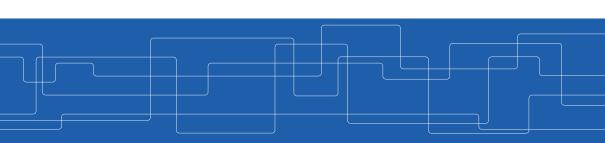


Classification with Separating Hyperplanes





Linear separation

Structural Risk Minimization

Support Vector Machines

Kernels

Non-separable Classes



Linear separation

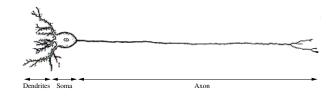
Structural Risk Minimization

Support Vector Machines

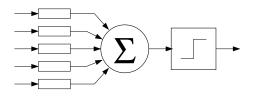
Kernels

Non-separable Classes





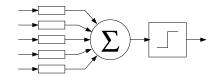
Neuron caricature, "artificial neuron"



- Weighted input signals
- Summing
- ► Thresholded output



What can a single "artificial neuron" compute?



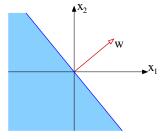
- \vec{x} Input in vector format
- w Weights in vector format
- y Output

$$y = \operatorname{sign}\left(\sum_{i} x_{i} w_{i}\right)$$



$$y = \operatorname{sign}\left(\sum_{i} x_{i} w_{i}\right)$$

Geometrical interpretation



Variable threshold \equiv Not anchored to origin Common trick: treat the variable threshold as an extra weight



Training a linear separator

What does learning mean here?

Learning means finding the best weights w_i

Several efficient algorithms exist:

- Error correctionPerceptron Learning
- Error/Loss minimization
 Delta Rule
 Logistic Regression



Perceptron Learning

- Incremental learning
- Weights only change when the output is wrong
- ▶ Update rule: $w_i \leftarrow w_i + \eta(t y)x_i$
- Always converges if the problem is solvable

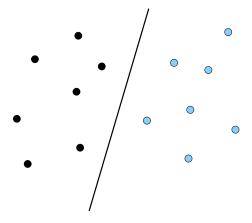


Delta Rule

- Minimize $\sum_{k \in \mathcal{D}} ||t_k \vec{w}^T \vec{x}_k||^2$
- ▶ Incremental version: Stochastic Gradient Descent For each sample: $\vec{w} \leftarrow \vec{w} \eta \ \mathrm{grad}_{\vec{w}} ||t \vec{w}^T \vec{x}||^2$
- $\qquad \qquad \mathbf{w}_i \leftarrow \mathbf{w}_i + \eta (t \vec{\mathbf{w}}^T \vec{\mathbf{x}}) \mathbf{x}_i$

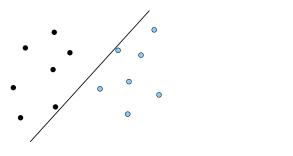


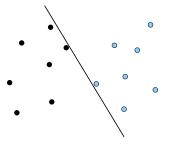
Linear Separation





Many acceptable solutions \rightarrow bad generalization





Structural Risk



Linear separation

Structural Risk Minimization

Support Vector Machines

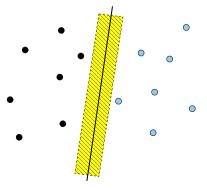
Kernels

Non-separable Classes



Hyperplane with margins

Training data points are at least a distance d from the plane



Less arbitrariness \rightarrow better generalization



- Wide margins restrict the possible hyperplanes to choose from
- Less risk to choose a bad hyperplane by accident
- Reduced risk for bad generalization

Minimization of the structural risk \equiv maximization of the margin

Out of all hyperplanes which solve the problem the one with widest margin will probably generalize best



Mathematical Formulation

Separating Hyperplane

$$\vec{\mathbf{w}}^T \vec{\mathbf{x}} = 0$$

Hyperplane with a margin

$$\vec{w}^T \vec{x} \ge 1$$
 when $t = 1$
 $\vec{w}^T \vec{x} \le -1$ when $t = -1$

Combined

$$t\vec{w}^T\vec{x} \ge 1$$

How wide is the margin?

