the data by a hyperplane?

Linear SVM Mathematically

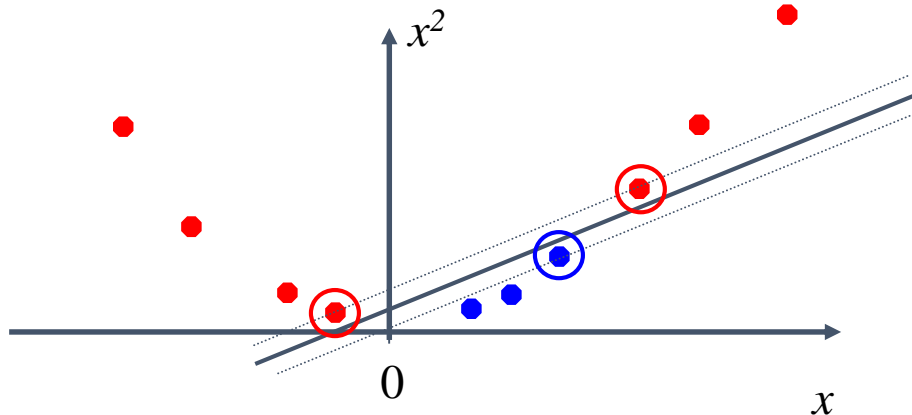
- Datasets that are linearly separable with some noise work out great:



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

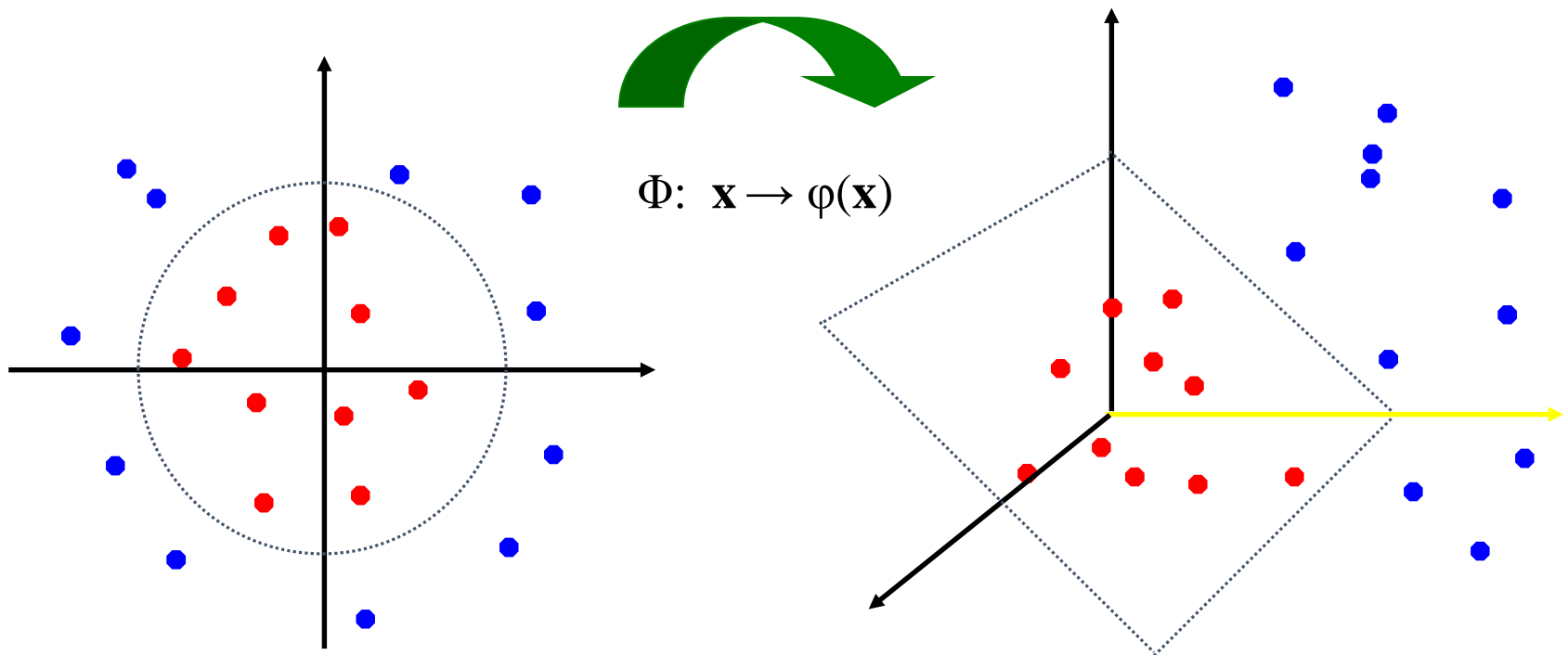


- How about... mapping data to a higher-dimensional space:



Non-linear SVMs: Feature spaces

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



What if we add a dimension which represents the distance to the origin?

The “Kernel Trick”

Avoids the **explicit mapping** to learn a nonlinear function or decision boundary!

- The linear classifier relies on dot product between vectors $K(x_i, x_j) = x_i^T x_j$
- If every data point is mapped into high-dimensional space via some transformation $\Phi: x \rightarrow \phi(x)$, the dot product becomes:

$$K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

- A *kernel function* is some function that corresponds to an inner product in some expanded feature space.
- Example:

2-dimensional vectors $x = [x_1 \ x_2]$; let $K(x_i, x_j) = (1 + x_i^T x_j)^2$,

Need to show that $K(x_i, x_j) = \phi(x_i)^T \phi(x_j)$:

$$\begin{aligned} K(x_i, x_j) &= (1 + x_i^T x_j)^2, \\ &= 1 + x_{i1}^2 x_{j1}^2 + 2 x_{i1} x_{j1} x_{i2} x_{j2} + x_{i2}^2 x_{j2}^2 + 2 x_{i1} x_{j1} + 2 x_{i2} x_{j2} \\ &= [1 \ x_{i1}^2 \ \sqrt{2} x_{i1} x_{i2} \ x_{i2}^2 \ \sqrt{2} x_{i1} \ \sqrt{2} x_{i2}]^T [1 \ x_{j1}^2 \ \sqrt{2} x_{j1} x_{j2} \ x_{j2}^2 \ \sqrt{2} x_{j1} \ \sqrt{2} x_{j2}] \\ &= \phi(x_i)^T \phi(x_j), \quad \text{where } \phi(x) = [1 \ x_1^2 \ \sqrt{2} x_1 x_2 \ x_2^2 \ \sqrt{2} x_1 \ \sqrt{2} x_2] \end{aligned}$$

Examples of Kernel Functions

- Linear: $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial of power p : $K(\mathbf{x}_i, \mathbf{x}_j) = (1 + \mathbf{x}_i^T \mathbf{x}_j)^p$
- Gaussian (radial-basis function network):

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|^2}{2\sigma^2}\right)$$

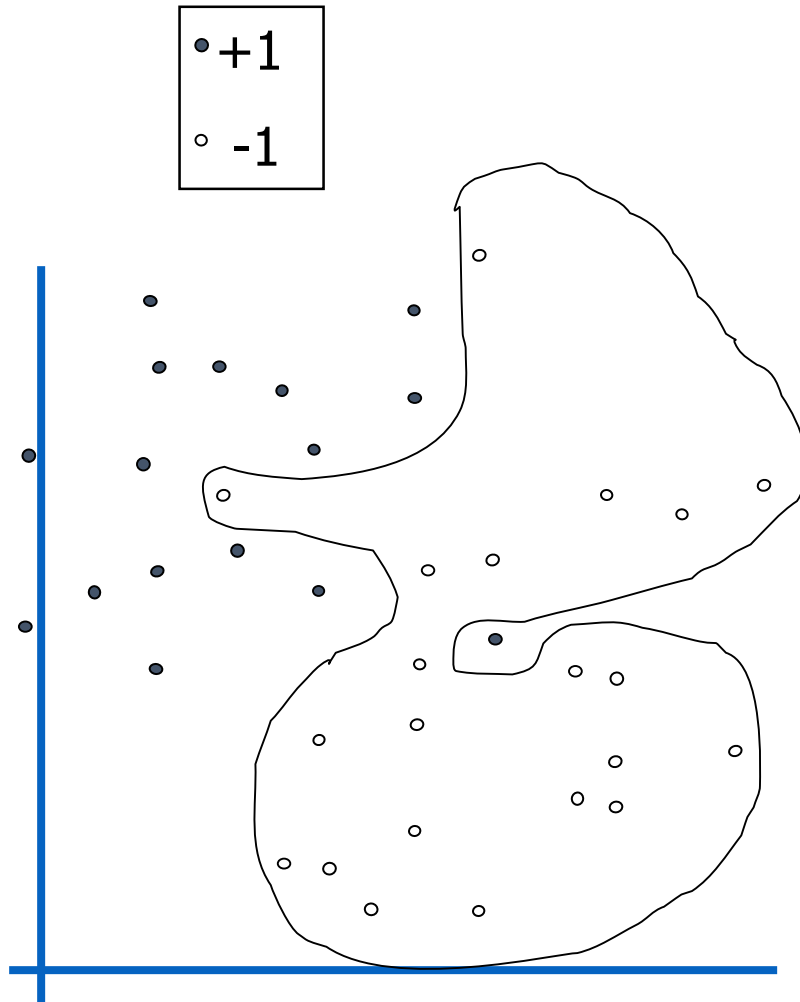
- Sigmoid: $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\beta_0 \mathbf{x}_i^T \mathbf{x}_j + \beta_1)$

