

# Introduction to Machine Learning

## NPFL 054

<http://ufal.mff.cuni.cz/course/npfl054>

Barbora Hladká  
hladka@ufal.mff.cuni.cz

Martin Holub  
holub@ufal.mff.cuni.cz

Charles University,  
Faculty of Mathematics and Physics,  
Institute of Formal and Applied Linguistics

# Outline

- Evaluation of binary classifiers (cntnd): ROC curve
- Support Vector Machines (SVM)

# Support Vector Machines

## Basic idea of SVM for binary classification tasks

We find a plane that separates the two classes in the feature space.

If it is not possible

- allow some training errors, or
- enrich the feature space so that finding a separating plane is possible

# Support Vector Machines

## Three key ideas

- Maximizing the margin
- Duality optimization task
- Kernels

## Key concepts needed

- Hyperplane
- Dot product
- Quadratic programming

# Hyperplane

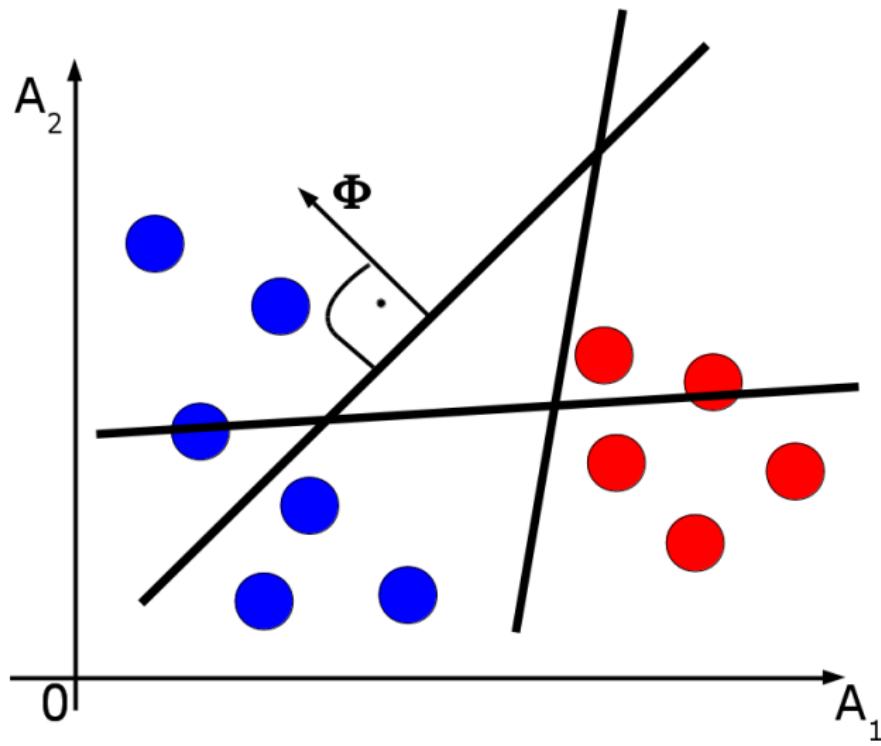
A **hyperplane** of an  $m$ -dimensional space is a **subspace with dimension  $m - 1$** .

## Mathematical definition

$$\Theta_0 + \Theta^T \mathbf{x} = 0, \text{ where } \Theta = \langle \Theta_1, \dots, \Theta_m \rangle$$

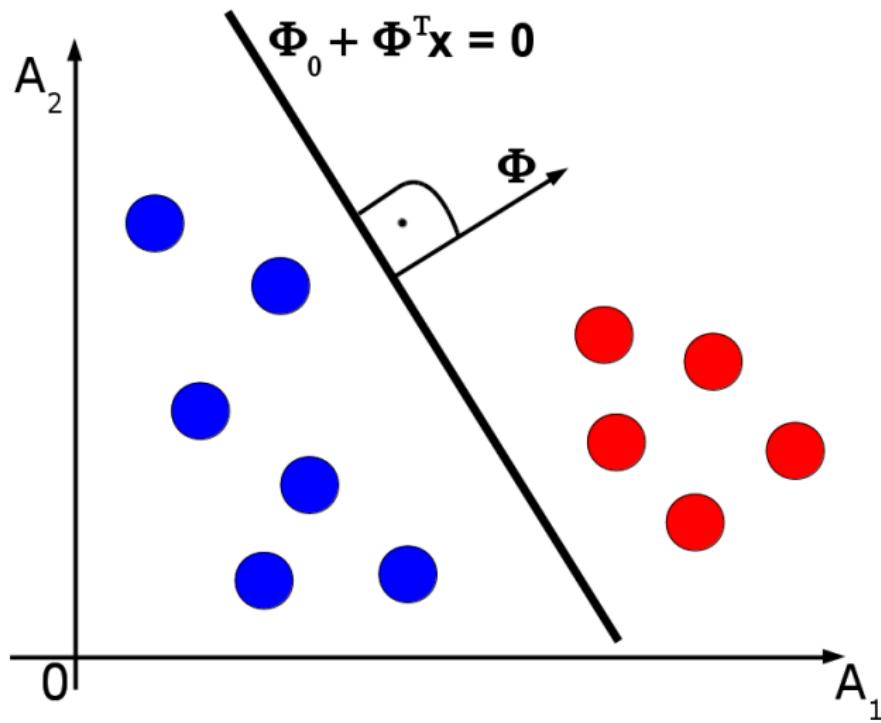
- If  $m = 2$ , a hyperplane is a line
- $\Theta$  is a normal vector
- If  $\bar{\mathbf{x}}$  satisfies the equation, then it lies on the hyperplane
- If  $\Theta_0 + \Theta^T \bar{\mathbf{x}} \neq 0$ , then  $\bar{\mathbf{x}}$  lies to one side of the hyperplane