1. Select two points, \vec{p} and \vec{q} , on the two margins:

$$\vec{\mathbf{w}}^T \vec{\mathbf{p}} = 1 \qquad \vec{\mathbf{w}}^T \vec{\mathbf{q}} = -1$$

2. Distance between \vec{p} and \vec{q} along \vec{w} .

$$2\mathbf{d} = \frac{\vec{\mathbf{w}}^T}{||\vec{\mathbf{w}}||}(\vec{\mathbf{p}} - \vec{\mathbf{q}})$$

3. Simplify:

$$2\mathbf{d} = \frac{\vec{\mathbf{w}}^T \vec{\mathbf{p}} - \vec{\mathbf{w}}^T \vec{\mathbf{q}}}{||\vec{\mathbf{w}}||} = \frac{1 - (-1)}{||\vec{\mathbf{w}}||} = \frac{2}{||\vec{\mathbf{w}}||}$$

Maximal margin corresponds to minimal length of the weight vector



Best Separating Hyperplane

Minimize

 $\vec{w}^T \vec{w}$

Constraints

$$t_i \vec{\mathbf{w}}^T \vec{\mathbf{x}}_i \ge 1 \qquad \forall i$$



Linear separation

Structural Risk Minimization

Support Vector Machines

Kernels

Non-separable Classes



Observation

Almost everything becomes linearly separable when represented in high-dimensional spaces

"Ordinary" low-dimensional data can be "scattered" into a high-dimensional space.

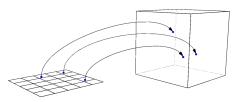
Two problems emerge

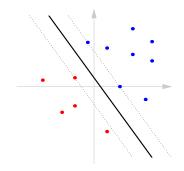
- 1. Many free parameters \rightarrow bad generalization
- 2. Extensive computations



Support Vector Machines

- 1. Transform the input to a suitable high-dimensional space
- Choose the unique separating hyperplane that has maximal margins
- 3. Classify new data using this hyperplane







Support Vector Machines

- Advantages
 - Very good generalization
 - Works well even with few training samples
 - Fast classification
- Disadvantages
 - Non-local weight calculation
 - Hard to implement efficiently



Linear separation

Structural Risk Minimization

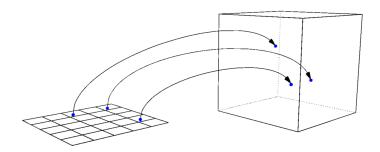
Support Vector Machines

Kernels

Non-separable Classes



Kernels: Only *pretend* that we transform the input data into a high-dimensional feature space!





Idea behind Kernels

Utilize the advantages of a high-dimensional space without actually representing anything high-dimensional

- Condition: The only operation done in the high-dimensional space is to compute scalar products between pairs of items
- ► Trick: The high-dimensional scalar product is computed using the original (low-dimensional) representation

Example

Points in 2D

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

Transformation to 4D

$$\phi(\vec{x}) = \begin{bmatrix} x_1^3 \\ \sqrt{3}x_1^2x_2 \\ \sqrt{3}x_1x_2^2 \\ x_3^3 \end{bmatrix}$$

$$\phi(\vec{x})^{T} \cdot \phi(\vec{y}) = x_{1}^{3} y_{1}^{3} + 3x_{1}^{2} y_{1}^{2} x_{2} y_{2} + 3x_{1} y_{1} x_{2}^{2} y_{2}^{2} + x_{2}^{3} y_{2}^{3}$$

$$= (x_{1} y_{1} + x_{2} y_{2})^{3}$$

$$= (\vec{x}^{T} \cdot \vec{y})^{3}$$

$$= \mathcal{K}(\vec{x}, \vec{y})$$



Common Kernels

Polynomials

$$\mathcal{K}(\vec{\mathbf{x}}, \vec{\mathbf{y}}) = (\vec{\mathbf{x}}^T \vec{\mathbf{y}} + 1)^{\mathbf{p}}$$

Radial Bases

$$\mathcal{K}(\vec{\mathbf{x}}, \vec{\mathbf{y}}) = \mathbf{e}^{-\frac{1}{2\rho^2}||\vec{\mathbf{x}} - \vec{\mathbf{y}}||^2}$$

