

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

- Support Vector Machines

Support Vector Machines

Key points & concepts

Key points

- 1 Maximizing the margin
- 2 Duality optimization task
- 3 Kernels

Key concepts

- 1 Hyperplane
- 2 Dot product
- 3 Quadratic programming

Support Vector Machines

Key idea for binary classification

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

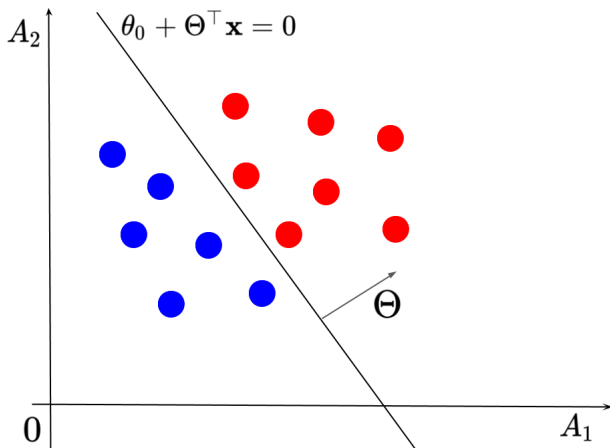
If it is not possible

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

A hyperplane is an $(n - 1)$ -dimensional linear subspace of an n dimensional vector space

$$\theta_0 + \theta_1 x_1 + \dots + \theta_m x_m = 0$$

Hyperplane



Support Vector Machines

Classification using hyperplane

Recall the classification rule using $\theta_0 + \theta_1 x_1 + \dots + \theta_m x_m = 0$

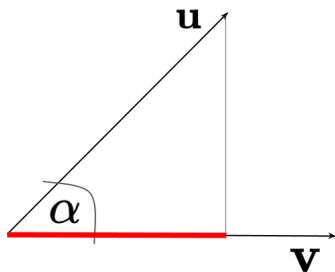
$$f(\mathbf{x}) = \begin{cases} 1 & \text{if } \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m \geq 0 \\ 0 & \text{if } \theta_0 + \theta_1 x_1 + \dots + \theta_m x_m < 0 \end{cases}$$

We can use the classification rule for different values of threshold.
Here threshold = 0.

Support Vector Machines

Dot product

- $\mathbf{u} = \langle u_1, \dots, u_m \rangle$, $\mathbf{v} = \langle v_1, \dots, v_m \rangle$
- algebraic definition $\mathbf{u} \cdot \mathbf{v} = u_1 v_1 + \dots + u_m v_m$
- geometric definition $\mathbf{u} \cdot \mathbf{v} = \|\mathbf{u}\| \cdot \|\mathbf{v}\| \cdot \cos \alpha$



$$\|\mathbf{v}\| = 1 \rightarrow \mathbf{u} \cdot \mathbf{v} = \|\mathbf{u}\| \cdot \cos \alpha$$

Support Vector Machines

Quadratic programming

Quadratic programming optimizes a quadratic objective function subject to constraints.

i.e. minimize/maximize

$$f(\mathbf{x}), \mathbf{x} \in \mathcal{D}$$

subject to

$$g_i(\mathbf{x}) \leq 0, i = 1, \dots, G$$

$$h_j(\mathbf{x}) = 0, j = 1, \dots, H$$

Support Vector Machines

Margins

Training examples $D = \{\langle \mathbf{x}_i, y_i \rangle, \mathbf{x}_i \in X, y_i \in \{-1, +1\}\}$