# Hyperplane



# Hyperplane

## Separating hyperplane

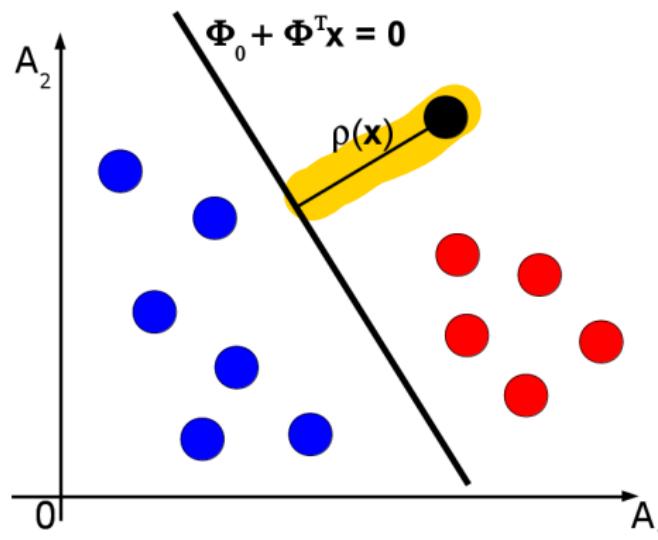


# Hyperplane

## Point-hyperplane distance

Distance of  $x$  to the hyperplane  $\Theta_0 + \Theta^T x = 0$

$$\rho(x) = \frac{|\Theta_0 + \Theta^T x|}{\|\Theta\|}$$



# Dot product

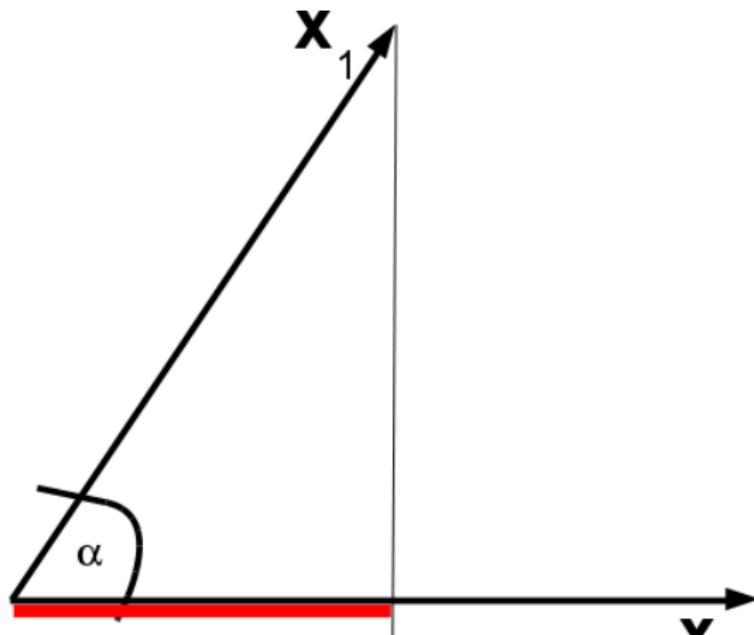
- $\mathbf{x} \in \mathcal{R}^m$
- length of  $\mathbf{x}$  (or 2-norm)  $||\mathbf{x}|| = \sqrt{\sum_{i=1}^m x_i^2}$
- dot product of two vectors  $\mathbf{x}_1, \mathbf{x}_2 \in \mathcal{R}^m$

$$\mathbf{x}_1 \cdot \mathbf{x}_2 = \sum_{i=1}^m x_{1i} x_{2i}$$

- $\mathbf{x}_1 \cdot \mathbf{x}_2 = ||\mathbf{x}_1|| \cdot ||\mathbf{x}_2|| \cdot \cos \alpha$
- geometric interpretation of  $\mathbf{x}_1 \cdot \mathbf{x}_2$ :  
the length of the projection of  $\mathbf{x}_1$  onto the unit vector  $\mathbf{x}_2$  ( $||\mathbf{x}_2|| = 1$ )
- $\mathbf{x} \cdot \mathbf{x} = ||\mathbf{x}||^2$

# Dot product

$$\|\mathbf{x}_2\| = 1$$



$$\mathbf{x}_1 \cdot \mathbf{x}_2 = \|\mathbf{x}_1\| \|\mathbf{x}_2\| \cos \alpha$$

# Quadratic programming

Quadratic programming is the problem of optimizing a quadratic function of several variables subject to linear constraints on these variables.

# Support Vector Machines

Binary classification task  $Y \in \{+1, -1\}$

$h$  has a form of

$$h(\mathbf{x}) = \text{sgn}(\Theta_0 + \boldsymbol{\Theta}^T \mathbf{x})$$

# Support Vector Machines

Binary classification task  $Y = \{+1, -1\}$

## Outline

- ① Large margin classifier (linear separability)
- ② Soft margin classifier (not linear separability)
- ③ Kernels (non-linear class boundaries)

# Support Vector Machines

Binary classification task  $Y = \{+1, -1\}$

Data set  $Data = \{\langle \mathbf{x}_i, y_i \rangle, \mathbf{x}_i \in X, y_i \in \{-1, +1\}\}$  is **linearly separable** if there exists a hyperplane so that all instances from  $Data$  are classified correctly.

# Support Vector Machines

Binary classification task  $Y = \{+1, -1\}$

Assume a hyperplane  $g$ :  $\Theta_0 + \Theta^T \mathbf{x} = 0$

- Margin of  $\mathbf{x}$  w.r.t.  $g$  is distance of  $\mathbf{x}$  to  $g$ :

$$\rho_g(\mathbf{x}) = \frac{|\Theta_0 + \Theta^T \mathbf{x}|}{\|\Theta\|}$$

- Functional margin of  $\mathbf{x}$ ,  $\langle \mathbf{x}, y \rangle \in Data$  w.r.t.  $g$  is

$$\bar{\rho}_g(\mathbf{x}, y) = y(\Theta_0 + \Theta^T \mathbf{x})$$

Is  $\mathbf{x}$  classified correctly or not?

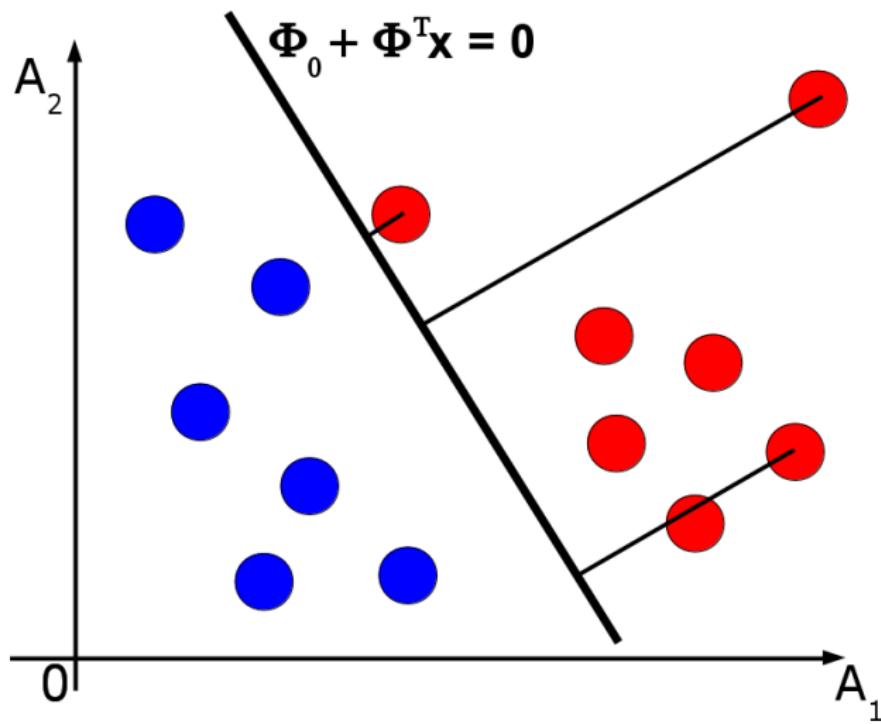
Large functional margin represents correct and confident classification.

- Geometric margin of  $\mathbf{x}$ ,  $\langle \mathbf{x}, y \rangle \in Data$  w.r.t.  $g$  is

$$\rho_g(\mathbf{x}, y) = \bar{\rho}_g(\mathbf{x}, y) / \|\Theta\|$$

i.e. functional margin scaled by  $\|\Theta\|$

## Geometric margin of $x$



# Support Vector Machines

Binary classification task  $Y = \{+1, -1\}$

- **Geometric margin** of  $Data$  w.r.t.  $g$  is

$$\rho_g(Data) = \operatorname{argmin}_{(\mathbf{x}, y) \in Data} \rho_g(\mathbf{x}, y)$$

# Support Vector Machines

Binary classification task  $Y = \{+1, -1\}$

- **Functional margin** of *Data* w.r.t.  $g$  is

$$\bar{\rho}_g(\text{Data}) = \operatorname{argmin}_{\langle \mathbf{x}, y \rangle \in \text{Data}} \bar{\rho}_g(\mathbf{x}, y)$$

Functional margin of training set is functional margin of closest points

# Large Margin Classifier

## Training data is linearly separable

We look for  $g^*$  so that

$$g^* = \operatorname{argmax}_g \rho_g(\text{Data})$$

# Large Margin Classifier

## Training data is linearly separable

$\Theta_0 + \Theta^T \mathbf{x}$  and  $k\Theta_0 + (k\Theta)^T \mathbf{x}$  define the same hyperplane, i.e. distances to separator are unchanged

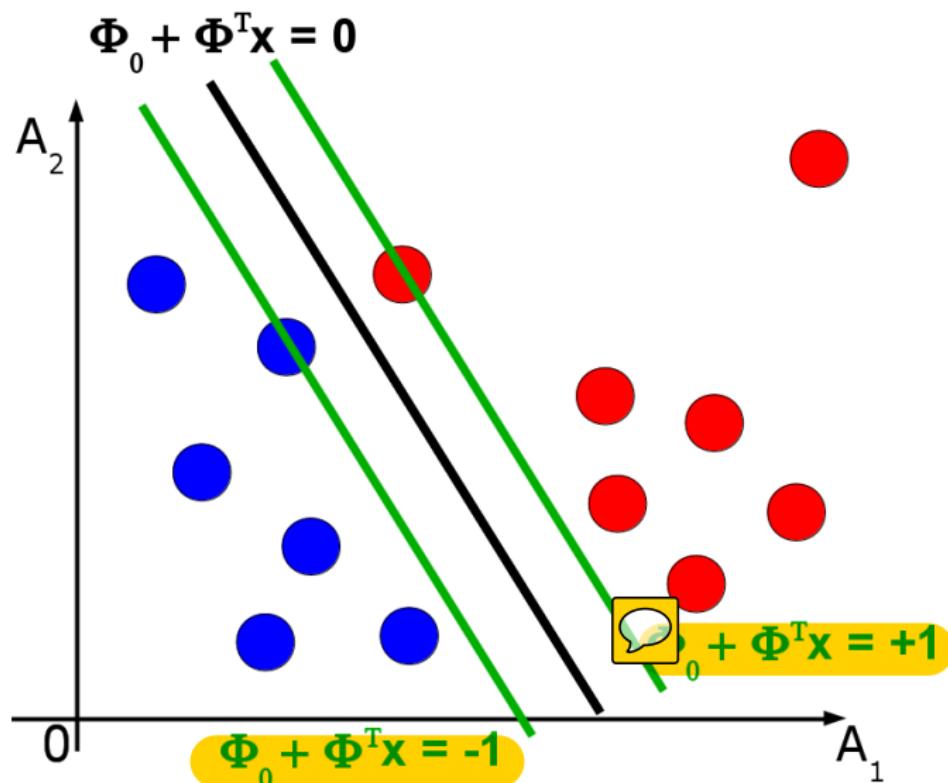
$$\frac{y_i(\Theta_0 + \Theta^T \mathbf{x}_i)}{\|\Theta\|} = \frac{y_i(k\Theta_0 + (k\Theta)^T \mathbf{x}_i)}{\|k\Theta\|}$$

Thus, we can choose  $\Theta$  so that  $\bar{\rho}_g(Data) = 1$ . Then

$$g^\star = \operatorname{argmax}_g \rho_g(Data) = \operatorname{argmax}_g \frac{1}{\|\Theta\|}$$

# Large Margin Classifier

## Training data is linearly separable



# Large Margin Classifier

## Training data is linearly separable

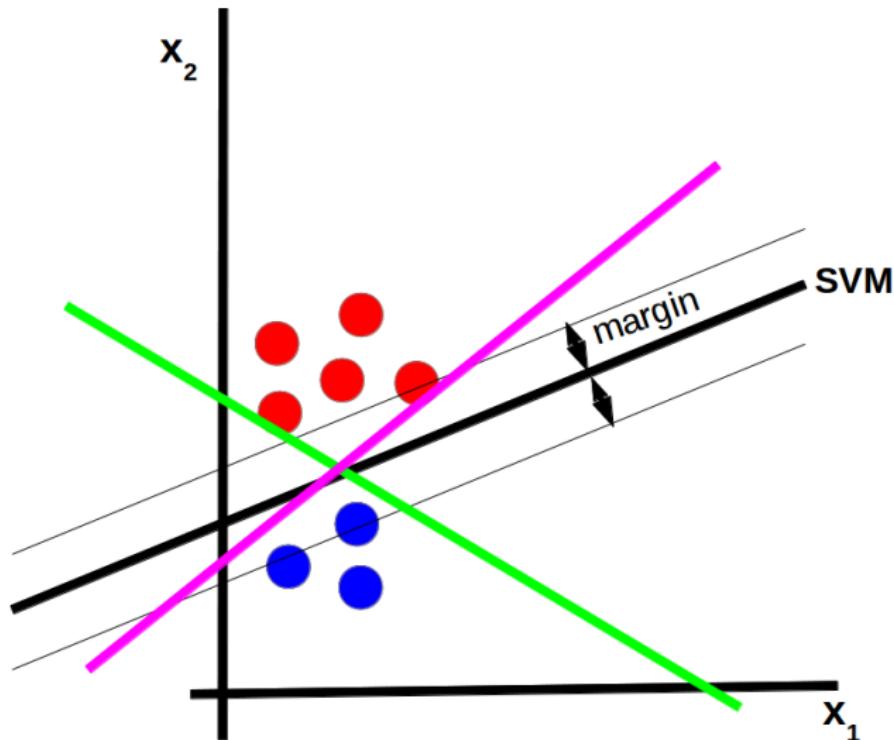
**Goal:** Orientate the separating hyperplane to be as far as possible from the closest instances of both classes.

$$\Theta^* = \operatorname{argmax}_{\Theta} \frac{1}{\|\Theta\|}$$

**Support vectors** are the instances touching the margins.

# Large Margin Classifier

## Training data is linearly separable



# Large Margin Classifier

## Training data is linearly separable

$$\Theta^* = \operatorname{argmax}_{\Theta} \frac{1}{||\Theta||} \equiv \operatorname{argmin}_{\Theta} \frac{1}{2} ||\Theta||^2$$

# Large Margin Classifier

## Training data is linearly separable

### Primal problem

Minimize

$$\frac{1}{2} \|\Theta\|^2$$

subject to constraints

$$y_i(\Theta_0 + \Theta^T x_i) \geq 1, i = 1, \dots, n$$

### Properties

- ① Convex optimization
- ② Unique solution for linearly separable training data

# Large Margin Classifier

## Training data is linearly separable

For each training example  $\langle \mathbf{x}_i, y_i \rangle$ , introduce Lagrange multiplier  $\alpha_i \geq 0$ . Let  $\boldsymbol{\alpha} = \langle \alpha_1, \dots, \alpha_n \rangle$ .

Primal Lagrangian  $L(\Theta, \Theta_0, \boldsymbol{\alpha})$  is given by

$$L(\Theta, \Theta_0, \boldsymbol{\alpha}) = \frac{1}{2} \|\Theta\|^2 - \sum_i \alpha_i [y_i(\Theta_0 + \Theta^T \mathbf{x}_i) - 1] \quad (1)$$

subject to constraints

$$\alpha_i [y_i(\Theta_0 + \Theta^T \mathbf{x}_i) - 1] = 0, i = 1, \dots, n$$

We wish to find the  $\Theta$  and  $\Theta_0$  which minimizes and the  $\boldsymbol{\alpha}$  which maximizes  $L$ .

# Large Margin Classifier

## Training data is linearly separable

1. Minimize  $L$  w.r.t.  $\Theta$

Thus differentiate  $L$  w.r.t.  $\Theta$  and  $\frac{\partial L}{\partial \Theta} = 0$

It gives

$$\Theta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (2)$$

2. Minimize  $L$  w.r.t.  $\Theta_0$

Thus differentiate  $L$  w.r.t.  $\Theta_0$  and  $\frac{\partial L}{\partial \Theta_0} = 0$

It gives

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (3)$$

# Large Margin Classifier

## Training data is linearly separable

3. Substitute (2) into the primal form (1). Then

$$L(\Theta, \Theta_0, \alpha) = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j$$

subject to  $\alpha_i \geq 0$  and  $\sum_i \alpha_i y_i = 0, i = 1 \dots n$

# Large Margin Classifier

## Training data is linearly separable

4. Solve the dual problem, i.e. maximize a quadratic function. Do quadratic programming.

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

subject to  $\alpha_i \geq 0$  and  $\sum_i \alpha_i y_i = 0, i = 1 \dots n$

5. Get  $\alpha^*$

6. Then  $\Theta^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$

# Large Margin Classifier

## Training data is linearly separable

- $\Theta^*$  is the solution to the primal problem
- $\alpha^*$  is the solution to the dual problem

$\Theta^*$  and  $\alpha^*$  satisfy the Karush-Kuhn-Tucker conditions where one of them is so called **Karush-Kuhn-Tucker dual complementarity**:

$$\alpha_i * (1 - y_i(\Theta_0 + \Theta^T \mathbf{x}_i)) = 0$$

- $y_i(\Theta_0 + \Theta^T \mathbf{x}_i) \neq 1$  ( $\mathbf{x}_i$  is not support vector)  $\Rightarrow \alpha_i = 0$
- $\alpha_i \neq 0 \Rightarrow y_i(\Theta_0 + \Theta^T \mathbf{x}_i) = 1$  ( $\mathbf{x}_i$  is support vector)

I.e., finding  $\Theta$  is equivalent to finding support vectors and their weights

# Large Margin Classifier

## Training data is linearly separable

**Prediction** for a new instance  $\mathbf{x}$

$$h(\mathbf{x}) = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \mathbf{x} + \Theta_0\right)$$

