



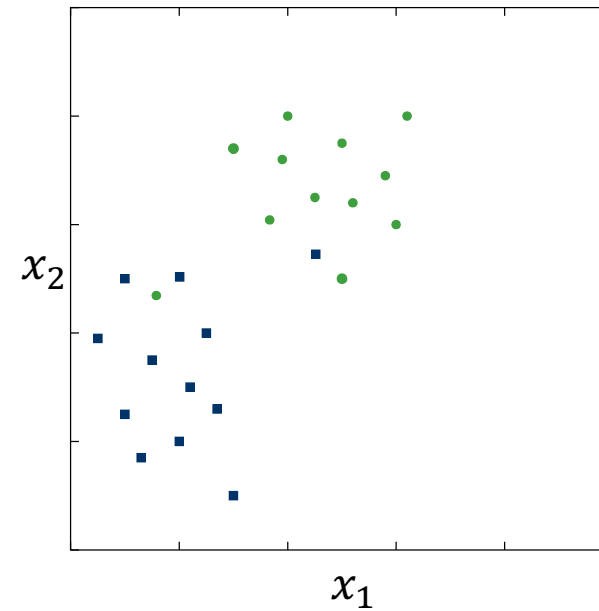
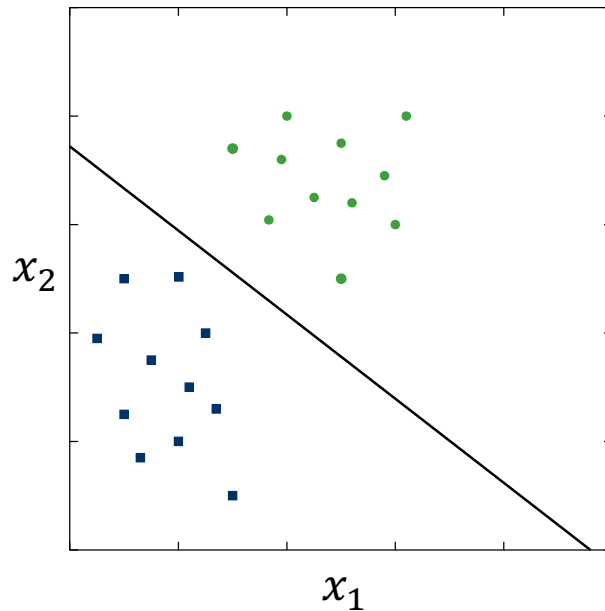
Machine Learning

Support Vector Machines

Prof. Dr. Jochen Schmidt

Motivation

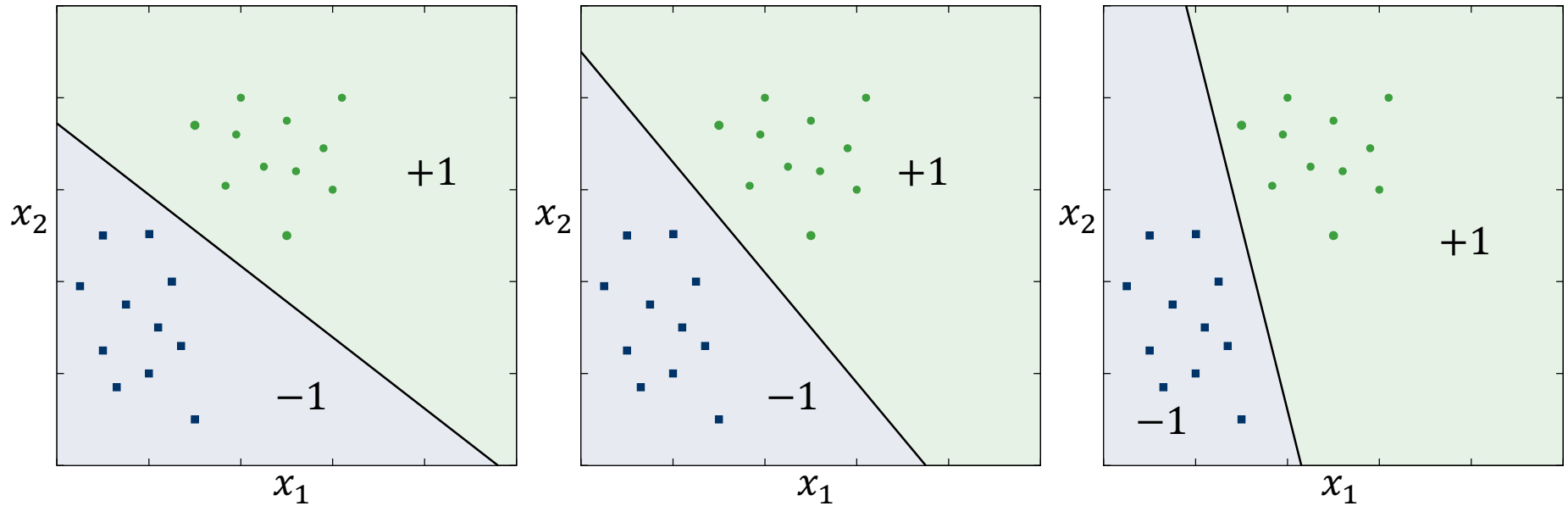
- assume two linearly separable classes
- compute linear decision boundary that
 - allows for separation of training data
 - generalizes well