Hyperplane $g: \theta_0 + \Theta^\top \mathbf{x} = 0$

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

$$\rho_g(\mathbf{x}) = \frac{|\theta_0 + \Theta^\top \mathbf{x}|}{\|\Theta\|}$$

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

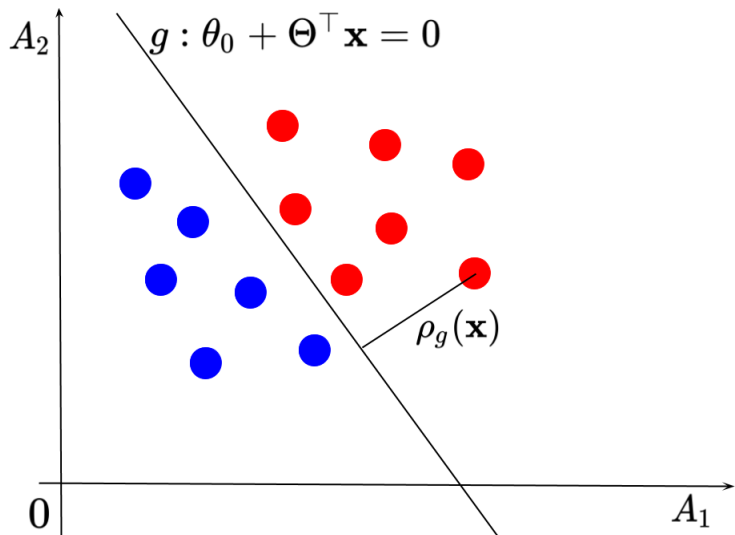
$$\bar{\rho}_g(\mathbf{x}, y) = y(\Theta_0 + \Theta^\top \mathbf{x})$$

- **Geometric margin of $\langle \mathbf{x}, y \rangle$** w.r.t. g is functional margin scaled by $\|\Theta\|$

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

Support Vector Machines

Margins



Support Vector Machines

Margins

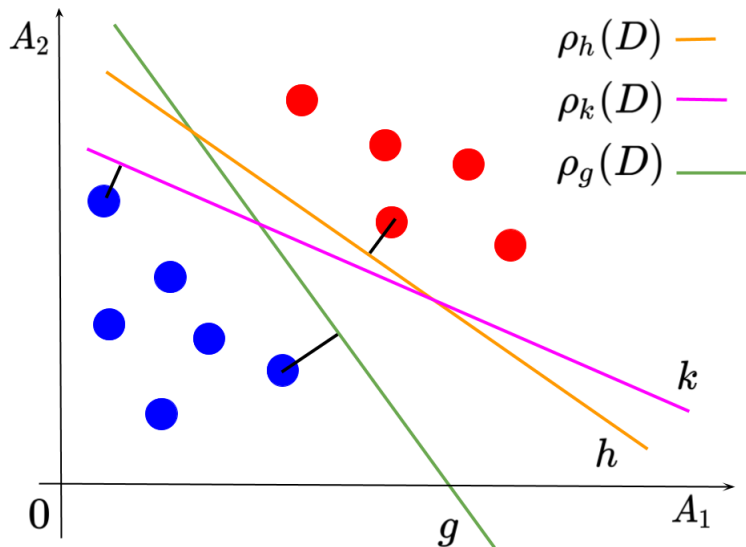
- **Functional margin of D w.r.t. g**

$$\bar{\rho}_g(D) = \min_{\langle \mathbf{x}, y \rangle \in D} \bar{\rho}_g(\mathbf{x}, y)$$

- **Geometric margin of D w.r.t. g**

$$\rho_g(D) = \min_{\langle \mathbf{x}, y \rangle \in D} \rho_g(\mathbf{x}, y)$$

Support Vector Machines Margins



Support Vector Machines

Margins

Data set $D = \{\langle \mathbf{x}_i, y_i \rangle, \mathbf{x}_i \in X, y_i \in \{-1, +1\}\}$ is **linearly separable** if there exists a hyperplane $g : \theta_0 + \Theta^\top \mathbf{x} = 0$ that separates the two classes completely, i.e.

$$\forall \langle \mathbf{x}, y \rangle \in D : \quad \overline{\rho}_g(\mathbf{x}, y) > 0$$

