

CS 1571 Introduction to AI
Lecture 26

Binary classification:
Support Vector Machines

Milos Hauskrecht
milos@cs.pitt.edu
5329 Sennott Square

CS 1571 Intro to AI

Supervised learning

Data: $D = \{D_1, D_2, \dots, D_n\}$ a set of n examples

$$D_i = \langle \mathbf{x}_i, y_i \rangle$$

$\mathbf{x}_i = (x_{i,1}, x_{i,2}, \dots, x_{i,d})$ is an input vector of size d

y_i is the desired output (given by a teacher)

Objective: learn the mapping $f : X \rightarrow Y$

$$\text{s.t. } y_i \approx f(\mathbf{x}_i) \text{ for all } i = 1, \dots, n$$

- **Regression:** Y is **continuous**
Example: earnings, product orders \rightarrow company stock price
- **Classification:** Y is **discrete**
Example: handwritten digit in binary form \rightarrow digit label

CS 1571 Intro to AI

Discriminant functions: review

- A classification model is typically defined using
 - discriminant functions
- Idea:
 - For each class i define a function $g_i(\mathbf{x})$ mapping $X \rightarrow \mathcal{R}$
 - When the decision on input \mathbf{x} should be made choose the class with the highest value of $g_i(\mathbf{x})$

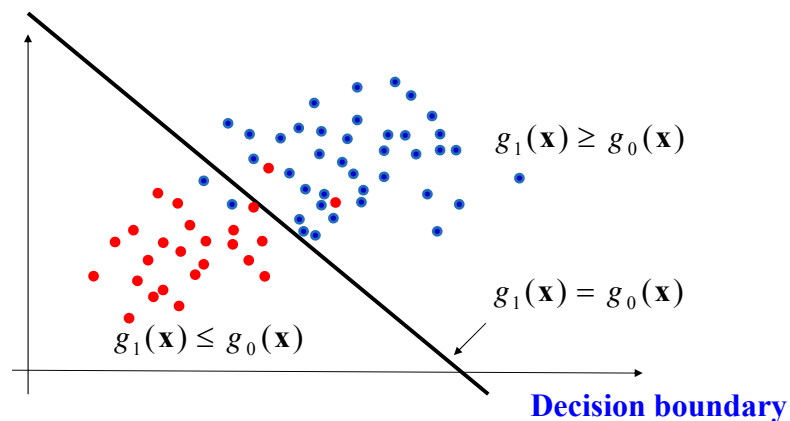
$$\text{class} = \arg \max_i g_i(x)$$

- Works for binary and multi-class classification

CS 1571 Introduction to AI

Discriminant functions: review

- Assume a binary classification problem with classes 0 and 1
- Discriminant functions $g_0(\mathbf{x})$ and $g_1(\mathbf{x})$



CS 1571 Introduction to AI

Logistic regression model: review

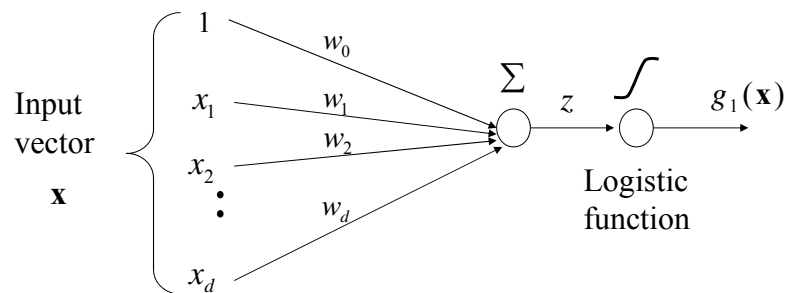
- **Model for binary (2 class) classification**

- **Defined by discriminant functions:**

$$g_1(\mathbf{x}) = g(\mathbf{w}^T \mathbf{x} + w_0) \quad g_0(\mathbf{x}) = 1 - g(\mathbf{w}^T \mathbf{x} + w_0)$$