Motivation

Many, many solutions...



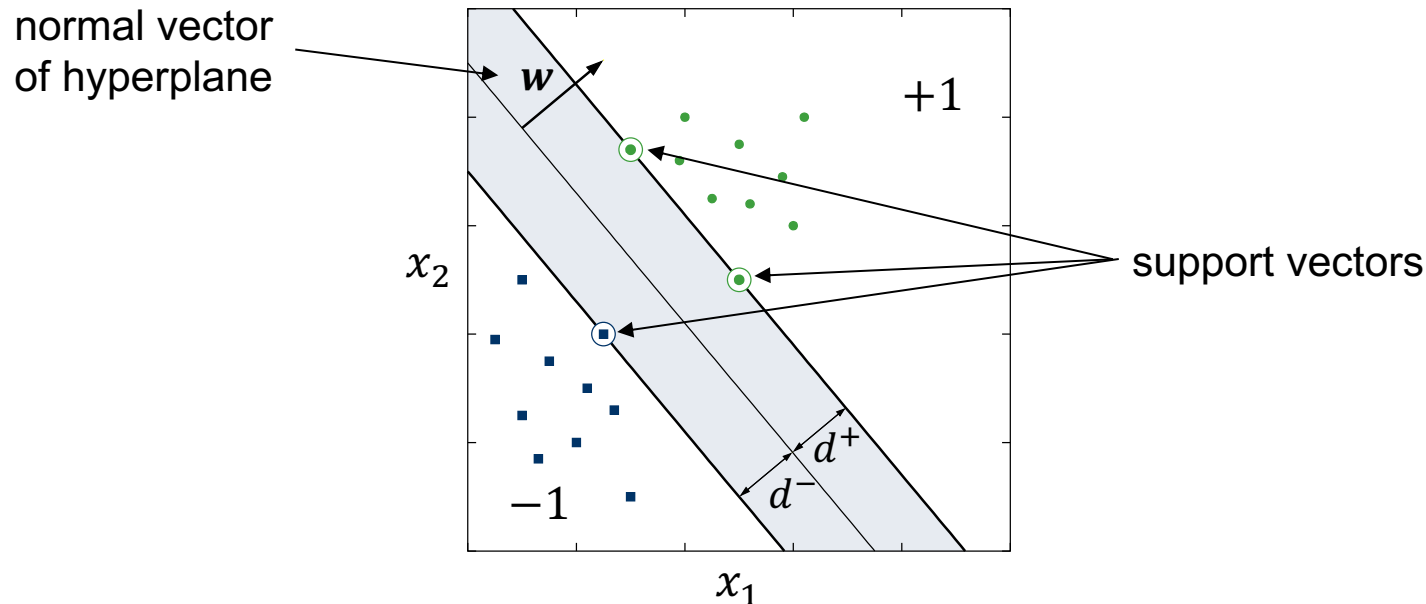
Optimal Separating Hyperplane

Vapnik 1996: Optimal separating hyperplane that

- separates two classes and
- maximizes the distance to the closest point from either class.

This results in

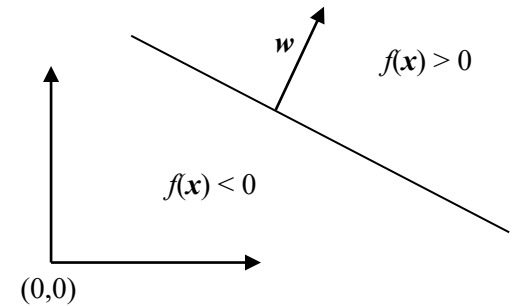
- unique solution for hyperplanes, and
- (in most cases) better generalization.