Support Vector Machines

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

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

Large Margin Classifier

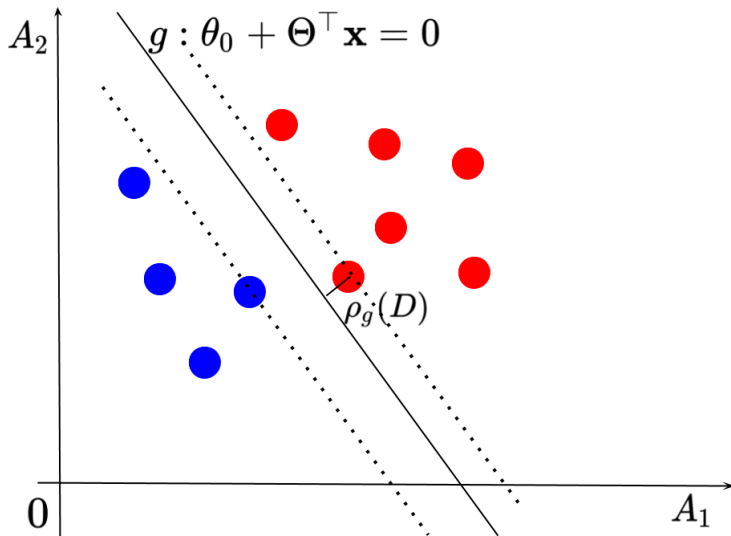
Training data is linearly separable

Optimization task

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

Large Margin Classifier

Training data is linearly separable



Large Margin Classifier

Training data is linearly separable

$\Theta_0 + \Theta^\top \mathbf{x}$ and $k\Theta_0 + (k\Theta)^\top \mathbf{x}$ define the same hyperplane.

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

Therefore we can scale Θ so that $\bar{\rho}_g(D) = 1$.

Then

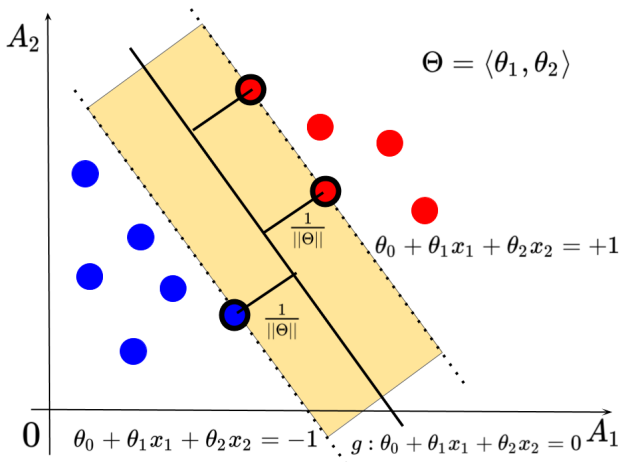
$$g^* = \operatorname{argmax}_g \rho_g(D) = \max_{\Theta} \frac{1}{\|\Theta\|}$$

Large Margin Classifier

Training data is linearly separable

Supporting hyperplanes $\theta_0 + \Theta^\top \mathbf{x} = -1$ and $\theta_0 + \Theta^\top \mathbf{x} = +1$.

Support vectors are the instances touching the supporting hyperplanes.



Large Margin Classifier

Training data is linearly separable

Optimization task

Minimize

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

subject to

$$y_i(\theta_0 + \Theta^\top \mathbf{x}_i) \geq 1, i = 1, \dots, n$$

quadratic function and linear constraints \rightarrow quadratic programming

Large Margin Classifier

Quadratic programming

Use the method of Lagrangian multipliers to find a minimum of a function subject to constraints.

