



4.1. Support Vector Machines

Multimedia Project SS 2016

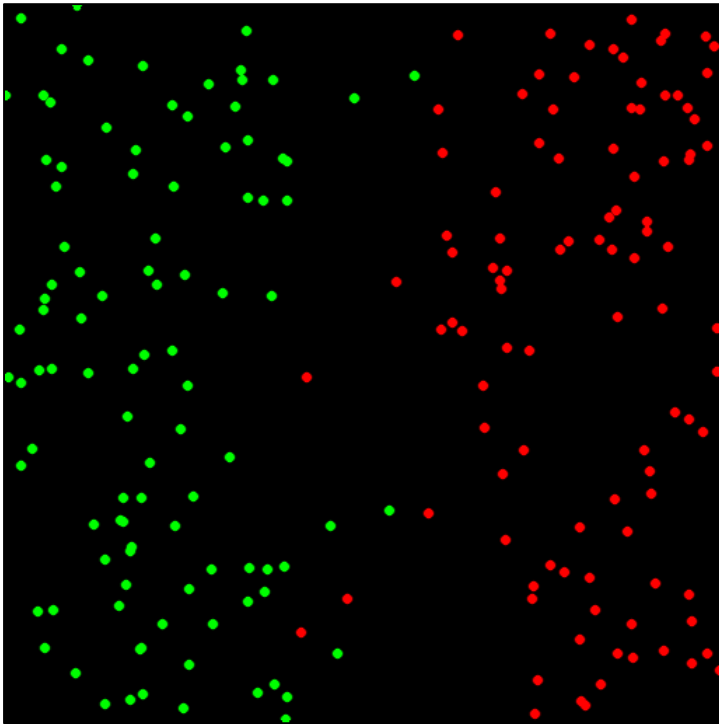


Multimedia Computing and Computer Vision Lab

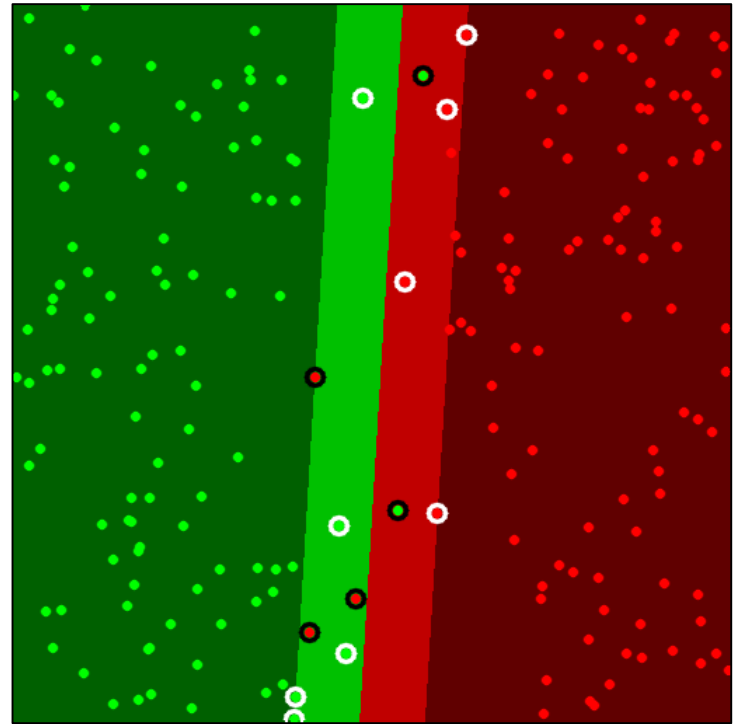
Anton Winschel

Support Vector Machines (SVM)

- Purpose: Binary classification with linear kernel
- Sparse solutions: Predictions only need the kernel function evaluated at a subset of training data points