Functional margin of separating hyperplane: $d^+ + d^-$

Optimal Separating Hyperplane

- Plane equation: $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + w_0$
 - normal vector: \mathbf{w}
 - point on plane: $f(\mathbf{x}) = 0$
 - point above plane: $f(\mathbf{x}) > 0$
 - point below plane: $f(\mathbf{x}) < 0$
 - "above" = in direction of plane normal



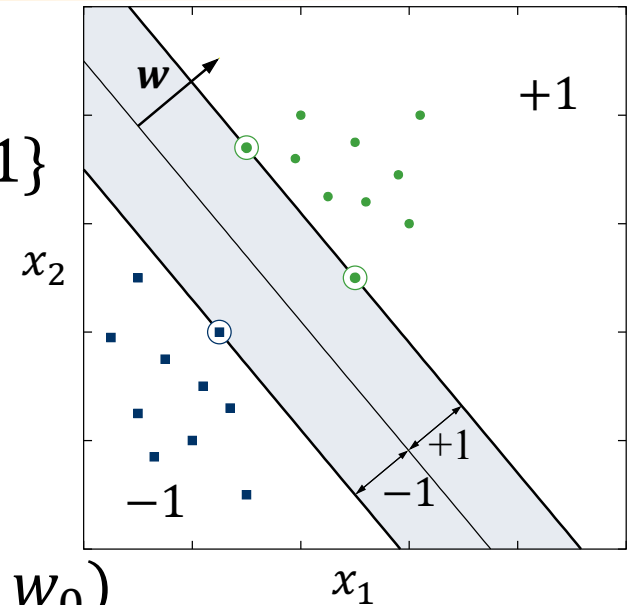
- Signed distance d of a point to hyperplane
 - normalize \mathbf{w} , such that $|\mathbf{w}| = 1$:

$$d = f'(\mathbf{x}) = \frac{1}{|\mathbf{w}|} f(\mathbf{x}) = \frac{1}{|\mathbf{w}|} \mathbf{w}^T \mathbf{x} + \frac{1}{|\mathbf{w}|} w_0$$

- distance of plane from origin: $-\frac{1}{|\mathbf{w}|} w_0$

SVM – Classification

- data point: x_i
- class of data point x_i is $y_i \in \{-1, +1\}$
- Classifier: $g(x_i) = \text{sgn}(\mathbf{w}^T \mathbf{x}_i + w_0)$
- Functional margin of x_i : $y_i (\mathbf{w}^T \mathbf{x}_i + w_0)$
 - can be increased/decreased by scaling plane equation
 - \rightarrow scale such that support vectors have distance $-1/+1$
 - then $g(x_i)$ is equivalent to: $y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$
- Functional margin for data set:
2 x minimum functional margin of all points: $\frac{2}{|\mathbf{w}|}$



SVM – Training


- Training =
 - find hyperplane maximizing the margin $\frac{2}{|\mathbf{w}|}$
 - subject to constraint $y_i (\mathbf{w}^T \mathbf{x}_i + w_0) \geq 1$ for all data points
- This is equivalent to
 - minimizing $\frac{1}{2} |\mathbf{w}|^2 = \frac{1}{2} \mathbf{w}^T \mathbf{w}$
 - subject to $y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 \geq 0$
- Remarks
 - this is a convex optimization problem
 - local optimum is always a global one – solution is unique
 - there exist efficient algorithms for convex optimization
 - standard libraries can be used

Optimization – Lagrangian

Solving the constrained convex optimization problem requires the Lagrangian, i.e. minimize

$$L(\mathbf{w}, w_0, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

Lagrange multipliers



Minimization:

- Compute first partial derivatives
- Set derivatives to zero

Here: The problem is reformulated as the so-called Lagrangian Dual, which is then maximized

Optimization – Lagrangian

$$L(\mathbf{w}, w_0, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1)$$

Partial derivative (1):

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\lambda})}{\partial \mathbf{w}} = \mathbf{w} - \sum_i \lambda_i y_i \mathbf{x}_i = 0$$

$$\Rightarrow \mathbf{w} = \sum_i \lambda_i y_i \mathbf{x}_i$$

Partial derivative (2):

$$\frac{\partial L(\mathbf{w}, w_0, \boldsymbol{\lambda})}{\partial w_0} = - \sum_i \lambda_i y_i = 0$$

the partial derivative for $\boldsymbol{\lambda}$ is not required for the following reformalization

Langrange Dual

$$\begin{aligned}
 L_D &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_i \lambda_i (y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1) \\
 &= \frac{1}{2} \mathbf{w}^T \mathbf{w} - \underbrace{(\sum_i \lambda_i y_i \mathbf{x}_i)^T \mathbf{w}}_{= \mathbf{w}^T} - \underbrace{\sum_i \lambda_i y_i w_0}_{= 0} + \sum_i \lambda_i \\
 &\quad \text{partial derivative (1)} \qquad \text{partial derivative (2)}
 \end{aligned}$$

$$= -\frac{1}{2} \mathbf{w}^T \mathbf{w} + \sum_i \lambda_i$$

$$= -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$$

This term is now maximized subject to $\lambda_i > 0$

SVM – Optimization Solution

Training: Maximize $L_D = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i$

Resulting hyperplane:

$$\mathbf{w} = \sum_i \lambda_i y_i \mathbf{x}_i \quad w_0 = y_k - \mathbf{w}^T \mathbf{x}_k \text{ for any } \mathbf{x}_k \text{ with non-zero } \lambda_i$$

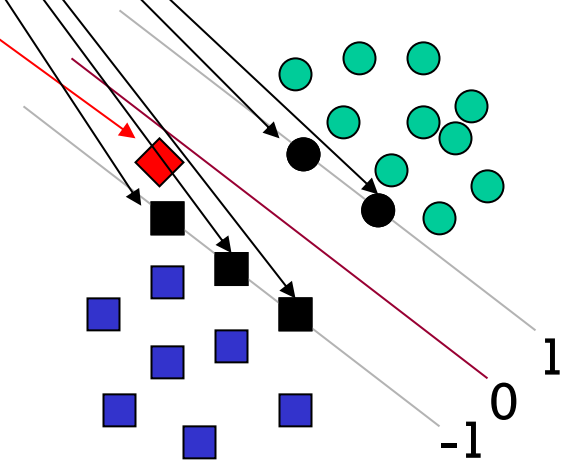
$$\text{Classification: } g(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + w_0) = \text{sgn} \left(\sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + w_0 \right)$$

- the \mathbf{x}_i with non-zero λ_i are the support vectors
- the feature vectors \mathbf{x}_i only appear in inner products
 - in training as well as classification phase
- \mathbf{w} is not required explicitly for classification
 - we need only the support vectors and their Lagrange multipliers

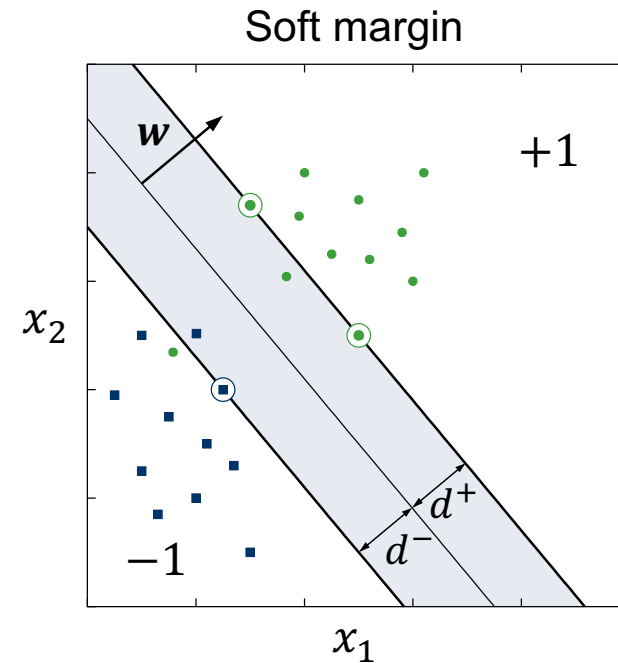
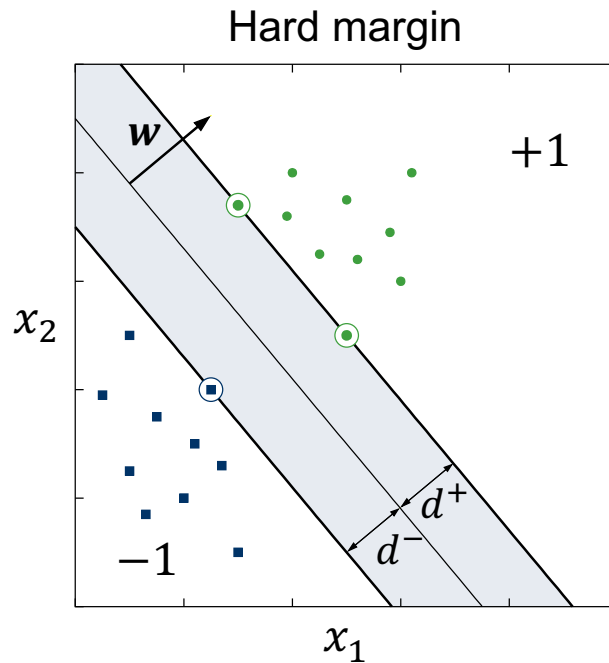
SVM – Classification with threshold

Classification: $g(\mathbf{x}) = \sum_i \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + w_0$

- Classification without threshold
 - decide for class based on $g(\mathbf{x}) < 0$ or $g(\mathbf{x}) > 0$
- Classification with confidence threshold t
 - $g(\mathbf{x}) < -t$: class -1
 - $g(\mathbf{x}) > t$: class $+1$
 - $-t < g(\mathbf{x}) < t$: reject



Hard and Soft Margin Problem



data are not linearly separable in this case

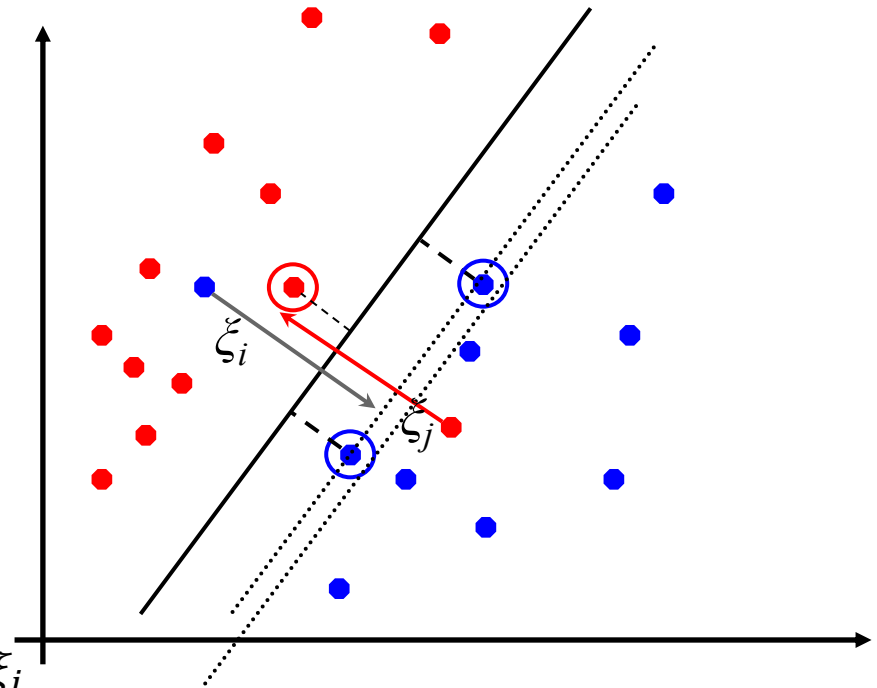
Soft Margin – Slack Variables

for not linearly separable data: allow some errors

- allow missclassification of difficult or noisy samples
- move data slightly, to where they belong
- pay a penalty

Training

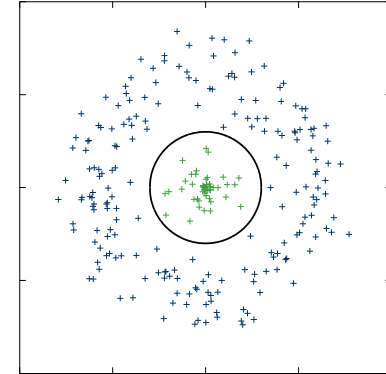
- minimize $\frac{1}{2} |\mathbf{w}|^2 + c \sum \xi_i$
- subject to
$$y_i (\mathbf{w}^T \mathbf{x}_i + w_0) - 1 + \xi_i \geq 0$$
$$\xi_i \geq 0$$
- reformulate as dual problem
 - neither the slack variables ξ_i
 - nor their Lagrange multipliers c appear there