Define the Lagrangian function $\mathcal{L}_P(\Theta, \theta_0, \alpha)$

$$\mathcal{L}_P(\Theta, \theta_0, \alpha) = \frac{1}{2} \|\Theta\|^2 - \sum_{i=1}^n \alpha_i (y_i (\theta_0 + \Theta^\top \mathbf{x}_i) - 1) \quad (1)$$

$$\alpha_i \geq 0, i = 1, \dots, n$$

Large Margin Classifier

Quadratic programming

Primal Lagrangian problem

$$\text{minimize}_{\Theta, \theta_0} \text{maximize}_{\alpha} \mathcal{L}_P(\Theta, \theta_0, \alpha) \quad (2)$$

Swap minimize and maximize (see Slater's condition)

$$\text{maximize}_{\alpha} \text{minimize}_{\Theta, \theta_0} \mathcal{L}_P(\Theta, \theta_0, \alpha) \quad (3)$$

subject to

$$\alpha_i \geq 0, i = 1, \dots, n$$

1. Minimize \mathcal{L}_P w.r.t. Θ

Therefore differentiate \mathcal{L}_P w.r.t. Θ and $\frac{\partial \mathcal{L}}{\partial \Theta} = 0$. It gives

$$\Theta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (4)$$

2. Minimize \mathcal{L}_P w.r.t. θ_0

Therefore differentiate \mathcal{L}_P w.r.t. θ_0 and $\frac{\partial \mathcal{L}}{\partial \theta_0} = 0$. It gives $\sum_{i=1}^n \alpha_i y_i = 0$

Large Margin Classifier

Quadratic programming

3. Substitute (4) into the primal form (1) and solve **Wolfe dual optimization problem**

$$\text{maximize}_{\alpha} \mathcal{L}_D(\alpha) = \text{maximize}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \cdot \mathbf{x}_j$$

subject to

$$\alpha_i \geq 0, i = 1 \dots n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

4. Get α^*
5. Get $\Theta^* = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$. Assume that \mathbf{x}_i is a support vector. Then $1 = y_i(\theta_0^* + \Theta^{*\top} \mathbf{x}_i) \rightarrow \theta_0^* = y_i - \Theta^{*\top} \mathbf{x}_i$

Large Margin Classifier

Quadratic programming

- α^* is the solution to the dual problem
- Θ^* is the solution to the primal problem
- the solutions α^* and Θ^* must satisfy the Karush-Kuhn-Tucker conditions where one of them is *KKT dual complementarity*:

$$\alpha_i^* \cdot (1 - y_i(\theta_0^* + \Theta^{*\top} \mathbf{x}_i)) = 0$$

- $y_i(\theta_0^* + \Theta^{*\top} \mathbf{x}_i) \neq 1$, i.e., \mathbf{x}_i is not support vector $\Rightarrow \alpha_i^* = 0$
- $\alpha_i^* \neq 0 \Rightarrow y_i(\theta_0^* + \Theta^{*\top} \mathbf{x}_i) = 1$, i.e., \mathbf{x}_i is support vector

I.e., finding Θ^* is equivalent to finding support vectors and their weights

Large Margin Classifier

Prediction

Prediction for a new instance \mathbf{x}

$$f(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i \cdot \mathbf{x} + \theta_0^*\right)$$

- 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 α_i , contributes more to the classification
- if y_i is positive, then the contribution is positive, otherwise negative

Soft Margin Classifier

Training data is not linearly separable

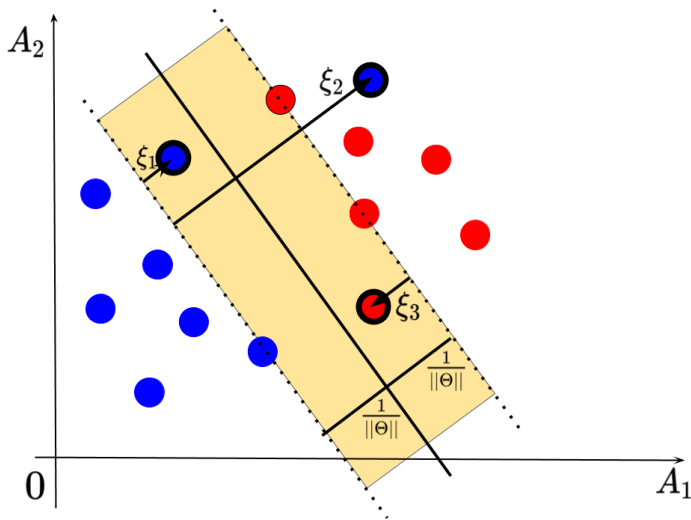
In a real problem it is unlikely that a hyperplane will exactly separate the data – even if a curved decision boundary is possible. So exactly separating the data is probably not desirable – if the data has noise and outliers, a smooth decision boundary that ignores a few data points is better than one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

