

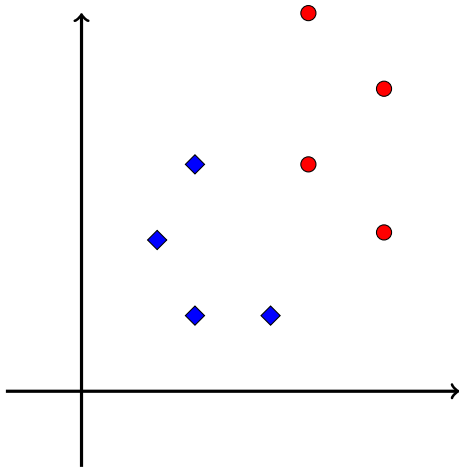
# Support Vector Machines (SVM) and Sequential Minimal Optimization (SMO)

Felix Brockherde

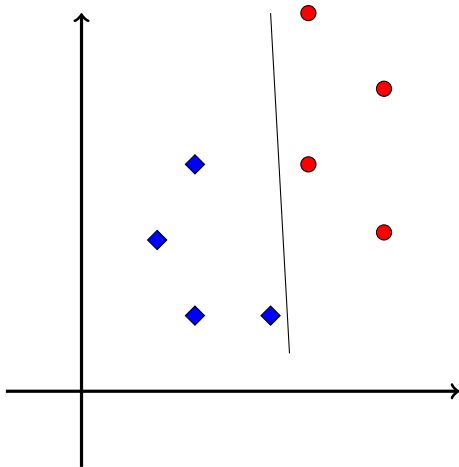
Technische Universität Berlin and Max Planck Institute of Microstructure Physics

25 Jun 2014

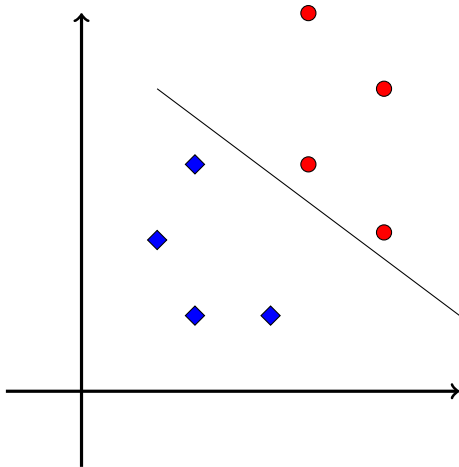
# Support Vector Machines (SVM)



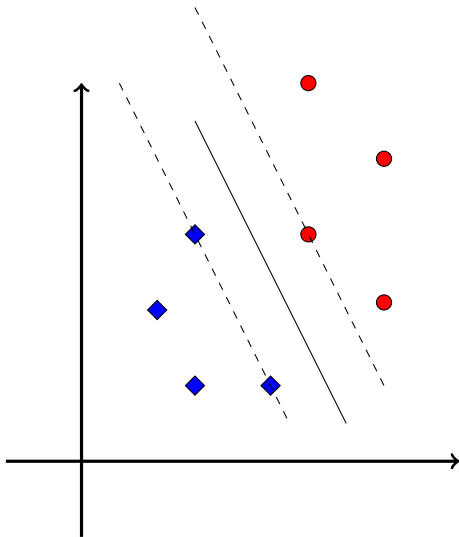
# Support Vector Machines (SVM)



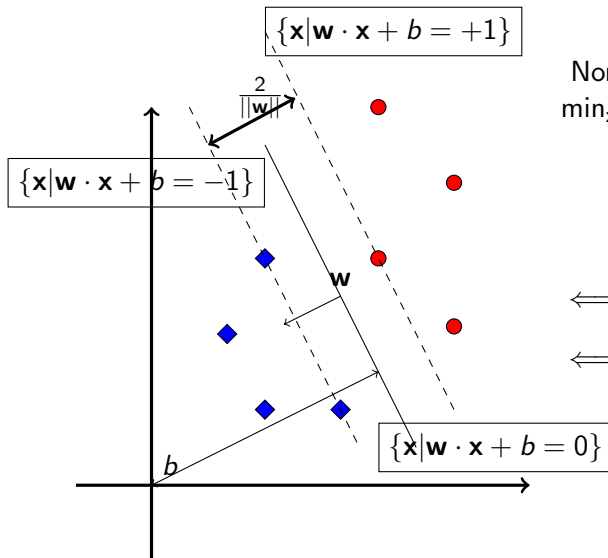
# Support Vector Machines (SVM)



# Support Vector Machines (SVM)



# Support Vector Machines (SVM)



Normalize  $\mathbf{w}$  so that  
 $\min_{x_i} \mathbf{w} \cdot \mathbf{x}_i + b = 1.$