Kernels / Non-linear Boundaries

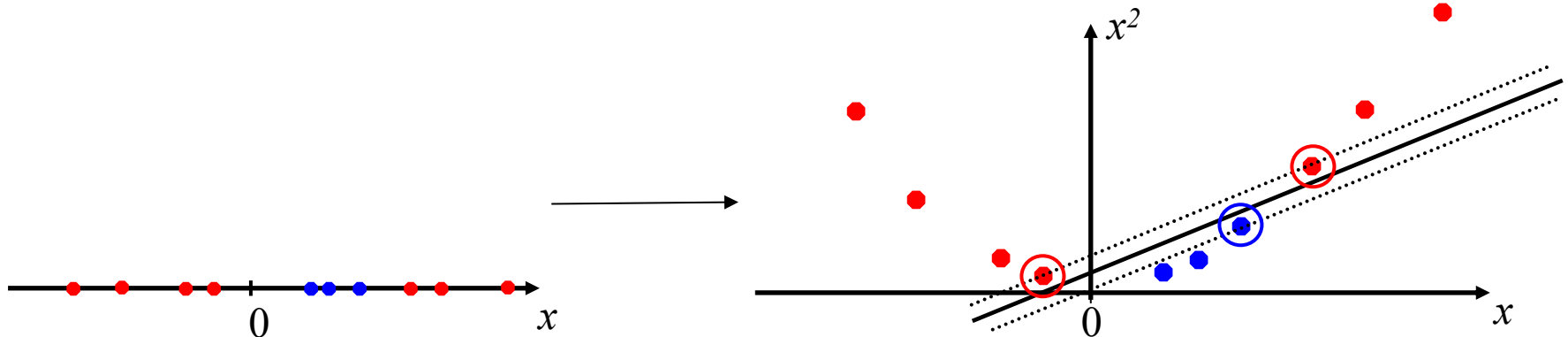
Limitations of linear decision boundaries

- too simple for most practical purposes
- non-linearly separable data cannot be classified
- noisy data cause problems



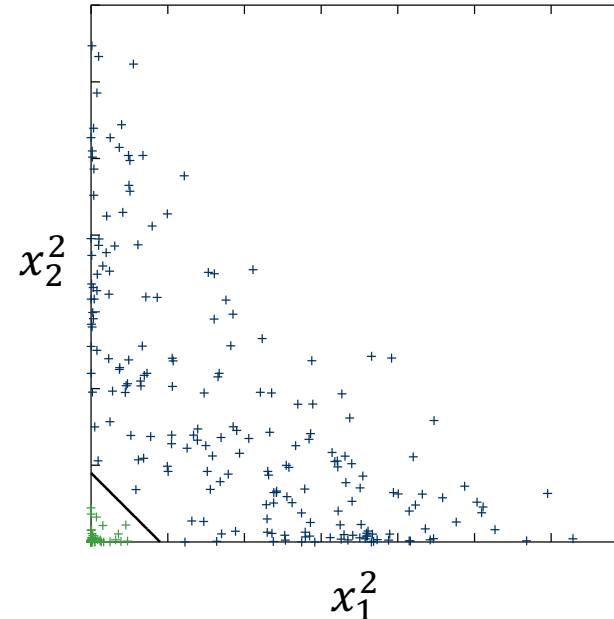
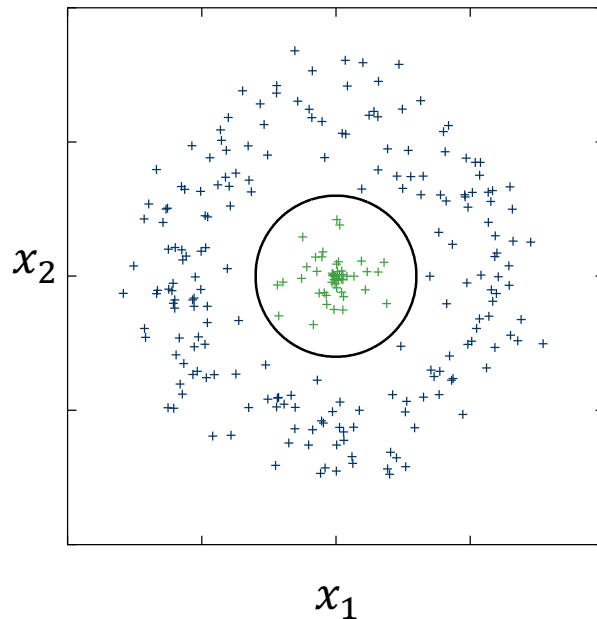
Possible solution

- Map data to higher dimensional feature space using non-linear feature transform,
- then use a linear classifier



Feature Transforms

Select a feature transform $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^D$ such that the resulting features $\phi(x_i)$ are linearly separable.