Therefore

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

Soft Margin Classifier

Training data is not linearly separable



Soft Margin Classifier

Slack variables

$$\xi_i \geq 0$$

- $\xi_i = 0$ if \mathbf{x}_i is correctly classified
- ξ_i is distance to y_i 's supporting hyperplane otherwise
 - margin violation – $0 < \xi_i \leq 1/\|\Theta\|$, see ξ_1, ξ_3 above
 - misclassification – $\xi_i > 1/\|\Theta\|$, see ξ_2 above

Soft Margin Classifier

Optimization problem

Minimize

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

subject to

$$\xi_i \geq 0, y_i(\theta_0 + \Theta^\top \mathbf{x}_i) \geq 1 - \xi_i, i = 1, \dots, n$$

$C \geq 0$ trade-off parameter

- small $C \Rightarrow$ large margin
relaxed model; misclassifications are not penalized
- large $C \Rightarrow$ narrow margin
misclassifications are penalized strongly
the model will not generalize much

Soft Margin Classifier

Quadratic programming

For each constraint $y_i(\theta_0 + \Theta^\top \mathbf{x}_i) \geq 1 - \xi_i$
introduce Lagrange multiplier $\alpha_i \geq 0$.

Let $\boldsymbol{\alpha} = \langle \alpha_1, \dots, \alpha_n \rangle$.

Primal Lagrangian $\mathcal{L}_P(\Theta, \theta_0, \xi, \boldsymbol{\alpha})$ is given by

$$\mathcal{L}_P(\Theta, \theta_0, \xi, \boldsymbol{\alpha}) = \frac{1}{2} \|\Theta\|^2 + C \sum_{i=1}^n \xi_i - \sum_{i=1}^n \alpha_i (y_i(\theta_0 + \Theta^\top \mathbf{x}_i) - 1 + \xi_i) \quad (5)$$

Soft Margin Classifier

Quadratic programming

Primal Lagrangian problem

$$\text{minimize}_{\Theta, \theta_0, \xi} \text{maximize}_{\alpha} \mathcal{L}_P(\Theta, \theta_0, \xi, \alpha) \quad (6)$$

$$\text{maximize}_{\alpha} \text{minimize}_{\Theta, \theta_0, \xi} \mathcal{L}_P(\Theta, \theta_0, \xi, \alpha) \quad (7)$$

subject to

$$\alpha_i \geq 0, \xi_i \geq 0, i = 1, \dots, n$$

1. Minimize \mathcal{L}_P w.r.t. Θ . Therefore $\frac{\partial \mathcal{L}}{\partial \Theta} = 0$. It gives

$$\Theta = \sum_{i=1}^n \alpha_i y_i \mathbf{x}_i \quad (8)$$

2. Minimize \mathcal{L}_P w.r.t. θ_0 . Therefore $\frac{\partial \mathcal{L}}{\partial \theta_0} = 0$. It gives

$$\sum_{i=1}^n \alpha_i y_i = 0 \quad (9)$$

Soft Margin Classifier

Quadratic programming

3. Minimize \mathcal{L}_P w.r.t. ξ . Therefore $\frac{\partial L}{\partial \xi} = 0$. It gives

$$C\xi - \alpha = 0 \quad (10)$$

4. Substitute (8) into the primal form (5) and solve **Wolfe dual opt. problem**

$$\text{maximize}_{\alpha} \mathcal{L}_D(\alpha) = \text{maximize}_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1, j=1}^n \alpha_i \alpha_j y_i y_j \mathbf{x}_i \mathbf{x}_j$$