- **where**

$g(z) = 1/(1 + e^{-z})$ is a logistic function



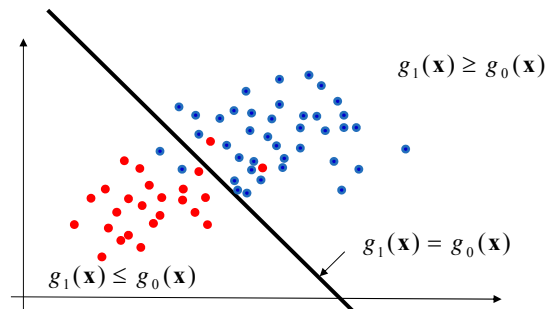
CS 2750 Machine Learning

Logistic regression model. Decision boundary

- **Logistic regression model defines a linear decision boundary**

$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$

Example: 2 classes (blue and red points)



CS 1571 Introduction to AI

Decision boundary

An alternative way to define discriminant functions with a linear decision boundary

Class 1:

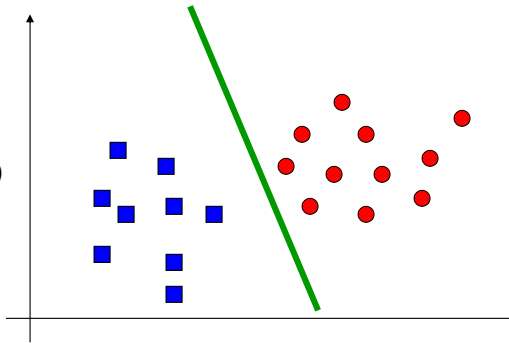
$$g_1(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$$

Class -1:

$$g_{-1}(\mathbf{x}) = -(\mathbf{w}^T \mathbf{x} + w_0)$$

Decision boundary:

$$\mathbf{w}^T \mathbf{x} + w_0 = 0$$



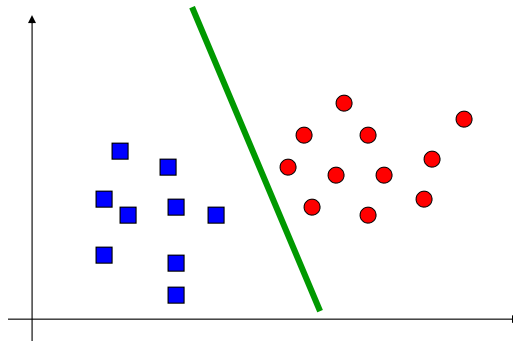
CS 2750 Machine Learning

Linearly separable classes

Linearly separable classes:

There is a **hyperplane** $\mathbf{w}^T \mathbf{x} + w_0 = 0$
that separates training instances with no error

Class (+1)
$\mathbf{w}^T \mathbf{x} + w_0 > 0$
Class (-1)
$\mathbf{w}^T \mathbf{x} + w_0 < 0$



CS 2750 Machine Learning

Learning linearly separable sets

Finding weights for linearly separable classes:

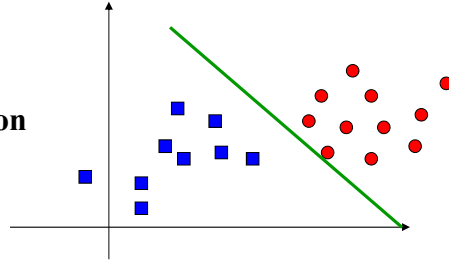
- **Linear program (LP) solution**
- It finds weights that satisfy the following constraints:

$$\mathbf{w}^T \mathbf{x}_i + w_0 \geq 0 \quad \text{For all } i, \text{ such that } y_i = +1$$

$$\mathbf{w}^T \mathbf{x}_i + w_0 \leq 0 \quad \text{For all } i, \text{ such that } y_i = -1$$

$$\text{Together: } y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 0$$

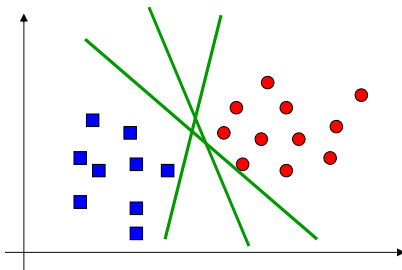
Property: if there is a hyperplane separating the examples, the linear program finds the solution