2 class training data



maximum margin classification

Part I: Hyperplane Classifiers

Maximum Margin Classifier

- Two class classification problem using linear decision model

$$y(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$$

$$y : \mathbb{R}^M \rightarrow \mathbb{R}$$

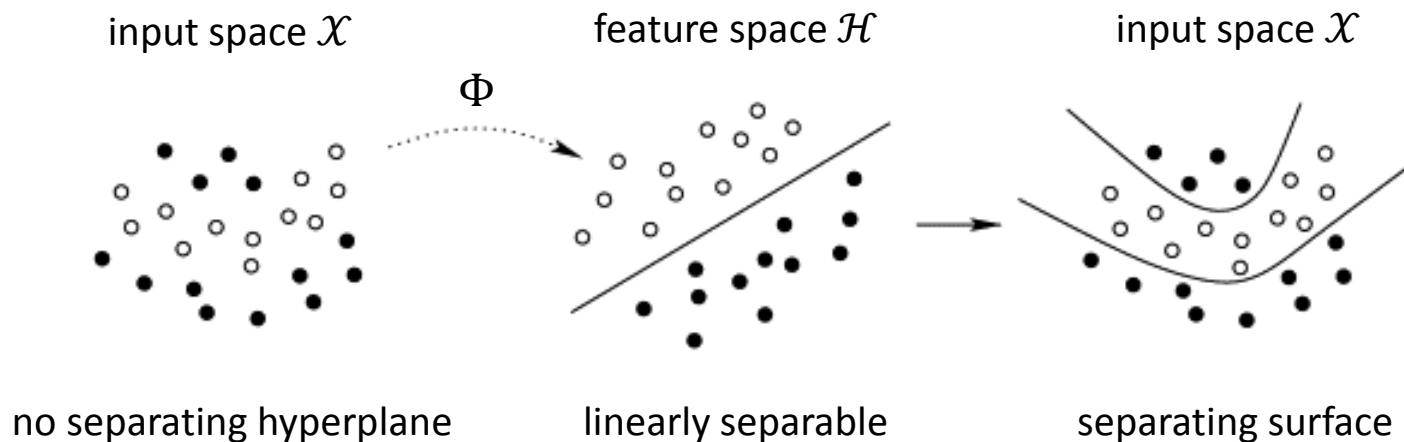
- Training data set
 - Examples $(\mathbf{x}_1, t_1), \dots, (\mathbf{x}_N, t_N) \in \mathbb{R}^M \times \{\pm 1\}$
 - Data points are represented by *pattern* \mathbf{x}_i (\rightarrow attributes), and target value (or *label*) t_i (\rightarrow class)
 - Finding hypothesis that *optimally* splits data points according to their classes
 - If there are multiple solutions, try to find the one that gives the smallest generalization error
 - Convex optimization problem

Maximum Margin Classifier

- Classification
 - New data point x is classified according to the sign of $y(x)$
 - Only requires a small subset of training points – all other points can be discarded

Non-Linear Kernels

- In case of data that is not linearly separable, non-linear kernel functions can be used to map data points into higher-dimensional feature space
- Potentially better classification results but: (i) higher computational cost at test time, (ii) longer training time and (iii) potential overfitting of model to data
- Linear separation in feature space \mathcal{H} gives non-linear decision surface in input space \mathcal{X} :



Kernel Trick

- Kernel trick: Replace dot product by non-linear kernel (e.g. Gaussian)
 - Requires algorithm that depends on data only through dot products
 - By replacing the dot product with a kernel function, data points from the input domain are mapped into (possibly very high-dimensional) feature space
- Generates non-linear decision boundary in the input space
- Prediction: Linear combination of kernel functions

$$y(\mathbf{x}) = \sum_{i=1}^N \alpha_i t_i k(\mathbf{x}, \mathbf{x}_i) + b, \text{ with } k(\mathbf{x}, \mathbf{x}') = \Phi(\mathbf{x})^T \Phi(\mathbf{x}')$$

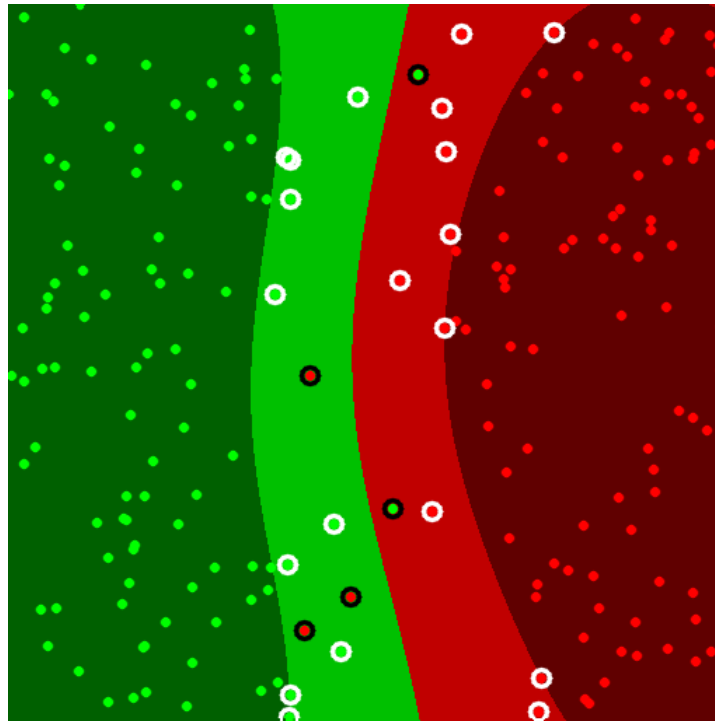
- Using identity mapping $\Phi(\mathbf{x}) := \mathbf{x}$ the kernel is a simple dot product:

$$k(\mathbf{x}, \mathbf{x}') = \langle \mathbf{x}, \mathbf{x}' \rangle$$

- For linear models, computing $y(\mathbf{x})$ is very efficient

Example: Non-Linear Decision Surfaces

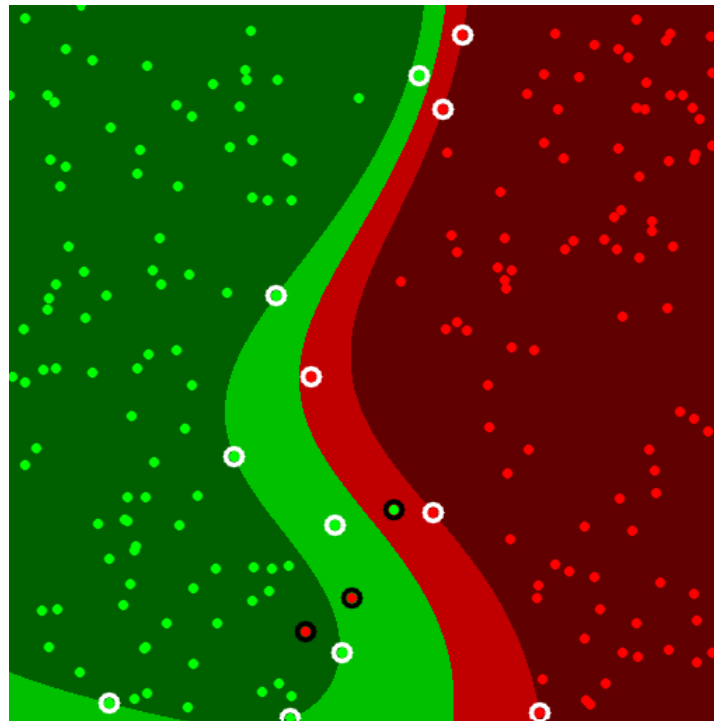
- RBF kernel (non-linear)



$$C = 1$$

Example: Non-Linear Decision Surfaces

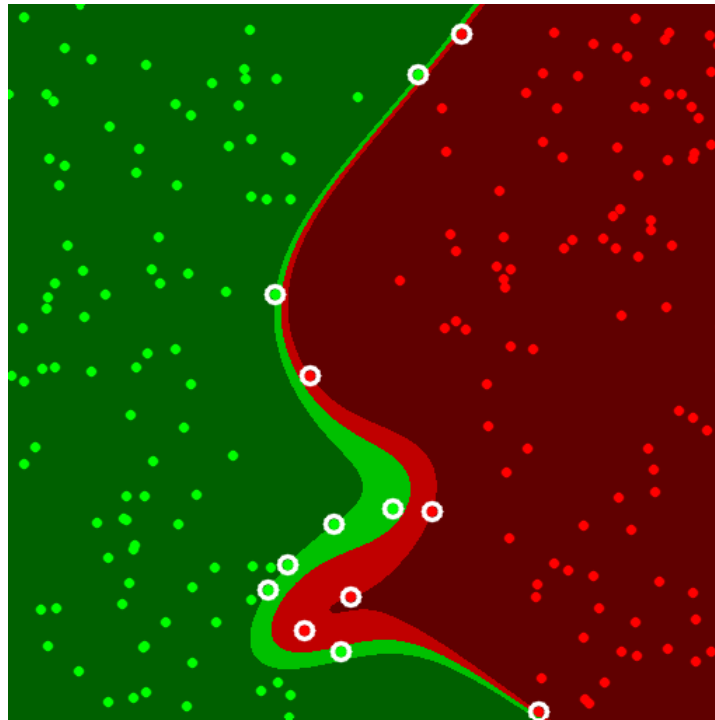
- RBF kernel (non-linear)



$C = 100$

Example: Non-Linear Decision Surfaces

- RBF kernel (non-linear)



$C = 10000$

Part II: Hard-Margin Formulation

Hyperplane Classifiers

- Hyperplane learning algorithm

- A *hyperplane* in dot product space \mathcal{H} is defined by

$$\{\mathbf{x} \in \mathcal{H} | \langle \mathbf{w}, \mathbf{x} \rangle + b = 0\}, \text{ where } \mathbf{w} \in \mathcal{H}, b \in \mathbb{R}$$

- A *linear decision function* corresponding to \mathbf{w} is given by

$$y(\mathbf{x}) = \langle \mathbf{w}, \mathbf{x} \rangle + b$$

- Assumption: N linearly separable training examples $\mathbf{x}_1, \dots, \mathbf{x}_N$, with target values t_1, \dots, t_N
 - There exists at least one \mathbf{w} and b s.t. $y(\mathbf{x}_i) > 0$ for all points having $t_i = +1$, and $y(\mathbf{x}_i) < 0$ for those having $t_i = -1$
 - If hyperplane separates data perfectly, it holds $t_i y(\mathbf{x}_i) > 0$ for all training examples
 - Classification according to $\text{sign}(y(\mathbf{x}))$