Nonlinear!

Nonlinear SVM - Recap

- SVM locates a separating hyperplane in the feature space and classify points in that space.
- It does not need to represent the space explicitly, simply by defining a kernel function.
- The kernel function plays the role of the dot product in the feature space.

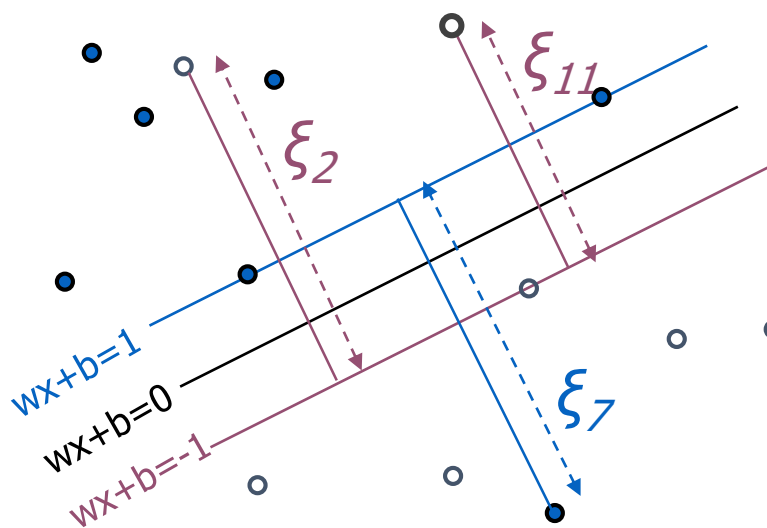
Data noise



- **Hard Margin:** So far we require all data points be classified correctly
 - No training error
- **What if the training set is noisy?**
 - use very powerful kernels

Soft Margin Classification

Slack variables ξ_i can be added to allow misclassification of difficult or noisy examples.



What should our quadratic optimization criterion be?

Minimize

$$\frac{1}{2} W \cdot W + C \sum_{k=1}^M \xi_k$$

- The old formulation:

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$ is minimized and for all $\{(\mathbf{x}_i, y_i)\}$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1$$

- The new formulation incorporating slack variables:

Find \mathbf{w} and b such that

$\Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum \xi_i$ is minimized and for all $\{(\mathbf{x}_i, y_i)\}$

$$y_i (\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{and} \quad \xi_i \geq 0 \text{ for all } i$$

- Parameter C can be viewed as a way to control overfitting.

Some Issues

- Choice of kernel
 - Gaussian or polynomial kernel is default
- Choice of kernel parameters
 - e.g. σ in Gaussian kernel
 - In the absence of reliable criteria, applications rely on the use of a validation set or cross-validation to set such parameters.
- Optimization criterion – Hard margin v.s. Soft margin
 - a lengthy series of experiments in which various parameters are tested

Questions?

A: Can I copy your HW?

B: Yeah, just change it up a bit so it doesn't look obvious you copied.

```
if(count > 10){  
    total += 15;  
}  
else{  
    total += 0;  
}
```

A: Is it okay?

```
total += (count>10?15:0)
```

B:

