# DS 4400

# Machine Learning and Data Mining I
# Spring 2021

Alina Oprea

Associate Professor

Khoury College of Computer Science

Northeastern University
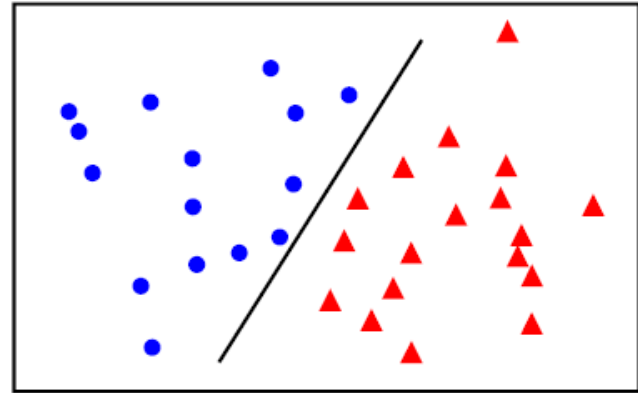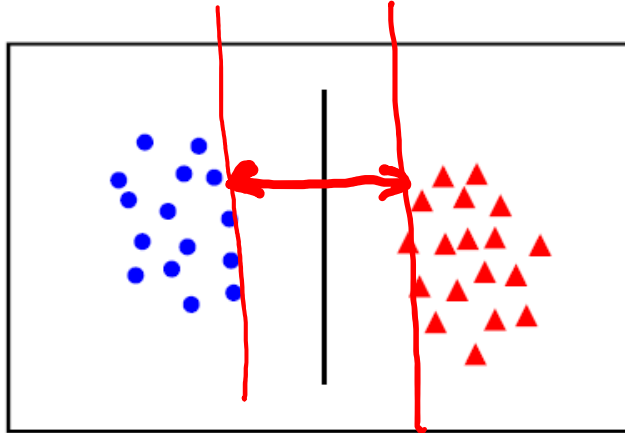
March 16 2021

# Announcements

- Midterm exams have been graded
- HW 4 is due next Friday, March 26
- Project milestone due on March 31
  - Template in Gradescope
- Final exam on Tuesday, April 6
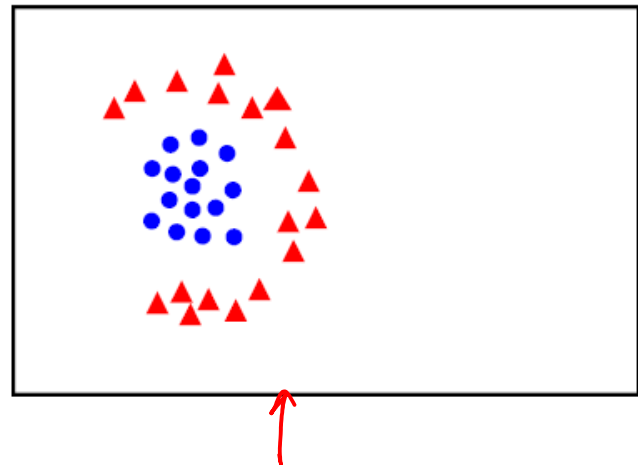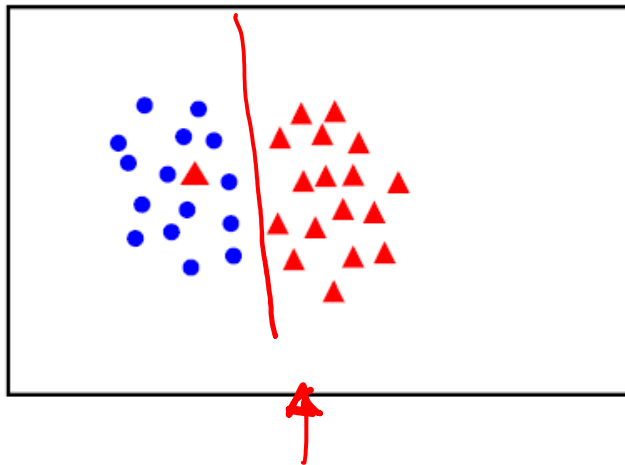  - Review on Thursday, April 1

# Outline

- Support Vector Machines
  - Non-linearly separable data
    - Support vector classifier
- Deep Learning
  - Motivation
  - Goals
- Deep Learning as representation learning
- Perceptron and its limitations

# Linear separability
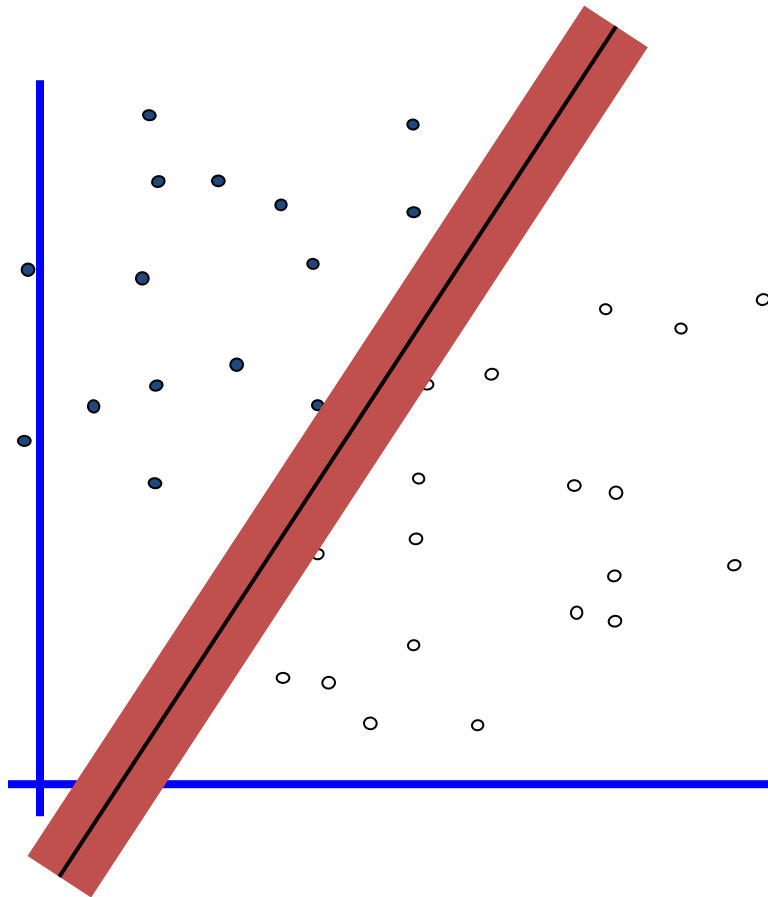
linearly
separable

not
linearly
separable

# Maximum Margin



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

Choose the maximum margin linear classifier: the linear classifier with the maximum margin.

# Maximum margin classifier

LINEARLY   SEPARABLE

- Training data $x_1, \ldots, x_N$ with $x_i = (x_{i1}, \ldots, x_{id})^{\mathrm{T}}$

- Labels are from 2 classes: $y_i \in \{-1, 1\}$

$M > 0$

$h_\theta(x) = \theta^{\mathrm{T}} x$

maximize M          $h_\theta(x_i)$

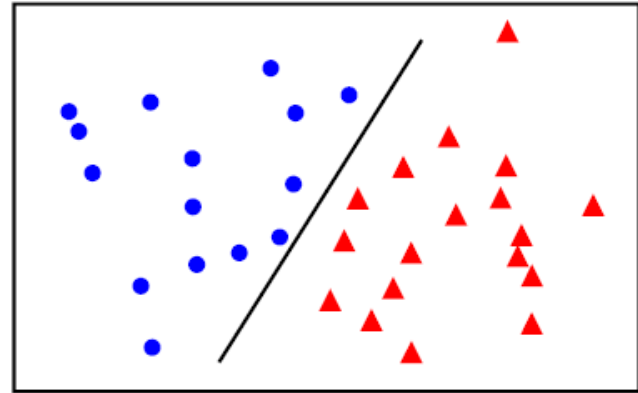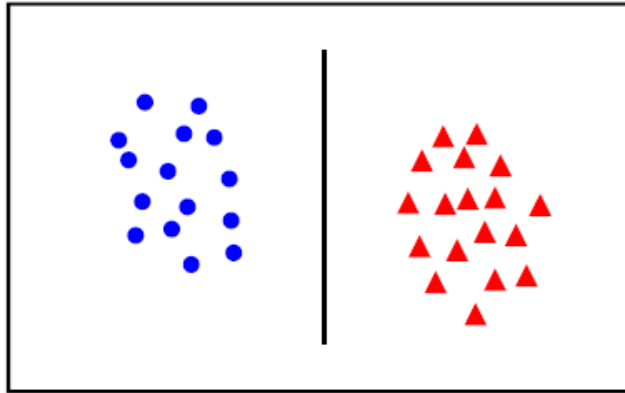$$y_i(\theta_0 + \theta_1 x_{i1} + \cdots \theta_d x_{id}) \geq M \; \forall i$$

$$\|\theta\|_2 = 1$$
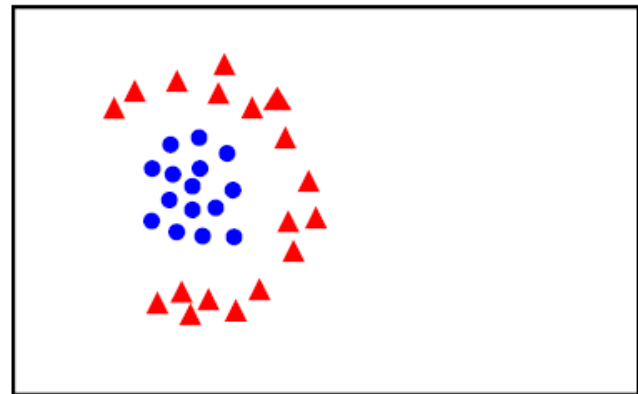
Normalization constraint
(to have unique solution)

Each point is on the
right side of hyper-
plane at distance $\geq M$

6

# Linear separability



linearly
separable

not
linearly
separable
(but almost)

# Support vector classifier (SVC)

- Allow for small number of mistakes on training data

- Soft margin classifier | Linear SVM

$$\max M$$

$y_i h_\theta(x_i)$

$$y_i\left(\theta_0 + \theta_1 x_{i1} + \cdots \theta_d x_{id}\right) \geq M(1 - \epsilon_i)\forall i$$

ERROR

$$\lVert\theta\rVert_2 = 1$$

$$\epsilon_i \geq 0, \sum_i \epsilon_i \leq C$$

Slack

$C =$ regularization parameter

Error Budget (hyper-parameter)

# Support vectors



Larger C

Lower C

Support vectors: all points within the margin of the classifier
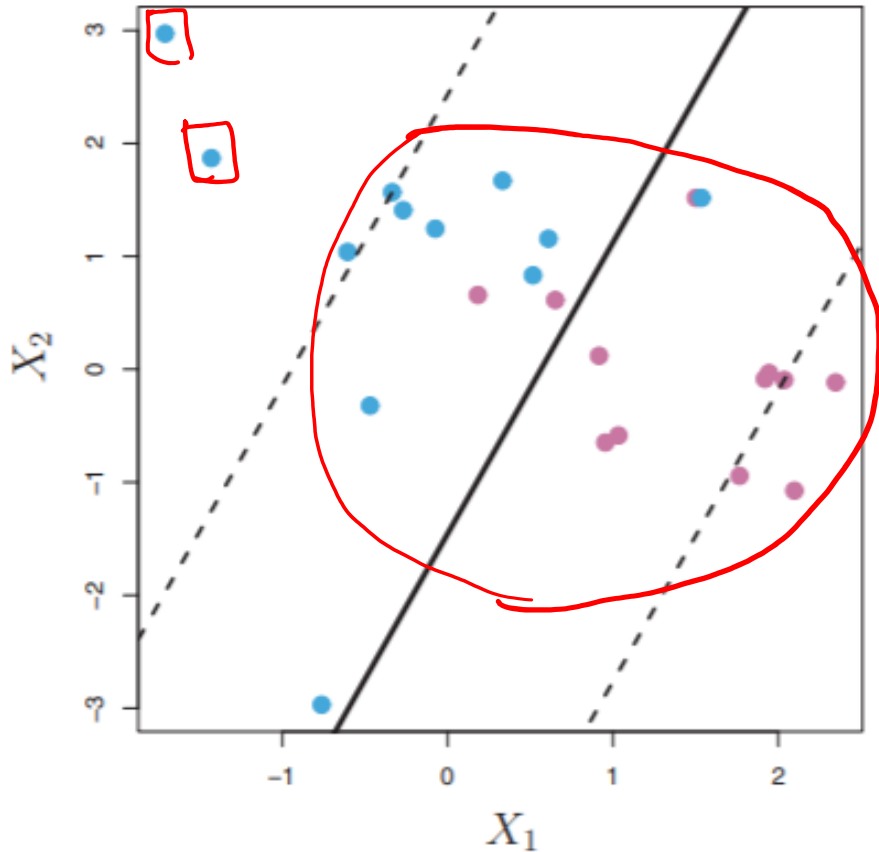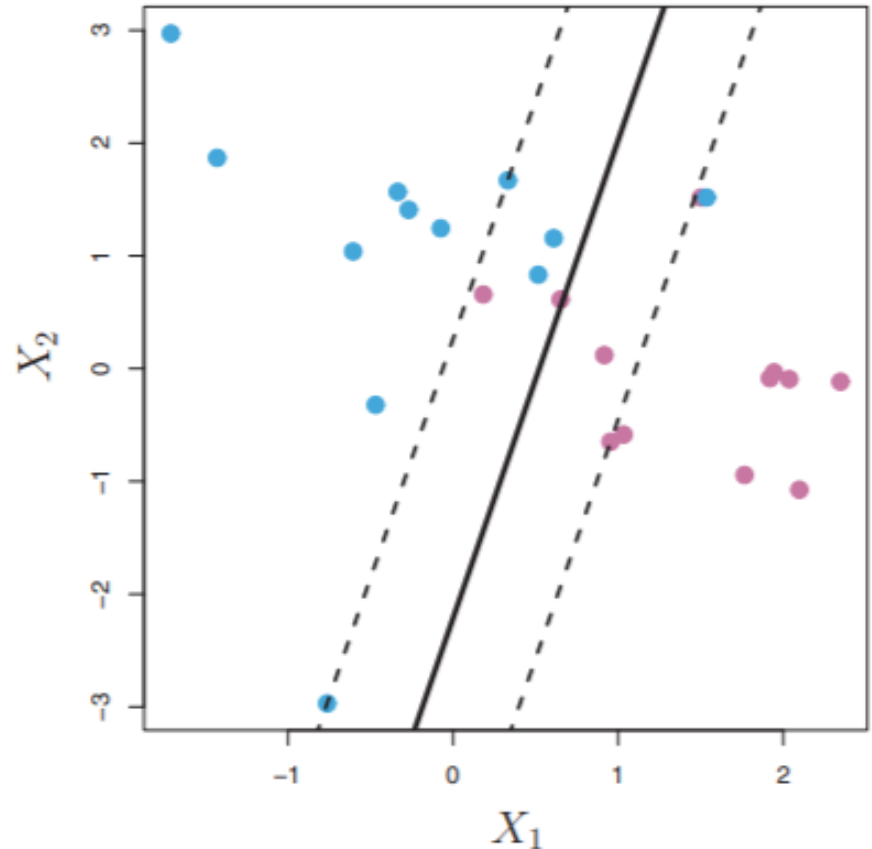
# Support vector classifier

$S$ = index for support vectors

$x_i$, $i \in S$ = support vectors $\in$ Training data

$\sum_{j=1}^{d} z_j x_{ij}$   LINEAR FUNCTION IN $z$

- Just like in separable case, gives solution of the form:

$$h\, f(z) = \theta_0 + \sum_{i \in S} \alpha_i \langle z, x_i \rangle$$

BIAS

$\alpha_i$ = weight for support vector $x_i$

Where $\alpha_i \neq 0$ for support vectors (and $\alpha_i = 0$ for all other training points)

Linear SVM - stores:

- This model is called

  – Support Vector Classifier (SVC)

  – Linear SVM

  – Soft-margin classifier

$\begin{cases} - \text{Support vectors } x_i ; i \in S \\ - \alpha_i \\ - \theta_0 \end{cases}$

INSTANCE LEARNER

# Logistic Regression

$$J(\theta) = -\sum_{i=1}^{N} [y_i \log h_\theta(x_i) + (1 - y_i)\log(1 - h_\theta(x_i))]$$

- Cost of a single instance:

$$\text{cost}\,(h_{\boldsymbol{\theta}}(\boldsymbol{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\boldsymbol{x})) & \text{if } y = 0 \end{cases}$$

- Can re-write objective function as

$$J(\boldsymbol{\theta}) = \sum_{i=1}^{n} \text{cost}\left( h_\theta(x_i), y_i \right)$$

Cross-entropy loss

# Regularized Logistic Regression

$$J(\theta) = -\sum_{i=1}^{N} [y_i \log h_\theta(x_i) + (1 - y_i) \log(1 - h_\theta(x_i))]$$

- We can regularize logistic regression exactly as before:

$$J_{\text{regularized}}(\boldsymbol{\theta}) = J(\boldsymbol{\theta}) + \lambda \sum_{j=1}^{d} \theta_j^2 \qquad \text{RIDGE}$$

$$= J(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}_{[1:d]}\|_2^2$$

L2 regularization

# Hinge Loss

- Linear SVM: $h_\theta(x_i) = \theta^T x_i \in R$ ; $y_i \in \{-1, 1\}$

- Optimization solution equivalent to:

- $J(\theta) = \sum_{i=1}^{N} \max(0, 1 - y_i h_\theta(x_i)) + \lambda \sum_{j=1}^{d} \theta_j^2$

$\underbrace{\qquad\qquad}_{\text{COST FOR TRAINING EXAMPLE } i}$ $\underbrace{\qquad}_{\text{RIDGE}}$

1) $y_i h_\theta(x_i) > 1 \implies |h_\theta(x_i)| > 1 \implies \text{COST} = 0$

2) $y_i h_\theta(x_i) < 1 \implies \text{COST} = 1 - y_i h_\theta(x_i)$

$\lambda = \dfrac{1}{C}$ ; $C = $ Error budget.

USE GRADIENT DESCENT.

13

# Connection to Logistic Regression
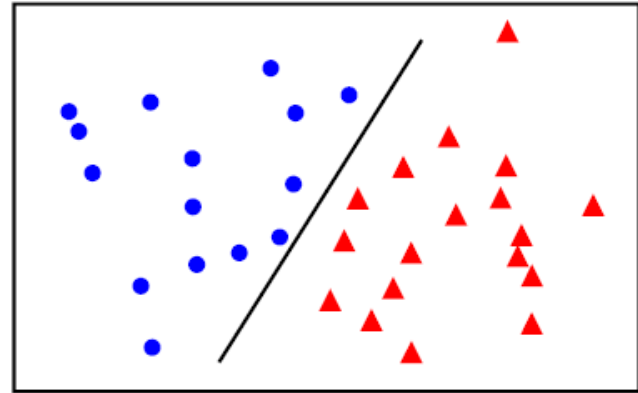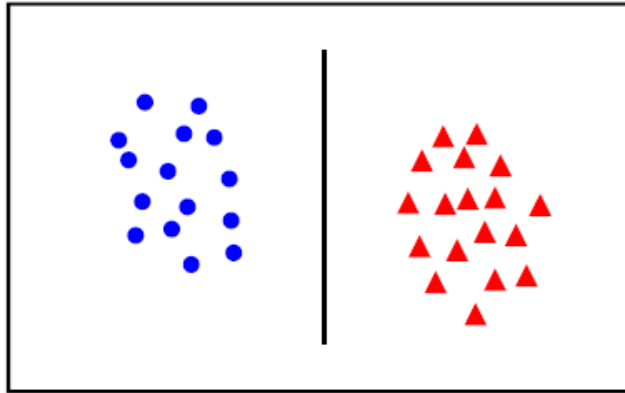
- Logistic regression
  - Cross-entropy loss
- SVM
  - Hinge loss



$$y_i(\theta^T x_i) = y_i \cdot h_\theta(x_i)$$

# Linear separability

linearly
separable

not
linearly
separable

(but almost)

(not even close!)
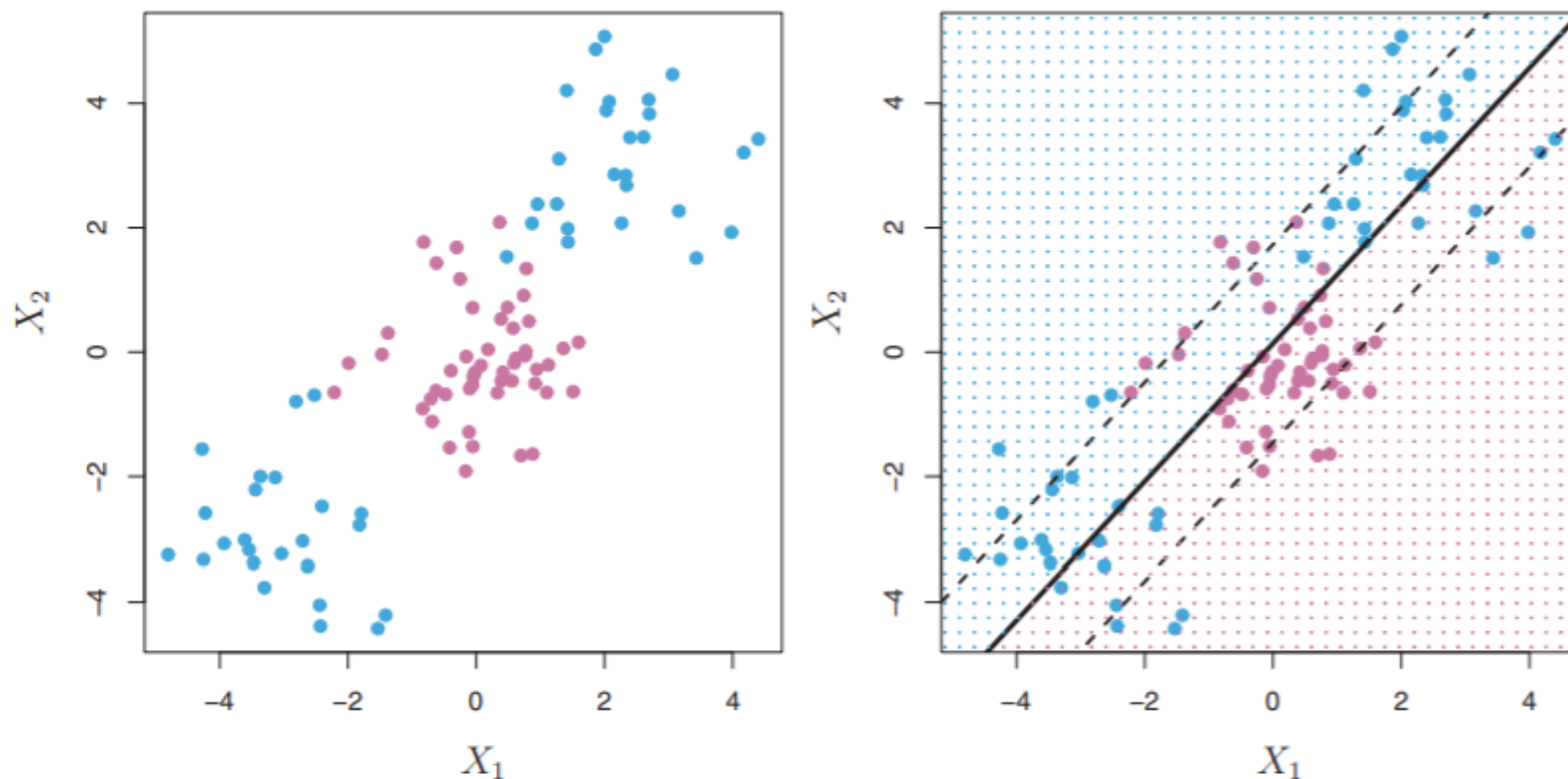
# Non-linear decision



**FIGURE 9.8.** Left: *The observations fall into two classes, with a non-linear boundary between them.* Right: *The support vector classifier seeks a linear boundary, and consequently performs very poorly.*

# More examples



Image from http://www.atrandomresearch.com/iclass/

# Kernels

- Support vector classifier <span style="color:red">**Linear** SVM</span>

<span style="color:red">→</span> – $h(z) = \theta_0 + \sum_{i \in S} \alpha_i < z, x_i >$

$$= \theta_0 + \sum_{i \in S} \alpha_i \sum_{j=1} z_j x_{ij}$$

- – S is set of support vectors

<span style="color:red">$K(x,y)$ : function for similarity | Kernel</span>

<span style="color:red">– symmetric : $K(x,y) = K(y,x)$</span>

<span style="color:red">– close to 0 if points are not similar.</span>

<span style="color:red">$$h(z) = \theta_0 + \sum_{i \in S} \alpha_i \, K(z, x_i)$$</span>

<span style="color:red">$K(x,y) = <x,y>$    LINEAR KERNEL</span>

# The Kernel Trick

"Given an algorithm which is formulated in terms of a positive definite kernel $K_1$, one can construct an alternative algorithm by replacing $K_1$ with another positive definite kernel $K_2$"
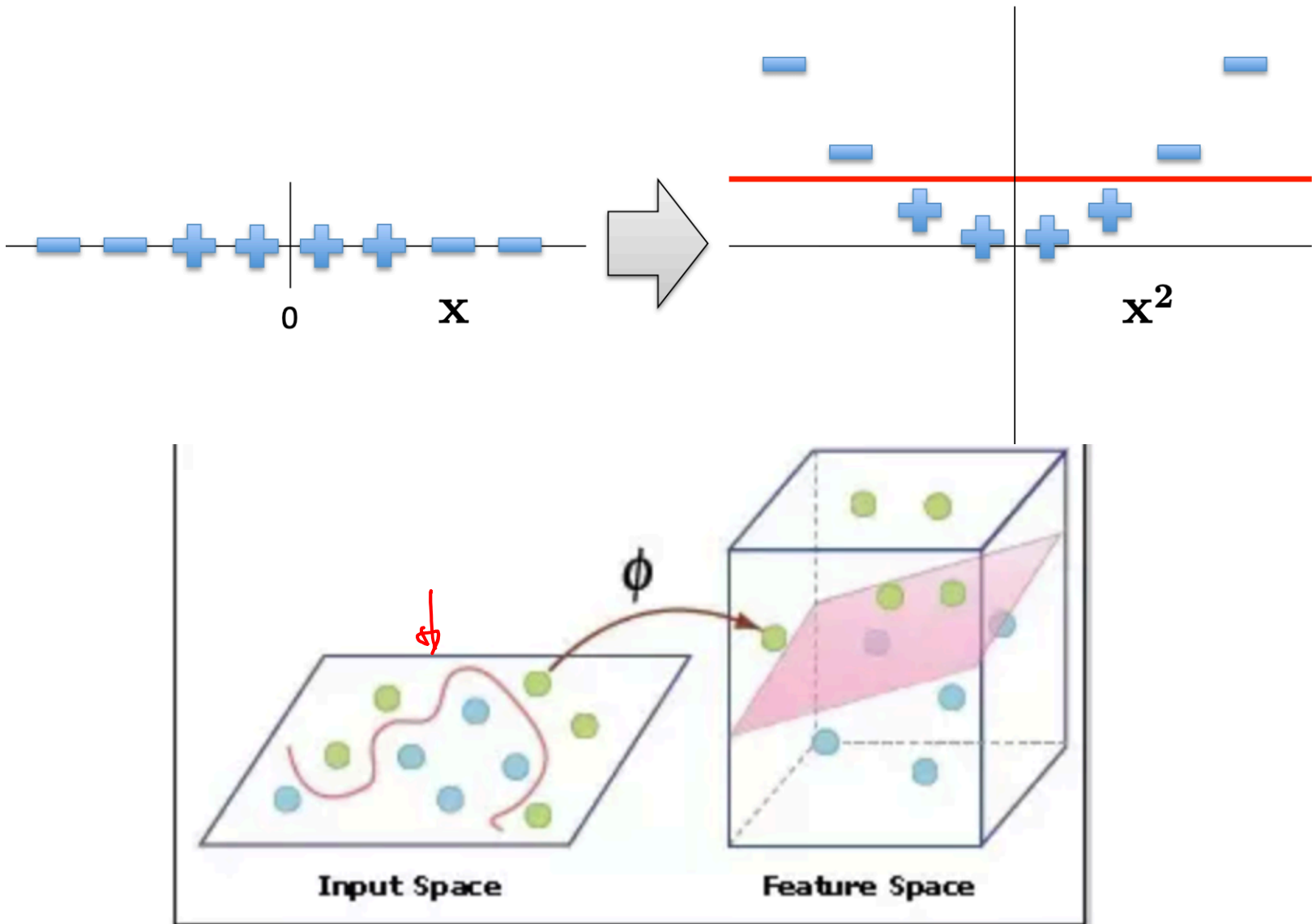
➢ SVMs can use the kernel trick

- Enlarge feature space
- Shape of the kernel changes the decision boundary

# Why Use Kernels

# SVM Classifier

- Select a kernel function

  - The Gram matrix $G_{ij} = K(\mathbf{x}_i, \mathbf{x}_j)$
    - Symmetric matrix
    - Positive semi-definite matrix:
      $\mathbf{z}^\top \mathbf{G} \mathbf{z} \geq 0$ for every non-zero vector $\mathbf{z} \in \mathbb{R}^n$

- Final SVM classifier is linear combination of kernel between testing point and support vectors

  $- h(z) = \theta_0 + \sum_{i \in S} \alpha_i K(z, x_i)$

  *NON-LINEAR DECISION BOUNDARY*

# Kernels

1) LINEAR
$$K(X,Y) = \langle X, Y \rangle = \sum_{i=1}^{d} x_i y_i$$

2) POLYNOMIAL
$$K(X,Y) = \left(1 + \sum_{i=1}^{d} x_i y_i\right)^p \quad ; \quad p = \text{degree of polynomial}$$

3) GAUSSIAN (RADIAL BASIS FUNCTION OR RBF)
$$K(X,Y) = e^{-\sum_{i=1}^{d} (x_i - y_i)^2 / 2\sigma^2}$$

If $X = Y \Rightarrow K(X,Y) = 1$

If $X$ and $Y$ are "far away" / $\|X - Y\|_2 \to \infty$

$$\Rightarrow K(X,Y) \to 0$$

# Examples of SVM classifiers

- Notation
  - S = index of support vectors
  - $\{x_i\}, i \in S$ = set of support vectors
- SVM with polynomial kernel
  - $h(z) = \theta_0 + \sum_{i \in S} \alpha_i \left(1 + \sum_{j=0}^{d} z_j x_{ij}\right)^p$
  - Hyper-parameter p (degree of polynomial)
- SVM with Gaussian / radial kernel
  - $h(z) = \theta_0 + \sum_{i \in S} \alpha_i e^{-\sum_{j=0}^{d}(z_j - x_{ij})^2 / 2\sigma^2}$
  - Hyper-parameter $\sigma$

# Gaussian / Radial Kernel

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

- Has value 1 when $\mathbf{x}_i = \mathbf{x}_j$
- Value falls off to 0 with increasing distance
- Note: Need to do feature scaling <u>before</u> using Gaussian Kernel

$\sigma^2 = 0.5$       $\sigma^2 = 1$       $\sigma^2 = 3$

lower bias,
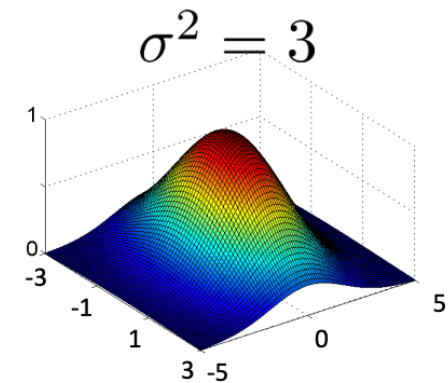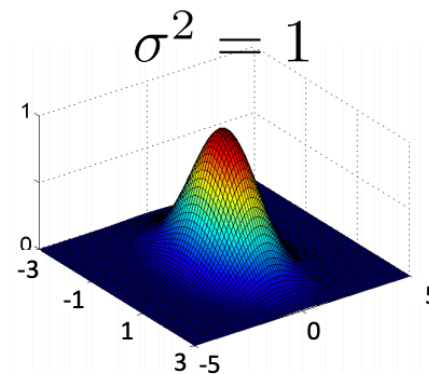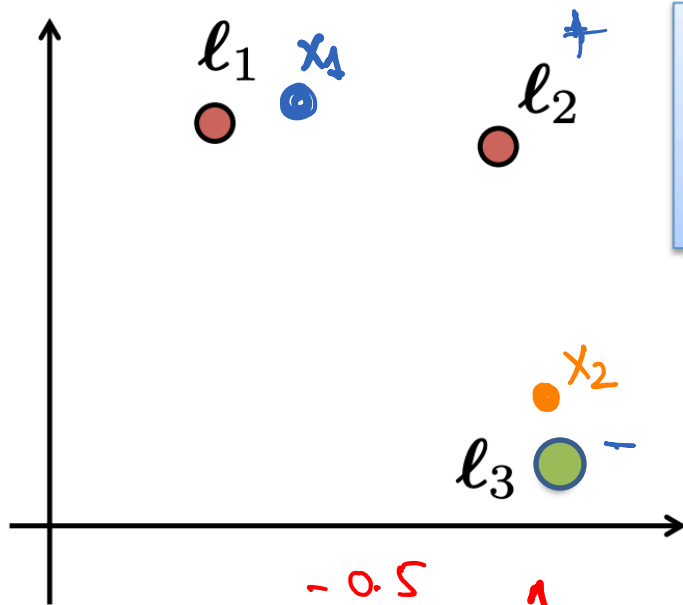higher variance

higher bias,
lower variance

# Gaussian / Radial Kernel Example

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:

$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \geq 0$

$h_\theta(x)$

$x_1$: $K(x_1, \ell_1) \simeq 1$ ; $K(x_1, \ell_2) \simeq 0$ ; $K(x_1, \ell_3) \simeq 0$
$\Rightarrow$ PREDICTION: $-0.5 + 1 \simeq 0.5 \geq 0 \Rightarrow$ PREDICT RED

$x_2$: $K(x_2, \ell_1) \simeq 0$ ; $K(x_2, \ell_2) \simeq 0$ ; $K(x_2, \ell_3) \simeq 1$
$\Rightarrow$ PREDICTION: $-0.5 + 0 + 0 + 0 < 0 \Rightarrow$ PREDICT GREEN

# Gaussian / Radial Kernel Example



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:
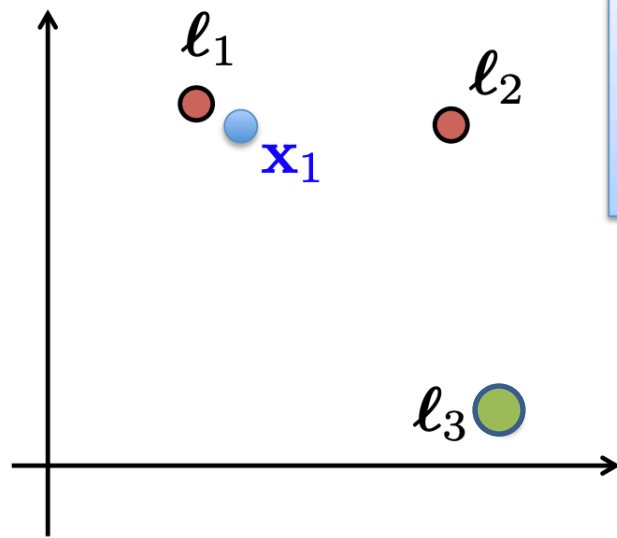
$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \geq 0$

- For $\mathbf{x}_1$, we have $K(\mathbf{x}_1, \boldsymbol{\ell}_1) \approx 1$, other similarities ≈ 0

$$\theta_0 + \theta_1(1) + \theta_2(0) + \theta_3(0)$$
$$= -0.5 + 1(1) + 1(0) + 0(0)$$
$$= 0.5 \geq 0 \text{ , so predict +1}$$

Based on example by Andrew Ng

# Gaussian / Radial Kernel Example



$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$
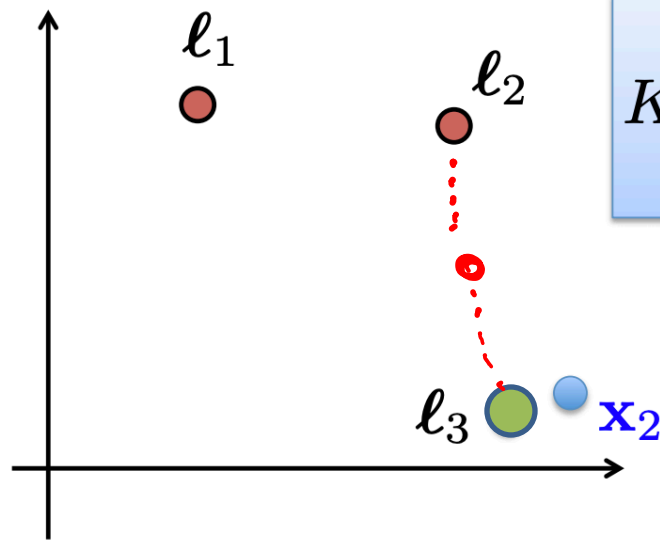
Imagine we've learned that:

$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \geq 0$

- For $\mathbf{x}_2$, we have $K(\mathbf{x}_2, \boldsymbol{\ell}_3) \approx 1$, other similarities $\approx 0$

$$\theta_0 + \theta_1(0) + \theta_2(0) + \theta_3(1)$$
$$= -0.5 + 1(0) + 1(0) + 0(1)$$
$$= -0.5 < 0 \text{, so predict -1}$$

Based on example by Andrew Ng

# Gaussian / Radial Kernel Example

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{2\sigma^2}\right)$$

Imagine we've learned that:
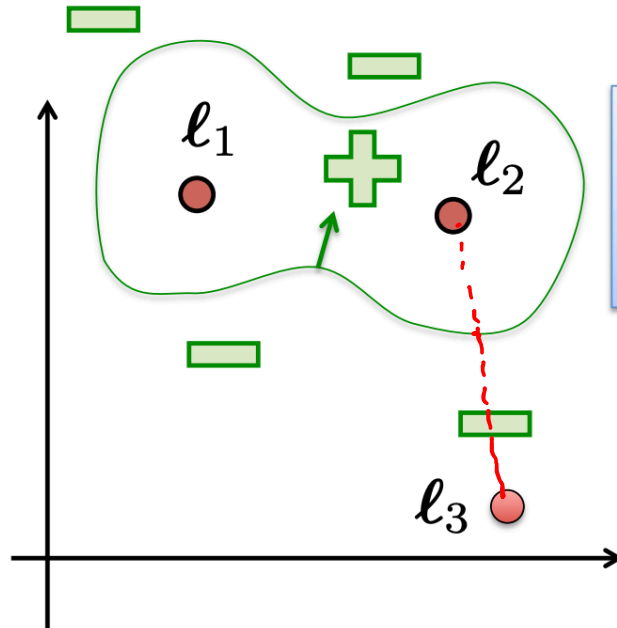
$$\boldsymbol{\theta} = [-0.5, 1, 1, 0]$$

Predict +1 if $\theta_0 + \theta_1 K(\mathbf{x}, \boldsymbol{\ell}_1) + \theta_2 K(\mathbf{x}, \boldsymbol{\ell}_2) + \theta_3 K(\mathbf{x}, \boldsymbol{\ell}_3) \geq 0$

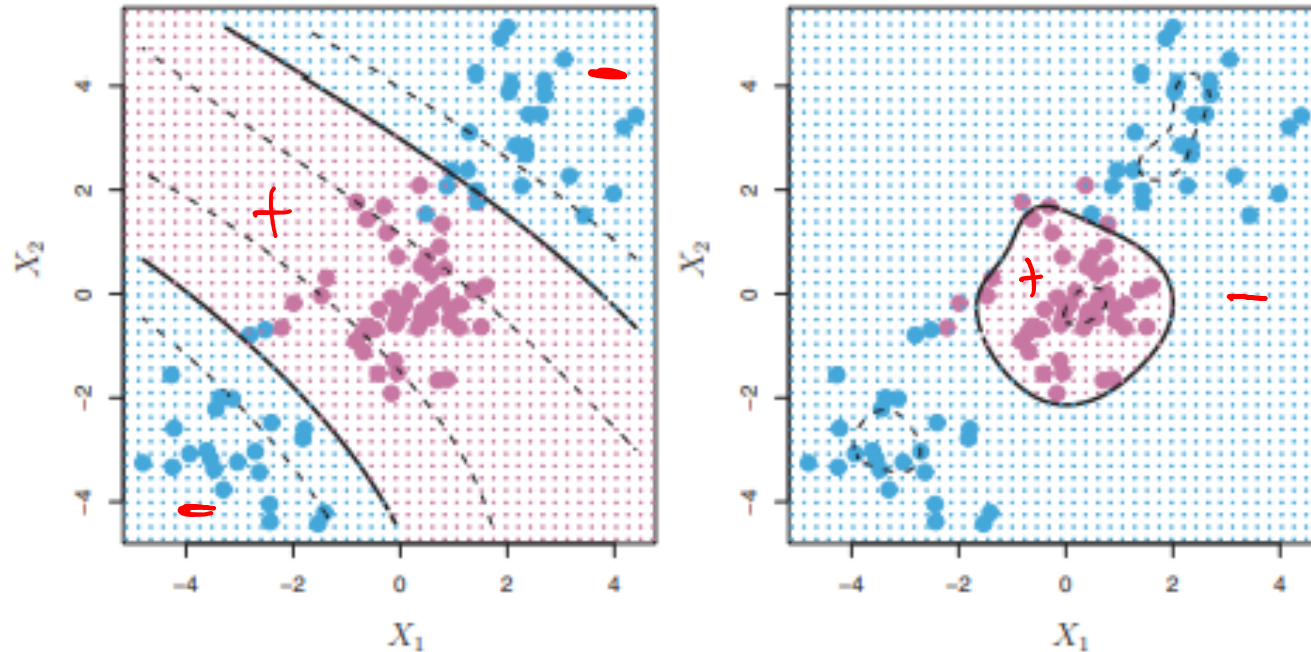Rough sketch of decision surface

# Kernel Example



**FIGURE 9.9.** Left: *An SVM with a polynomial kernel of degree 3 is applied to the non-linear data from Figure 9.8, resulting in a far more appropriate decision rule.* Right: *An SVM with a radial kernel is applied. In this example, either kernel is capable of capturing the decision boundary.*

# Advantages of Kernels

- Generate non-linear features
- More flexibility in decision boundary
- Generate a family of SVM classifiers
- Testing is computationally efficient
  - Cost depends only on support vectors and kernel operation

- Disadvantages
  - Kernels need to be tuned (additional hyper-parameters)

# When to use different kernels?

- If data is (close to) linearly separable, use linear SVM

- Radial or polynomial kernels preferred for non-linear data

- Training radial or polynomial kernels takes longer than linear SVM

- Other kernels
  - Sigmoid
  - Hyperbolic Tangent

# Kernels in ML

- **Kernel ridge regression**
  - Non-linear regression method
- **Kernel Density Estimate (KDE)**
  - Unsupervised method for learning density estimation using kernel functions
- **Kernel PCA**
  - Unsupervised learning for dimensionality reduction
- **Kernel clustering**
  - Create clusters of similar points
  - Similarity between points computed with kernel function

# Review SVM

- SVMs find optimal linear separator
- The kernel trick makes SVMs learn non-linear decision surfaces

- Strength of SVMs:
  - Good theoretical and empirical performance
  - Supports many types of kernels

- Disadvantages of SVMs:
  - "Slow" to train/predict for huge data sets (but relatively fast!)
  - Need to choose the kernel (and tune its parameters)

# Comparing SVM with other classifiers

- SVM is resilient to outliers
  - Similar to Logistic Regression
  - LDA or kNN are not
- SVM can be trained with Gradient Descent
  - Hinge loss cost function
- Supports regularization
  - Can add penalty term (ridge or Lasso) to cost function
- Linear SVM is most similar to Logistic Regression

# Acknowledgements

- Slides made using resources from:
  - Andrew Ng
  - Eric Eaton
  - David Sontag
  - Andrew Moore
  - Yann LeCun
- Thanks!