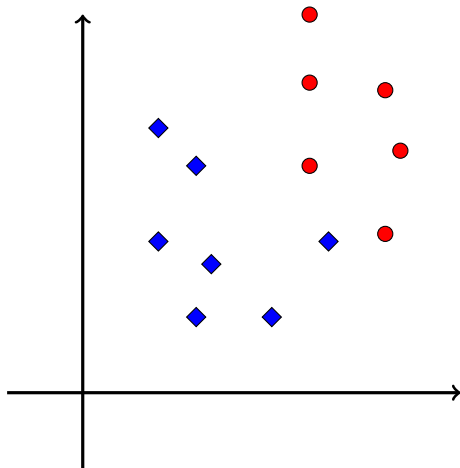
$$\mathbf{w} \cdot \mathbf{x}_1 + b = +1$$

$$\mathbf{w} \cdot \mathbf{x}_2 + b = -1$$

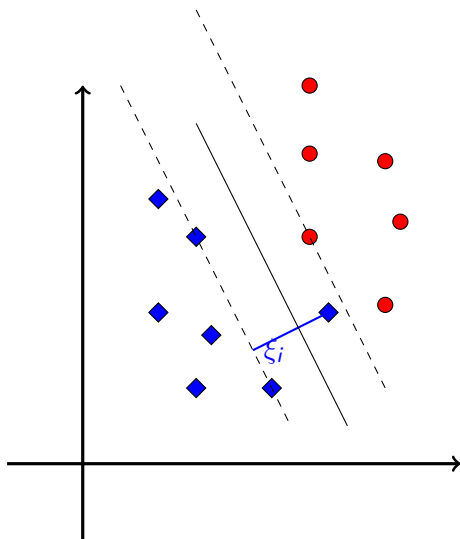
$$\iff \mathbf{w} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = 2$$

$$\iff \frac{\mathbf{w}}{\|\mathbf{w}\|} \cdot (\mathbf{x}_1 - \mathbf{x}_2) = \frac{2}{\|\mathbf{w}\|}$$

# Slack variables



# Slack variables



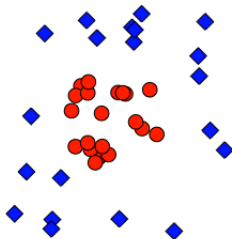
Introduce slack variables  $\xi_i$ :

$$\min_{\mathbf{w}, b, \xi_i} \quad \|\mathbf{w}\|^2 + C \sum_{i=1}^N \xi_i$$

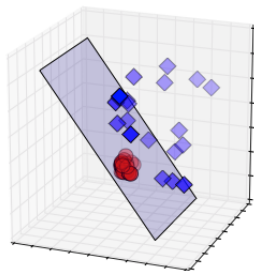
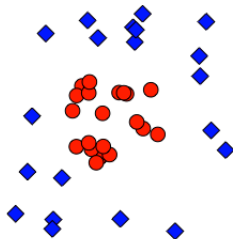
$$\text{subject to} \quad y_i(\mathbf{w} \cdot \mathbf{x}_i + b) \geq 1 - \xi_i$$
$$\xi_i \geq 0$$



# Non-linear hyperplanes



# Non-linear hyperplanes



Map into a higher dimensional feature space:

$$\Phi : \mathbb{R}^2 \rightarrow \mathbb{R}^3$$

$$(x_1, x_2) \mapsto (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

**Primal**

$$\begin{aligned} \min_{\mathbf{w}, b, \xi_i} \quad & ||\mathbf{w}'||^2 + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i(\mathbf{w} \cdot \Phi(\mathbf{x}_i) + b) \geq 1 - \xi_i \text{ and } \xi_i \geq 0 \text{ for } i = 1 \dots N \end{aligned}$$

**Dual**

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \text{ and } C \geq \alpha_i \geq 0 \text{ for } i = 1 \dots N \end{aligned}$$

Data points  $x_i$  only appear in scalar products  $(\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j))$ .

# The Kernel Trick

Replace scalar products with kernel function (?):

$$k(x, y) = \Phi(x) \cdot \Phi(y)$$

- ▶ Compute kernel matrix  $K_{ij} = k(x_i, x_j)$ , **i.e. never use  $\Phi$  directly**
- ▶ Underlying mapping  $\Phi$  can be unknown
- ▶ Kernels can be adopted to specific task, e.g. using prior knowledge (kernels for graphs, trees, strings, ...)