- similarity between  $\mathbf{x}$  and support vector  $\mathbf{x}_i$ : a support vector that is more similar contributes more to the classification
- support vector that is more important, i.e. has larger  $\alpha_i$ , contributes more to the classification
- if  $y_i$  is positive, than the contribution is positive, otherwise negative

# Soft Margin Classifier

## Training data is not linearly separable

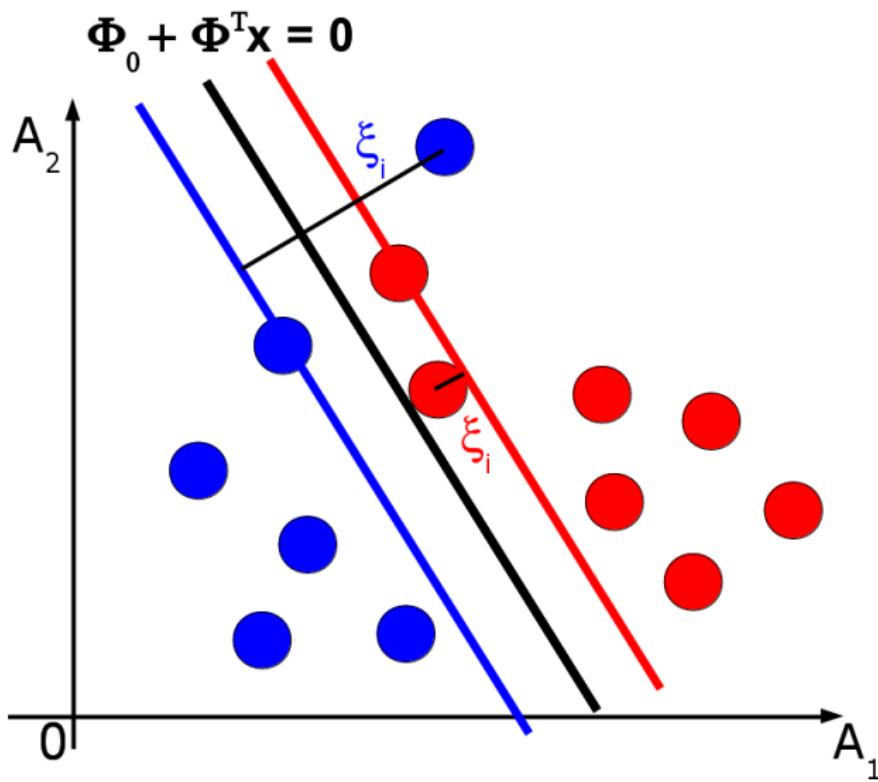
In a real problem it is unlikely that a line will exactly separate the data – even if a curved decision boundary is possible. So exactly separating the data is probably not desirable – if the data has noise and outliers, a smooth decision boundary that ignores a few data points is better than one that loops around the outliers.

Thus

minimize  $\|\Theta\|^2$  AND the number of training mistakes

# Soft Margin Classifier

Training data is not linearly separable



# Soft Margin Classifier

Training data is not linearly separable

Introducing slack variables  $\xi_i \geq 0$

- $\xi_i = 0$  if  $x_i$  is correctly classified
- $\xi_i$  is distance to "its supporting hyperplane" otherwise
  - $0 < \xi_i \leq 1/\|\Theta\|$ : margin violation
  - $\xi_i > 1/\|\Theta\|$ : misclassification

# Soft Margin Classifier

Training data is not linearly separable

## Primal problem

Minimize

$$\frac{1}{2} \|\Theta\|^2 + C \sum_{i=1}^n \xi_i$$

subject to constraint

$$y_i(\Theta_0 + \Theta^T \mathbf{x}_i) \geq 1 - \xi_i, i = 1, \dots, n$$

- $C \geq 0$  trade-off parameter
  - small  $C \Rightarrow$  large margin
  - large  $C \Rightarrow$  narrow margin

# Soft Margin Classifier

Training data is not linearly separable

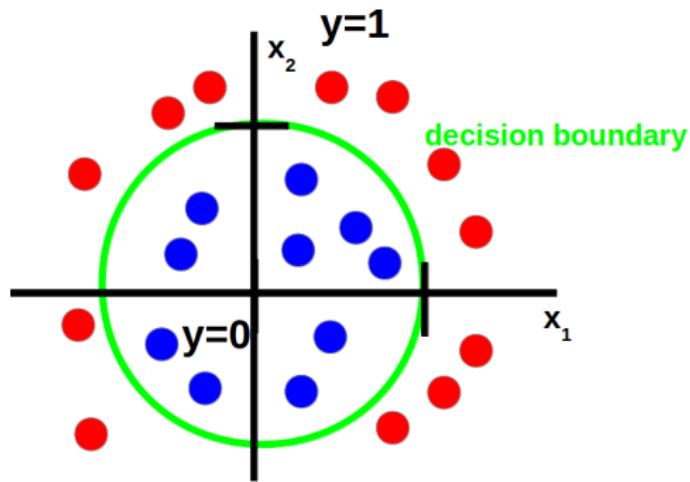
- Do quadratic programming as for Large Margin Classifier
- **Prediction** for a new instance  $\mathbf{x}$

$$h(\mathbf{x}) = \text{sgn}\left(\sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \cdot \mathbf{x} + \Theta_0\right)$$

# Support Vector Machines

## Non-linear boundary

If the examples are separated by a nonlinear region



# Support Vector Machines

## Non-linear boundary

### Recall polynomial regression

Polynomial regression is an extension of linear regression where the relationship between features and target value is modelled as a  $d$ -th order polynomial.