CS 2750 Machine Learning

Optimal separating hyperplane

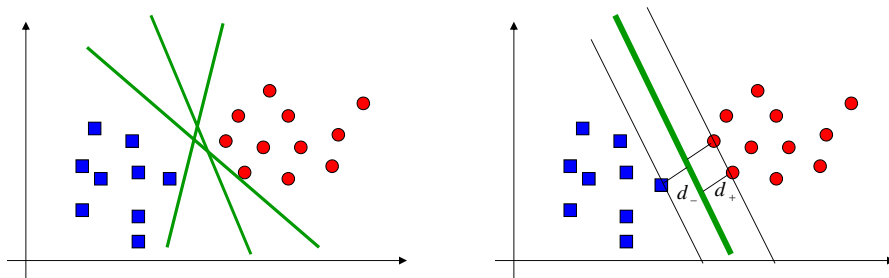
- **Problem:**
- There are multiple hyperplanes that separate the data points
- Which one to choose?



CS 2750 Machine Learning

Optimal separating hyperplane

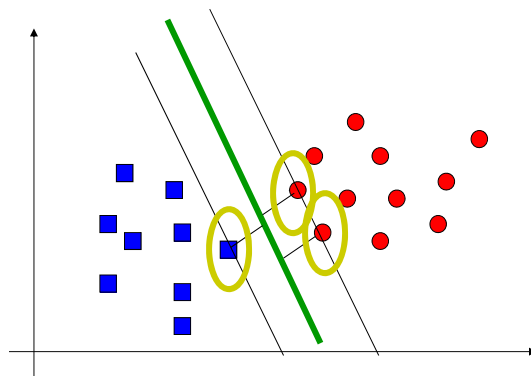
- **Problem:** multiple hyperplanes that separate the data exists
 - Which one to choose?
- **Maximum margin** choice: maximum distance of $d_+ + d_-$
 - where d_+ is the shortest distance of a positive example from the hyperplane (similarly d_- for negative examples)



CS 2750 Machine Learning

Maximum margin hyperplane

- For the maximum margin hyperplane only examples on the margin matter (only these affect the distances)
- These are called **support vectors**



CS 2750 Machine Learning

Finding maximum margin hyperplanes

- **Assume** that examples in the training set are (\mathbf{x}_i, y_i) such that $y_i \in \{+1, -1\}$
- **Assume** that all data satisfy:

$$\mathbf{w}^T \mathbf{x}_i + w_0 \geq 1 \quad \text{for} \quad y_i = +1$$

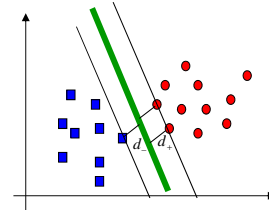
$$\mathbf{w}^T \mathbf{x}_i + w_0 \leq -1 \quad \text{for} \quad y_i = -1$$

- The inequalities can be combined as:

$$y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \geq 0 \quad \text{for all } i$$

- Equalities define two hyperplanes:

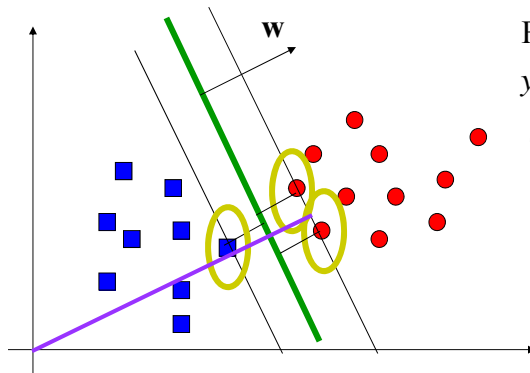
$$\mathbf{w}^T \mathbf{x}_i + w_0 = 1 \quad \mathbf{w}^T \mathbf{x}_i + w_0 = -1$$