subject to

$$\alpha_i \geq 0, \quad i = 1, \dots, n$$

$$\sum_{i=1}^n \alpha_i y_i = 0$$

5. Get α^*

6. Get Θ^*

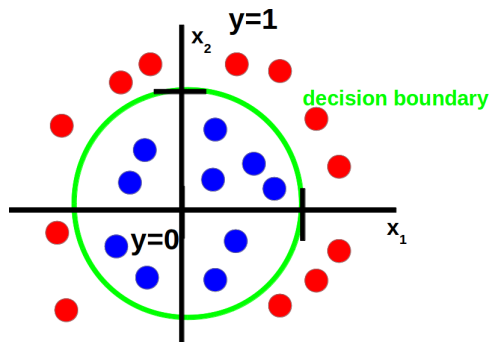
Soft Margin Classifier Prediction

Prediction for a new instance \mathbf{x}

$$f(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i \cdot \mathbf{x} + \theta_0^*\right)$$

Non-linear boundary

If the examples are separated by a nonlinear region



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

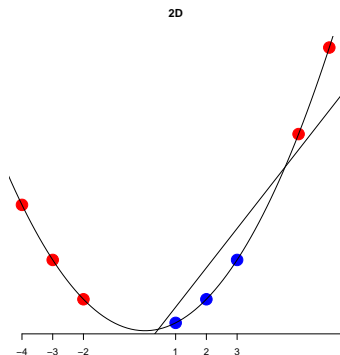
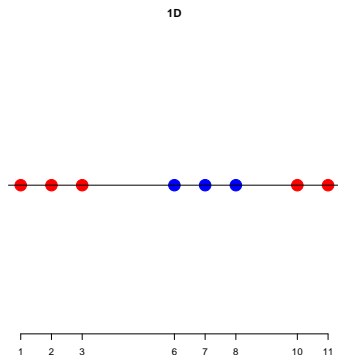
$$y = \Theta_0 + \Theta_1 x_1 + \Theta_2 x_1^2 + \dots + \Theta_d x_1^d$$

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

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

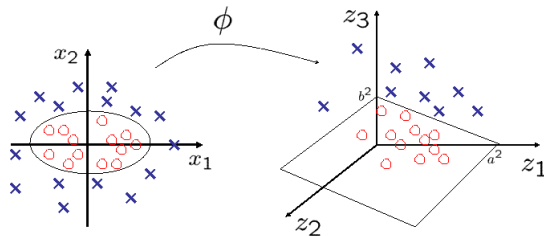
Non-linear boundary

- for example $\phi(x_1) = [x_1 - 5, (x_1 - 5)^2]$



Idea

- Apply Large/Soft margin classifier not to the original features but to the features obtained by the feature mapping $\phi(\mathbf{x}) : \mathcal{R}^m \rightarrow \mathcal{F}$
- Large/Soft margin classifier uses dot product $\mathbf{x}_i \cdot \mathbf{x}_j$.
Replace it with $\phi(\mathbf{x}_i) \cdot \phi(\mathbf{x}_j)$.



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

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

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

However, finding ϕ could be expensive.

Kernel trick

- No need to know what ϕ 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 $K(\mathbf{x}_i, \mathbf{x}_j)$:

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

Kernels

Prediction

Prediction for a new instance \mathbf{x}

$$f(\mathbf{x}) = \text{sign}\left(\sum_{i=1}^n \alpha_i^* y_i K(\mathbf{x}_i, \mathbf{x}) + \theta_0^*\right)$$

Common Kernel functions

- **Linear**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \mathbf{x}_i \cdot \mathbf{x}_j$$

- **Polynomial**

$$K(\mathbf{x}_i, \mathbf{x}_j) = (\gamma \mathbf{x}_i \cdot \mathbf{x}_j + c)^d$$

- smaller degree can generalize better

- higher degree can fit (only) training data better

- **Radial basis function**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \exp(-\gamma(\|\mathbf{x}_i - \mathbf{x}_j\|^2))$$

- very robust

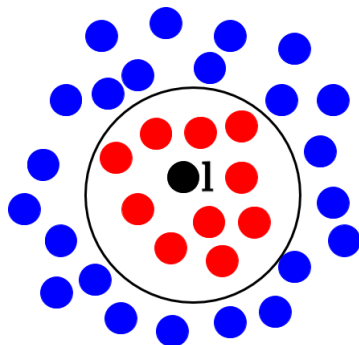
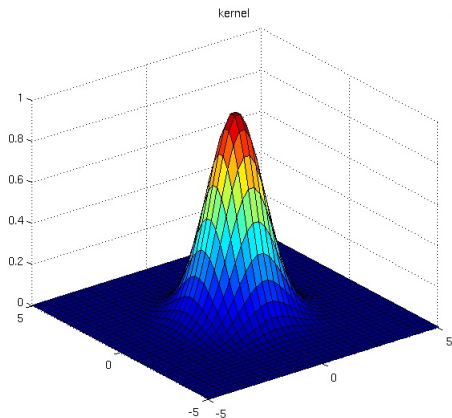
- use it when polynomial kernel is weak to fit data

- **Sigmoid**

$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\gamma \mathbf{x}_i \cdot \mathbf{x}_j + c), \text{ where } \tanh(x) = \frac{\exp(x) - \exp(-x)}{\exp(x) + \exp(-x)}$$

Radial Basis Function Kernel

- $K(\mathbf{x}_i, \mathbf{l}) = e^{-\gamma \|\mathbf{x}_i - \mathbf{l}\|^2}$



Source: <http://www.cs.toronto.edu/~duvenaud/cookbook/index.html>

Regularized logistic regression

$$f(\mathbf{x}) = \frac{1}{1 + e^{-\Theta^\top \mathbf{x}}}$$

$$L(\Theta) = - \sum_{i=1}^n y_i \log P(y_i | \mathbf{x}_i; \Theta) + (1 - y_i) \log(1 - P(y_i | \mathbf{x}_i; \Theta))$$

$$\Theta_R^* = \operatorname{argmin}_{\Theta} [L(\Theta) + \lambda \cdot \text{penalty}(\Theta)]$$

Regularized logistic regression

Ridge regression

$$\begin{aligned}\Theta_R^* &= \operatorname{argmin}_{\Theta} - \left[\sum_{i=1}^n y_i \log(f(\mathbf{x}_i)) + (1 - y_i) \log(1 - f(\mathbf{x}_i)) \right] + \lambda \sum_{j=1}^m \theta_j^2 = \\ &= \operatorname{argmin}_{\Theta} \left[\sum_{i=1}^n y_i (-\log(f(\mathbf{x}_i))) + (1 - y_i) (-\log(1 - f(\mathbf{x}_i))) \right] + \lambda \sum_{j=1}^m \theta_j^2 = \\ &= \operatorname{argmin}_{\Theta} \left[\sum_{i=1}^n y_i L_1(\Theta) + (1 - y_i) L_0(\Theta) \right] + \lambda \sum_{j=1}^m \theta_j^2\end{aligned}$$

Regularized logistic regression

Ridge regression

Since

$$\mathbf{A} + \lambda \mathbf{B} \equiv \mathbf{C} \mathbf{A} + \mathbf{B}, \mathbf{C} = \frac{1}{\lambda}$$

then

$$\Theta_R^* = \operatorname{argmin}_{\Theta} \left[\sum_{j=1}^m \theta_j^2 + \mathbf{C} \left[\sum_{i=1}^n y_i L_1(\Theta) + (1 - y_i) L_0(\Theta) \right] \right]$$