#### Simple regression

$$y = \Theta_0 + \Theta_1 x_1$$

#### Polynomial regression

$$y = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_1^2 + \dots + \Theta_d x_1^d$$

It is still a linear model with features  
 $A_1, A_1^2, \dots, A_1^d$ .

This defines a feature mapping  $\phi(x_1) = [x_1, x_1^2, \dots, x_1^d]$

# Support Vector Machines

## Kernels

### Idea

Apply Large/Soft margin classifier not to the orginal features but to the features obtained by the feature mapping  $\phi$

$$\phi(\mathbf{x}) : \mathcal{R}^m \rightarrow \mathcal{F}$$

Large/Soft margin classifier uses dot product  $\mathbf{x}_i \cdot \mathbf{x}_j$ . Now, replace it with  $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$ .

# Support Vector Machines

## Kernels

However, finding  $\phi$  could be expensive.

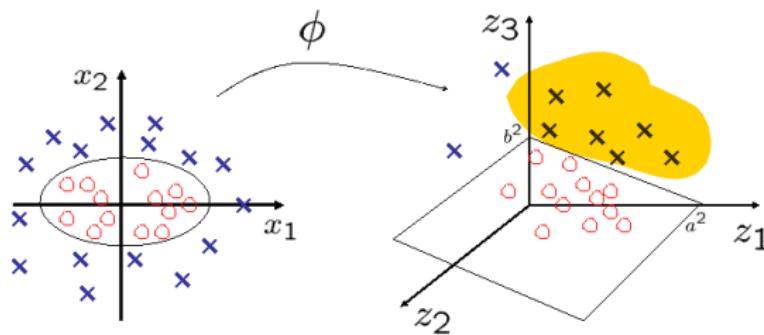
### Kernel trick

- No need to know what  $\phi$  is and what the feature space is, i.e. no need to explicitly map the data to the feature space
- Define a kernel function  $K : \mathcal{R}^m \times \mathcal{R}^m \rightarrow \mathcal{R}$
- Replace the dot product  $\mathbf{x}_i \cdot \mathbf{x}_j$  with a Kernel function  $K(\mathbf{x}_i, \mathbf{x}_j)$ :

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

# Support Vector Machines

## Kernels



$$\phi : (x_1, x_2) \longrightarrow (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

$$\left(\frac{x_1}{a}\right)^2 + \left(\frac{x_2}{b}\right)^2 = 1 \longrightarrow \frac{z_1}{a^2} + \frac{z_3}{b^2} = 1$$

Source: <http://omega.albany.edu:8008/machine-learning-dir/notes-dir/ker1/ker1-1.html>

# Support Vector Machines

## Kernels

### Common kernel functions

- Linear:  $K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i^T \mathbf{x}_j$
- Polynomial:  $K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i^T \mathbf{x}_j + c)^d$
- Radial basis function:  $K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma(||\mathbf{x}_i - \mathbf{x}_j||))^2$
- Sigmoid:  $K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\gamma \mathbf{x}_i^T \mathbf{x}_j + c)$ , where  $\tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$

# Support Vector Machines

## Multiclass classification tasks

### One-to-one

- Train  $\binom{K}{2}$  SVM binary classifiers
- Classify  $\mathbf{x}$  using each of the  $\binom{K}{2}$  classifiers. The instance is assigned to the class which is the most frequent class assigned in the pairwise classification.

# Support Vector Machines

## Multiclass classification tasks

### One-to-all

- Train  $K$  SVM binary classifiers. Each of them, doing classification of  $k$ -th class (+1) to the others (-1), is characterized by the hypothesis parameters  $\Theta_k, k = 1, \dots, K$
- The instance  $\mathbf{x}$  is assigned to the class  $k^* = \max_k \Theta_k^T \mathbf{x}$

# Evaluation of binary classifiers

## Sensitivity vs. specificity

### Confusion matrix

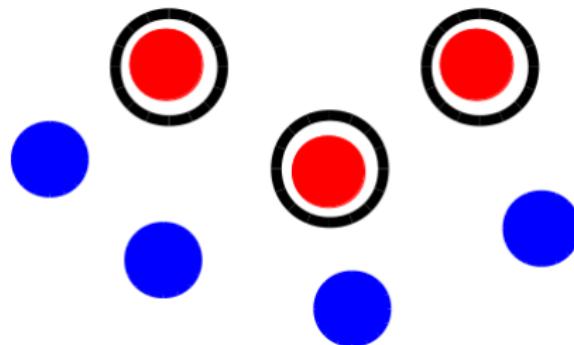
		Predicted class	
		Positive	Negative
True class	Positive	True Positive (TP)	False Negative (FN)
	Negative	False Positive (FP)	True Negative (TN)