Structural Risk Minimization

Minimize

$$\vec{w}^T \vec{w}$$

Constraints

$$t_i \vec{\mathbf{w}}^T \vec{\mathbf{x}}_i \geq 1 \qquad \forall i$$

▶ Non-linear transformation ϕ of input \vec{x}

New formulation

Minimize

$$\frac{1}{2}\vec{\mathbf{w}}^T\vec{\mathbf{w}}$$

Constraints

$$t_i \vec{\mathbf{w}}^T \phi(\vec{\mathbf{x}}_i) \ge 1 \quad \forall i$$



Structural Risk Minimization

Minimize

$$\frac{1}{2}\vec{\mathbf{w}}^T\vec{\mathbf{w}}$$

Constraints

$$t_i \vec{\mathbf{w}}^T \phi(\vec{\mathbf{x}}_i) \ge 1 \qquad \forall i$$

Lagranges Multiplier Method

$$L = \frac{1}{2} \vec{\mathbf{w}}^T \vec{\mathbf{w}} - \sum_{i} \alpha_i \left[t_i \vec{\mathbf{w}}^T \phi(\vec{\mathbf{x}}_i) - 1 \right]$$

Minimize w.r.t. \vec{w} , maximize w.r.t. $\alpha_i \ge 0$

$$\frac{\partial L}{\partial \vec{\mathbf{w}}} = 0$$

$$L = \frac{1}{2} \vec{w}^T \vec{w} - \sum_i \alpha_i \left[t_i \vec{w}^T \phi(\vec{x}_i) - 1 \right]$$
$$\frac{\partial L}{\partial \vec{w}} = 0 \implies \vec{w} - \sum_i \alpha_i t_i \phi(\vec{x}_i) = 0$$
$$\vec{w} = \sum_i \alpha_i t_i \phi(\vec{x}_i)$$



Use

$$\vec{\mathbf{w}} = \sum_{i} \alpha_{i} t_{i} \phi(\vec{\mathbf{x}}_{i})$$

to eliminate \vec{w}

$$L = \frac{1}{2} \vec{\mathbf{w}}^T \vec{\mathbf{w}} - \sum_{i} \alpha_i \left[t_i \vec{\mathbf{w}}^T \phi(\vec{\mathbf{x}}_i) - 1 \right]$$

$$L = \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j t_i t_j \phi(\vec{x}_i)^T \phi(\vec{x}_j) - \sum_{i,j} \alpha_i \alpha_j t_i t_j \phi(\vec{x}_i)^T \phi(\vec{x}_j) + \sum_i \alpha_i$$
$$L = \sum_i \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j t_i t_j \phi(\vec{x}_i)^T \phi(\vec{x}_j)$$



The Dual Problem

Maximize

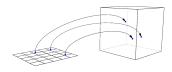
$$\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} t_{i} t_{j} \phi(\vec{\mathbf{x}}_{i})^{\mathsf{T}} \phi(\vec{\mathbf{x}}_{j})$$

Under the constraints

$$\alpha_i \geq 0 \quad \forall i$$

- \vec{w} has disappeared
- $\phi(\vec{x})$ only appear in scalar product pairs





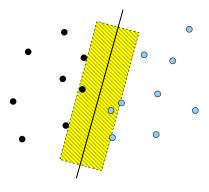
- 1. Choose a suitable kernel function
- 2. Compute α_i (solve the maximization problem)
- 3. \vec{x}_i corresponding to $\alpha_i \neq 0$ are called support vectors
- 4. Classify new data points via

$$\sum_{i} \alpha_{i} t_{i} \mathcal{K}(\vec{\mathbf{x}}, \vec{\mathbf{x}}_{i}) > 0$$



None-Separable Training Samples

Allow for Slack





Re-formulation of the minimization problem

Minimize

$$\frac{1}{2}\vec{\mathbf{w}}^T\vec{\mathbf{w}} + C\sum_i \xi_i$$

Constraints

$$t_i \vec{\mathbf{w}}^T \phi(\vec{\mathbf{x}}_i) \geq 1 - \xi_i$$

 ξ_i are called *slack variables*



Dual Formulation with Slack

Maximize

$$\sum_{i} \alpha_{i} - \frac{1}{2} \sum_{i,j} \alpha_{i} \alpha_{j} t_{i} t_{j} \phi(\vec{\mathbf{x}}_{i})^{\mathsf{T}} \phi(\vec{\mathbf{x}}_{j})$$

With constraints

$$0 \le \alpha_i \le C \quad \forall i$$

Otherwise, everything remains as before