where

$$L_1(\Theta) = -\log \frac{1}{1 + e^{-\Theta^T \mathbf{x}}}$$

$$L_0(\Theta) = -\log \left(1 - \frac{1}{1 + e^{-\Theta^T \mathbf{x}}} \right)$$

Regularized logistic regression

Ridge regression and Soft Margin Classifier

$$\Theta_R^* = \operatorname{argmin}_{\Theta} \left[\sum_{j=1}^m \theta_j^2 + C \sum_{i=1}^n \log(1 + e^{-\bar{y}_i \Theta^T \mathbf{x}_i}) \right]$$

where

$$\bar{y}_i = \begin{cases} -1 & \text{if } y_i = 0 \\ +1 & \text{if } y_i = 1 \end{cases}$$

Soft Margin Classifier is equivalent to the regularization problem:

$$\Theta^* = \operatorname{argmin}_{\Theta} \sum_{j=1}^m \theta_j^2 + C \sum_{i=1}^n \xi_i$$

Soft Margin Classifier

$$\Theta^* = \operatorname{argmin}_{\Theta} \sum_{j=1}^m \theta_j^2 + C \sum_{i=1}^n \xi_i$$

$\xi_i \geq 0$ is equivalent to $\xi_i = \max(0, 1 - y_i \Theta^\top \mathbf{x}_i)$, i.e.

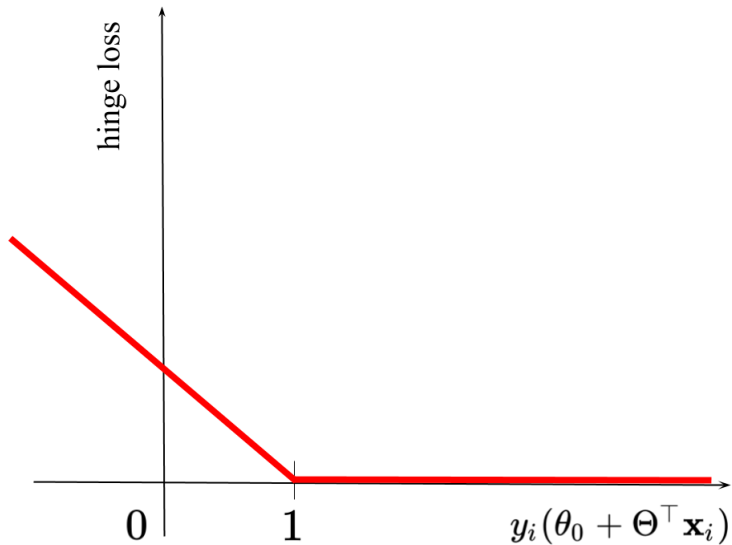
$$\Theta^* = \operatorname{argmin}_{\Theta} \left[\sum_{j=1}^m \theta_j^2 + C \sum_{i=1}^n \max(0, 1 - y_i \Theta^\top \mathbf{x}_i) \right]$$

s.t. $\Theta^\top \mathbf{x}_i \geq 1 - \xi_i$ if $y_i = +1$ and $\Theta^\top \mathbf{x}_i \leq -1 + \xi_i$ if $y_i = -1$

Hinge loss = $\max(0, 1 - y_i \Theta^\top \mathbf{x})$

- 1 $y_i \Theta^\top \mathbf{x}_i > 1$: no contribution to loss
- 2 $y_i \Theta^\top \mathbf{x}_i = 1$: no contribution to loss
- 3 $y_i \Theta^\top \mathbf{x}_i < 1$: contribution to loss

Soft Margin Classifier



Summary of Examination Requirements

- Hyperplane, margin, functional margin, geometric margin of example and data set
- Large margin classifier
linearly separable data, supporting hyperplanes, support vectors, optimization task, prediction function
- Soft margin classifier
not linearly separable data, supporting hyperplanes, support vectors, slack variables, optimization task, hyperparameter C , prediction function
- Kernel trick
feature mapping, Kernel functions, prediction function