CS 2750 Machine Learning

Finding the maximum margin hyperplane

- **Geometrical margin:** $\rho_{\mathbf{w}, w_0}(\mathbf{x}, y) = y(\mathbf{w}^T \mathbf{x} + w_0) / \|\mathbf{w}\|_{L_2}$
 - measures the distance of a point \mathbf{x} from the hyperplane
 - \mathbf{w} - normal to the hyperplane $\|\cdot\|_{L_2}$ - Euclidean norm



For points satisfying:
 $y_i(\mathbf{w}^T \mathbf{x}_i + w_0) - 1 = 0$

The distance is $\frac{1}{\|\mathbf{w}\|_{L_2}}$

Width of the margin:

$$d_+ + d_- = \frac{2}{\|\mathbf{w}\|_{L_2}}$$

CS 2750 Machine Learning

Maximum margin hyperplane

- We want to maximize $d_+ + d_- = \frac{2}{\|\mathbf{w}\|_{L2}}$
- We do it by minimizing

$$\|\mathbf{w}\|_{L2}^2 / 2 = \mathbf{w}^T \mathbf{w} / 2$$

\mathbf{w}, w_0 - variables

- But we also need to enforce the constraints on points:

$$[y_i(\mathbf{w}^T \mathbf{x} + w_0) - 1] \geq 0$$

CS 2750 Machine Learning

Maximum margin hyperplane

- **Solution:** Incorporate constraints into the optimization
- **Optimization problem** (Lagrangian)

$$J(\mathbf{w}, w_0, \alpha) = \|\mathbf{w}\|^2 / 2 - \sum_{i=1}^n \alpha_i [y_i(\mathbf{w}^T \mathbf{x} + w_0) - 1]$$

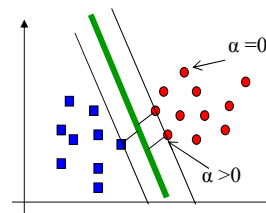
$\alpha_i \geq 0$ - **Lagrange multipliers**

- **Minimize** with respect to \mathbf{w}, w_0 (primal variables)
- **Maximize** with respect to α (dual variables)

What happens to α :

if $y_i(\mathbf{w}^T \mathbf{x} + w_0) - 1 > 0 \Rightarrow \alpha_i \rightarrow 0$
 else $\Rightarrow \alpha_i > 0$

Active constraint



CS 2750 Machine Learning

Max margin hyperplane solution

- Set derivatives to 0 (Kuhn-Tucker conditions)

$$\nabla_{\mathbf{w}} J(\mathbf{w}, w_0, \alpha) = \mathbf{w} - \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i = \bar{\mathbf{0}}$$

$$\frac{\partial J(\mathbf{w}, w_0, \alpha)}{\partial w_0} = -\sum_{i=1}^n \alpha_i y_i = 0$$

- Now we need to solve for Lagrange parameters (Wolfe dual)

$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j) \quad \leftarrow \text{maximize}$$

Subject to constraints

$$\alpha_i \geq 0 \quad \text{for all } i, \quad \text{and} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

- Quadratic optimization problem:** solution $\hat{\alpha}_i$ for all i

CS 2750 Machine Learning

Maximum margin solution

- The resulting parameter vector $\hat{\mathbf{w}}$ can be expressed as:

$$\hat{\mathbf{w}} = \sum_{i=1}^n \hat{\alpha}_i y_i \mathbf{x}_i \quad \hat{\alpha}_i \text{ is the solution of the optimization}$$

- The parameter w_0 is obtained from $\hat{\alpha}_i [y_i (\hat{\mathbf{w}}^T \mathbf{x}_i + w_0) - 1] = 0$