## Common kernels

**Gaussian Kernel:**  $k(x, y) = \exp\left(-\frac{\|x-y\|^2}{2\sigma^2}\right)$

Linear Kernel:  $k(x, y) = x \cdot y$

Polynomial Kernel:  $k(x, y) = (x \cdot y + c)^d$

# The Support Vectors in SVM

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j (\Phi(\mathbf{x}_i) \cdot \Phi(\mathbf{x}_j)) \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad C \geq \alpha_i \geq 0 \quad \text{for } i = 1 \dots N \end{aligned}$$

## KKT conditions

$$y_i [\mathbf{w} \Phi(\mathbf{x}_i) + b] > 1 \implies a_i = 0 \longrightarrow x_i \text{ irrelevant}$$

$$y_i [\mathbf{w} \Phi(\mathbf{x}_i) + b] = 1 \implies \text{on/in margin} \longrightarrow x_i \text{ Support Vector}$$

Old model  $f(x) = w \cdot \Phi(x_i) + b$  becomes via  $w = \sum_{i=1}^N \alpha_i y_i \Phi(x_i)$ :

$$f(x) = \sum_{i=1}^N \alpha_i y_i k(x_i, x) + b \longrightarrow f(x) = \sum_{x_i \in \text{SV}} \alpha_i y_i k(x_i, x) + b$$

# Quadratic Programming (QP)

Reminder: The SVM optimization problem

$$\begin{aligned} \max_{\alpha} \quad & \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i,j=1}^N \alpha_i \alpha_j y_i y_j k(x_i, x_j) \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \quad \text{and} \quad C \geq \alpha_i \geq 0 \quad \text{for } i = 1 \dots N \end{aligned}$$

## Quadratic Programming

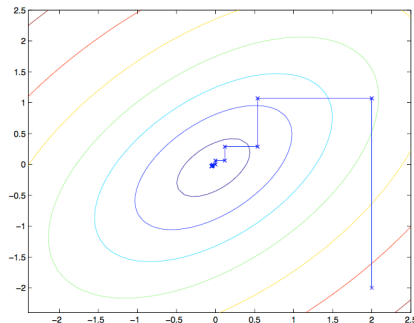
$$\begin{aligned} \min_x \quad & \frac{1}{2} x^T P x + q^T x \\ \text{s.t.} \quad & Gx \preceq h \\ & Ax = b \end{aligned}$$

⇒ We can solve the SVM problem with a QP solver.

# Coordinate Descent

Imagine a multivariate function  $W(x_1, x_2, \dots, x_N)$  and the problem

$$\operatorname{argmax}_{\mathbf{x}} W(x_1, x_2, \dots, x_N)$$



## Coordinate Descent Algorithm

```
while not converged do
  for  $i = 1 \dots N$  do
     $x_i \leftarrow \operatorname{argmax}_{x_i} W(x_1, \dots, x_i, \dots, x_N)$ 
  end for
end while
```

# Sequential Minimal Optimization (SMO) Idea

We have to solve the SVM optimization problem subject to the condition

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

I.e. we can *not* optimize one single variable individually.

Solution: Optimize two variables  $\alpha_i$  and  $\alpha_j$  with  $y_i \neq y_j$  while keeping the other  $\alpha$  fixed.

From KKT conditions, we get (optimizing  $\alpha_1$  and  $\alpha_2$  WLOG)

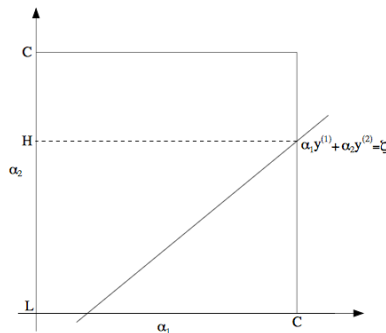
$$\alpha_1 y_1 + \alpha_2 y_2 = - \sum_{i=3}^N \alpha_i y_i$$

where the right handside is fixed  $\zeta = - \sum_{i=3}^N \alpha_i y_i$ .



# Box constraints

- ▶ From the second constraint  $0 \leq \alpha_i \leq C$ , we know that  $\alpha_1$  and  $\alpha_2$  have to lie inside the pictured box:
- ▶ We can see, that there exists a upper and lower bound for  $\alpha_2$ :  
 $L \leq \alpha_2 \leq H$ .



- ▶ If we write  $\alpha_1 = (\zeta - \alpha_2 y_2) y_1$ , we have to solve a quadratic function in one variable ( $\alpha_2$ ).
- ▶ The new  $\alpha_2$  then has to be clipped to the box constraints  $L$  and  $H$  (!)
- ▶ Find  $\alpha_1$  via  $\alpha_1 = (\zeta - \alpha_2 y_2) y_1$ . Then pick two new  $\alpha$  and repeat.