Applied feature transform in example: $\phi(x_i) = (x_1^2, x_2^2)^T$

Feature Transforms – Example

Assume the decision boundary is defined by the quadratic function

$$f(\mathbf{x}) = a_0 + a_1x_1^2 + a_2x_2^2 + a_3x_1x_2 + a_4x_1 + a_5x_2$$

This is obviously non-linear.

Using the following mapping, we get features having a linear decision boundary:

$$\phi(\mathbf{x}) = \begin{pmatrix} 1 \\ x_1^2 \\ x_2^2 \\ x_1x_2 \\ x_1 \\ x_2 \end{pmatrix}$$

Feature Transforms – SVM

The feature transforms can be easily incorporated into SVMs:

Replace $\mathbf{x}_i^T \mathbf{x}$ by $\phi^T(\mathbf{x}_i)\phi(\mathbf{x}) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle$ ————— $\langle \cdot \rangle$ notation for inner product

Classification/Decision boundary:

$$g(\mathbf{x}) = \sum_i \lambda_i y_i \phi^T(\mathbf{x}_i) \phi(\mathbf{x}) + w_0 = \sum_i \lambda_i y_i \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle + w_0$$

Training: Lagrange dual problem

Maximize
$$L_D = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle + \sum_i \lambda_i$$

subject to $\lambda_i > 0$

Note:

- the actual transform $\phi(\mathbf{x})$ is never required stand-alone
- it only appears as inner products \rightarrow Kernel Trick

Kernel-Trick

- in SVM training/classification, data appear only in the form of inner products $\langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$
- a Kernel-function is a function computing this inner product directly:
$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$$
 - i.e., without first transforming the features using $\phi(\mathbf{x})$
 - it can be computed in the original low-dimensional space!
- checking whether $K(\mathbf{x}_i, \mathbf{x}_j) = \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle$ holds is often cumbersome
 - Mercer's theorem:
Any positive semi-definite symmetric function is a kernel

Common Kernel Functions

- linear:
$$K(\mathbf{x}_i, \mathbf{x}_j) = \langle \mathbf{x}_i, \mathbf{x}_j \rangle$$
- polynomial:
$$K(\mathbf{x}_i, \mathbf{x}_j) = (\langle \mathbf{x}_i, \mathbf{x}_j \rangle + 1)^d$$
- Laplacian radial basis function (RBF):
$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_1}{\sigma^2}}$$
- Gaussian radial basis function (RBF):
$$K(\mathbf{x}_i, \mathbf{x}_j) = e^{-\frac{\|\mathbf{x}_i - \mathbf{x}_j\|_2^2}{\sigma^2}}$$
- sigmoid:
$$K(\mathbf{x}_i, \mathbf{x}_j) = \tanh(\alpha \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \beta)$$

Notes

- SVMs are perceptrons
 - training is different
 - compare sigmoid kernel to perceptron computation
 - Multilayer Perceptrons
 - are SVMs maximizing the margin in hidden layer space
 - all hidden units are SVMs
- Multiclass-SVM
 - split into multiple binary classifications
 - one-vs-all
 - one binary SVM per class, separating this class from all others
 - winner-takes all strategy (winner = class with highest value)
 - one-vs-one
 - train binary SVMs for each pair of classes
 - each SVM votes: max-wins strategy
- SVMs in scikit-learn:
<https://scikit-learn.org/stable/modules/svm.html>

References

slides based on

- slides of the lecture *Pattern Recognition* taught at the FAU Erlangen-Nuremberg, courtesy of D. Hahn, J. Hornegger, S. Steidl and E. Nöth.
- Ray Mooney: Support Vector Machines. Slides, University of Texas at Austin.
- Ch. Manning, P. Nayak: Introduction to Information Retrieval, Lecture 14: Support vector machines and machine learning documents. Stanford University.