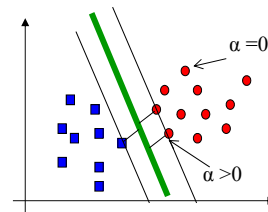
Solution properties

- $\hat{\alpha}_i = 0$ for all points that are not on the margin

- The decision boundary:**

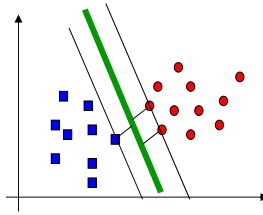
$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0 = 0$$

The decision boundary defined by support vectors only



CS 2750 Machine Learning

Support vector machines



- The decision boundary:

$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0$$

- Classification decision:

$$\hat{y} = \text{sign} \left[\sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0 \right]$$

CS 2750 Machine Learning

Support vector machines: solution property

- Decision boundary defined by the set of support vectors SV and their alpha values
 - Support vectors = a subset of datapoints in the training data that define the margin

$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0$$

- Classification decision:

$$\hat{y} = \text{sign} \left[\sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0 \right]$$

- **Note** that we do not have to explicitly compute $\hat{\mathbf{w}}$
 - This will be important for the nonlinear (kernel) case

CS 2750 Machine Learning

Support vector machines: inner product

- Decision on a new \mathbf{x} depends on the **inner product between two examples**

- The decision boundary:**

$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0$$

- Classification decision:**

$$\hat{y} = \text{sign} \left[\sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0 \right]$$

- Similarly, the optimization depends on $(\mathbf{x}_i^T \mathbf{x}_j)$

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

CS 2750 Machine Learning

Inner product of two vectors

- The decision boundary for the SVM and its optimization depend on inner product of two datapoints (vectors):**

$$\mathbf{x}_i^T \mathbf{x}_j$$

$$\mathbf{x}_i = \begin{pmatrix} 2 \\ 5 \\ 6 \end{pmatrix} \quad \mathbf{x}_j = \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix}$$

$$(\mathbf{x}_i^T \mathbf{x}_j) = \quad ?$$

CS 2750 Machine Learning

Inner product of two vectors

- The decision boundary for the SVM and its optimization depend on the inner product of two data points (vectors):

$$\mathbf{x}_i^T \mathbf{x}_j$$

$$\mathbf{x}_i = \begin{pmatrix} 2 \\ 5 \\ 6 \end{pmatrix} \quad \mathbf{x}_j = \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix}$$

$$(\mathbf{x}_i^T \mathbf{x}) = (2 \quad 5 \quad 6) * \begin{pmatrix} 2 \\ 3 \\ 1 \end{pmatrix} = 2 * 2 + 5 * 3 + 6 * 1 = 25$$

CS 2750 Machine Learning

Inner product of two vectors

- The decision boundary for the SVM and its optimization depend on the inner product of two data points (vectors):

$$\mathbf{x}_i^T \mathbf{x}_j$$

- The inner product is also equal

$$(\mathbf{x}_i^T \mathbf{x}) = \|\mathbf{x}_i\| * \|\mathbf{x}_i\| \cos \theta$$

If the angle in between them is 0 then: $(\mathbf{x}_i^T \mathbf{x}) = \|\mathbf{x}_i\| * \|\mathbf{x}_i\|$

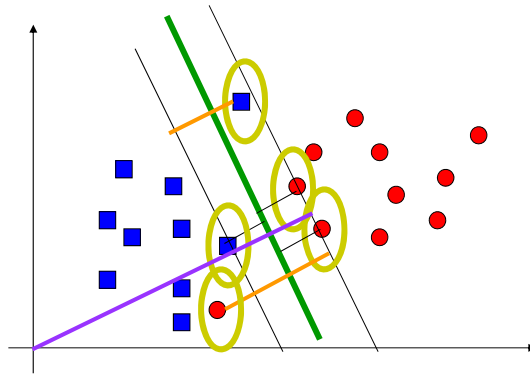
If the angle between them is 90 then: $(\mathbf{x}_i^T \mathbf{x}) = 0$