Measure	Formula
Precision	$TP/(TP+FP)$
Recall/Sensitivity	$TP/(TP+FN)$
Specificity	$TN/(TN+FP)$
Accuracy	$(TP+TN)/(TP+FP+TN+FN)$

# Evaluation of binary classifiers

## Sensitivity vs. specificity

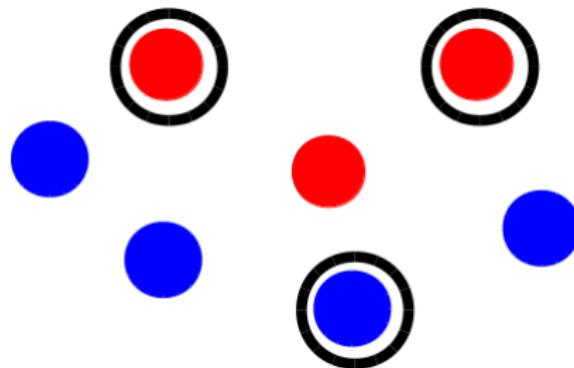
Perfect classifier



# Evaluation of binary classifiers

## Sensitivity vs. specificity

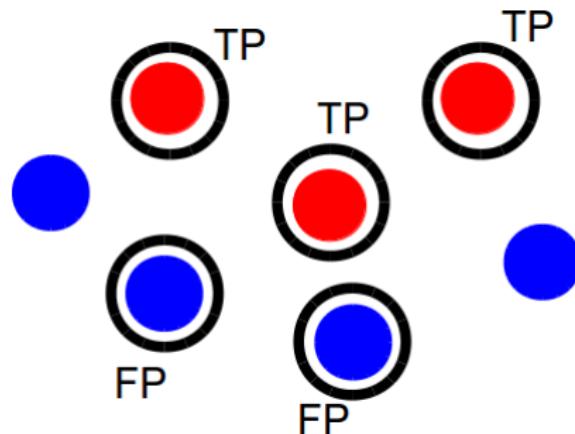
Reality



# Evaluation of binary classifiers

## Sensitivity vs. specificity

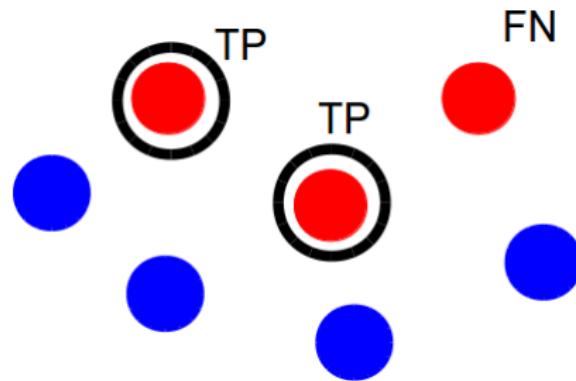
100% sensitive classifier



# Evaluation of binary classifiers

## Sensitivity vs. specificity

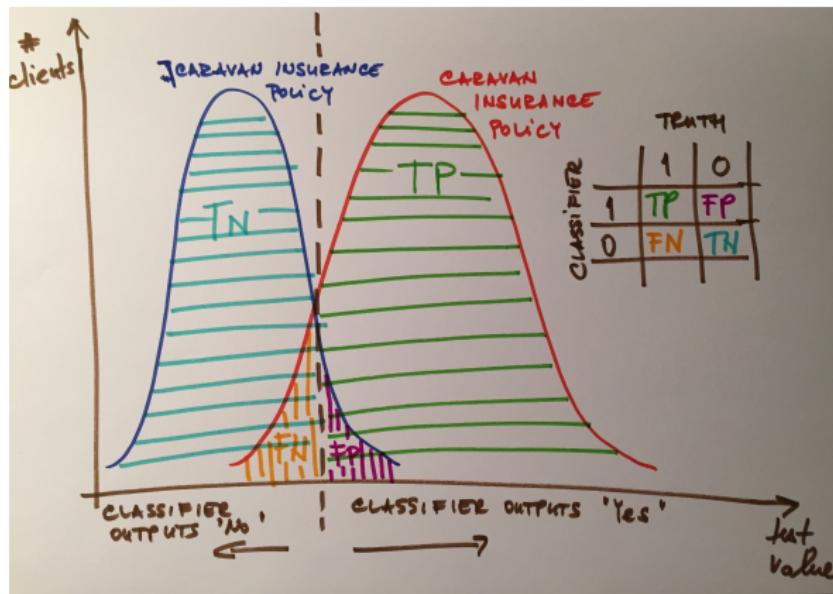
100% specific classifier



# Evaluation of binary classifiers

## Sensitivity vs. specificity

### Sensitivity vs. specificity



# Evaluation of binary classifiers

## ROC curve