The inner product measures how similar the two vectors are

CS 2750 Machine Learning

Extension to a linearly non-separable case

- **Idea:** Allow some flexibility on crossing the separating hyperplane



CS 2750 Machine Learning

Extension to the linearly non-separable case

- Relax constraints with variables $\xi_i \geq 0$
$$\mathbf{w}^T \mathbf{x}_i + w_0 \geq 1 - \xi_i \quad \text{for} \quad y_i = +1$$
$$\mathbf{w}^T \mathbf{x}_i + w_0 \leq -1 + \xi_i \quad \text{for} \quad y_i = -1$$
- Error occurs if $\xi_i \geq 1$, $\sum_{i=1}^n \xi_i$ is the upper bound on the number of errors
- Introduce a penalty for the errors

$$\text{minimize} \quad \|\mathbf{w}\|^2 / 2 + C \sum_{i=1}^n \xi_i$$

Subject to constraints

C – set by a user, larger C leads to a larger penalty for an error

CS 2750 Machine Learning

Support vector machines: solution

- The solution of the linearly non-separable case has the same properties as the linearly separable case.
 - The decision boundary is defined only by a set of support vectors (points that are on the margin or that cross the margin)
 - The decision boundary and the optimization can be expressed in terms of the inner product in between pairs of examples

$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0$$

$$\hat{y} = \text{sign} [\hat{\mathbf{w}}^T \mathbf{x} + w_0] = \text{sign} \left[\sum_{i \in SV} \hat{\alpha}_i y_i (\mathbf{x}_i^T \mathbf{x}) + w_0 \right]$$

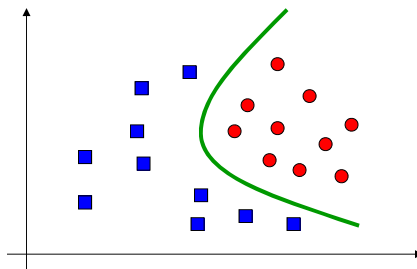
$$J(\alpha) = \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j=1}^n \alpha_i \alpha_j y_i y_j (\mathbf{x}_i^T \mathbf{x}_j)$$

CS 2750 Machine Learning

Nonlinear decision boundary

So far we have seen how to learn a linear decision boundary

- **But what if the linear decision boundary is not good.**
- **How can we learn a non-linear decision boundaries with the SVM?**



CS 2750 Machine Learning

Nonlinear decision boundary

- The non-linear case can be handled by using a set of features. Essentially we map input vectors to (larger) feature vectors

$$\mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x})$$

- Note that feature expansions are typically high dimensional
 - Examples: polynomial expansions
- Given the nonlinear feature mappings, we can use the linear SVM on the expanded feature vectors

$$(\mathbf{x}^T \mathbf{x}') \longrightarrow \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\varphi}(\mathbf{x}')$$

- **Kernel function**

$$K(\mathbf{x}, \mathbf{x}') = \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\varphi}(\mathbf{x}')$$

CS 2750 Machine Learning

Nonlinear case

- The linear case requires to compute $(\mathbf{x}^T \mathbf{x}')$
- The non-linear case can be handled by using a set of features. Essentially we map input vectors to (larger) feature vectors

$$\mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x})$$

- Note that feature expansions are typically high dimensional
 - Examples: polynomial expansions

- Now we can use SVM formalism on feature vectors

$$(\mathbf{x}^T \mathbf{x}') \longrightarrow \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\varphi}(\mathbf{x}')$$

- **Kernel function**

$$K(\mathbf{x}, \mathbf{x}') = \boldsymbol{\varphi}(\mathbf{x})^T \boldsymbol{\varphi}(\mathbf{x}')$$

CS 2750 Machine Learning

Support vector machines: solution for nonlinear decision boundaries

- The decision boundary:

$$\hat{\mathbf{w}}^T \mathbf{x} + w_0 = \sum_{i \in SV} \hat{\alpha}_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0$$

- Classification:

$$\hat{y} = \text{sign} [\hat{\mathbf{w}}^T \mathbf{x} + w_0] = \text{sign} \left[\sum_{i \in SV} \hat{\alpha}_i y_i K(\mathbf{x}_i, \mathbf{x}) + w_0 \right]$$

- Decision on a new \mathbf{x} requires to compute **the kernel function defining the similarity between the examples**
- Similarly, the optimization depends on the kernel

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

CS 2750 Machine Learning

Kernel trick

The non-linear case maps input vectors to (larger) feature space

$$\mathbf{x} \rightarrow \boldsymbol{\phi}(\mathbf{x})$$

- Note that feature expansions are typically high dimensional
 - Examples: polynomial expansions
- **Kernel function** defines the inner product in the expanded high dimensional feature vectors and let us use the SVM
$$(\mathbf{x}^T \mathbf{x}') \longrightarrow K(\mathbf{x}, \mathbf{x}') = \boldsymbol{\phi}(\mathbf{x})^T \boldsymbol{\phi}(\mathbf{x}')$$
- **Problem:** after expansion we need to perform inner products in a very high dimensional space
- **Kernel trick:**
 - If we choose the kernel function wisely we can compute linear separation in the high dimensional feature space implicitly by working in the original input space !!!!

CS 2750 Machine Learning

Kernel function example

- Assume $\mathbf{x} = [x_1, x_2]^T$ and a feature mapping that maps the input into a quadratic feature set

$$\mathbf{x} \rightarrow \boldsymbol{\varphi}(\mathbf{x}) = [x_1^2, x_2^2, \sqrt{2}x_1x_2, \sqrt{2}x_1, \sqrt{2}x_2, 1]^T$$

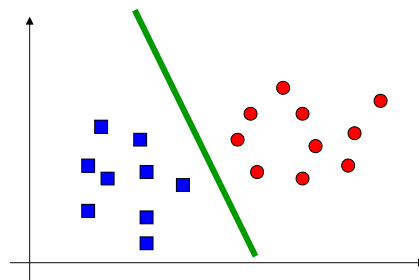
- Kernel function for the feature space:

$$\begin{aligned} K(\mathbf{x}', \mathbf{x}) &= \boldsymbol{\varphi}(\mathbf{x}')^T \boldsymbol{\varphi}(\mathbf{x}) \\ &= x_1^2 x_1'^2 + x_2^2 x_2'^2 + 2x_1x_2x_1'x_2' + 2x_1x_1' + 2x_2x_2' + 1 \\ &= (x_1x_1' + x_2x_2' + 1)^2 \\ &= (1 + (\mathbf{x}^T \mathbf{x}'))^2 \end{aligned}$$

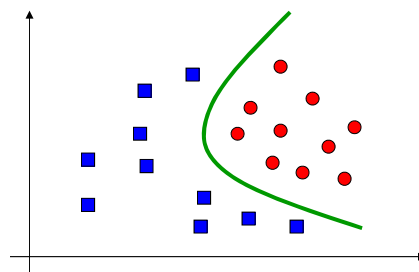
- The computation of the linear separation in the higher dimensional space is performed implicitly in the original input space

CS 2750 Machine Learning

Kernel function example



Linear separator
in the feature space



Non-linear separator
in the input space

CS 2750 Machine Learning

Kernel functions

- **Linear kernel**

$$K(\mathbf{x}, \mathbf{x}') = \mathbf{x}^T \mathbf{x}'$$

- **Polynomial kernel**

$$K(\mathbf{x}, \mathbf{x}') = \left[1 + \mathbf{x}^T \mathbf{x}'\right]^k$$

- **Radial basis kernel**

$$K(\mathbf{x}, \mathbf{x}') = \exp\left[-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2\right]$$

CS 2750 Machine Learning

Kernels

- **Kernel – similarity between pairs of objects**
- Kernels can be defined for more complex objects:
 - Strings
 - Graphs
 - Images

CS 2750 Machine Learning