

Pattern Recognition (PR)

Prof. Dr.-Ing. Andreas Maier
Pattern Recognition Lab (CS 5), Friedrich-Alexander-Universität Erlangen-Nürnberg
Winter Term 2020/21



This is a printable version of the slides of the lecture

Pattern Recognition (PR)
Winter term 2020/21
Friedrich-Alexander University of Erlangen-Nuremberg.

These slides are released under Creative Commons License Attribution CC BY 4.0.

Please feel free to reuse any of the figures and slides, as long as you keep a reference to the source of these slides at <https://lme.tf.fau.de/teaching/> acknowledging the authors Niemann, Hornegger, Hahn, Steidl, Nöth, Seitz, Rodriguez, Das and Maier.

Erlangen, January 8, 2021
Prof. Dr.-Ing. Andreas Maier

Support Vector Machines II



Hard Margin Problem

The hard margin SVM optimization problem is formulated as:

$$\text{minimize} \quad \frac{1}{2} \|\boldsymbol{\alpha}\|_2^2$$

$$\text{subject to} \quad \forall i: \quad y_i \cdot (\boldsymbol{\alpha}^T \mathbf{x}_i + \alpha_0) - 1 \geq 0$$

Soft Margin Problem

The soft margin SVM optimization problem is formulated as:

$$\text{minimize} \quad \frac{1}{2} \|\alpha\|_2^2 + \mu \sum_i \xi_i$$

$$\text{subject to} \quad \forall i: \quad -(y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1 + \xi_i) \leq 0,$$

$$\forall i: \quad -\xi_i \leq 0$$

Lagrangian

The solution of the **constrained convex optimization problem** requires the Lagrangian:

$$\begin{aligned} L(\alpha, \alpha_0, \xi, \lambda, \mu) = & \frac{1}{2} \|\alpha\|_2^2 + \mu \sum_i \xi_i - \sum_i \mu_i \xi_i \\ & - \sum_i \lambda_i (y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1 + \xi_i) \end{aligned}$$

Lagrangian

The solution of the **constrained convex optimization problem** requires the Lagrangian:

| meta- parameter | Lagrangian multiplier |
|-----------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------|
| $L(\alpha, \alpha_0, \xi, \lambda, \mu) = \frac{1}{2} \ \alpha\ _2^2 + \mu \sum_i \xi_i - \sum_i \mu_i \xi_i$ $- \sum_i \lambda_i (y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1 + \xi_i)$ | |

Lagrangian

The solution of the **constrained convex optimization problem** requires the Lagrangian:

| meta- parameter | Lagrangian multiplier |
|---------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|--------------------------|
| $L(\alpha, \alpha_0, \xi, \lambda, \mu) = \frac{1}{2} \ \alpha\ _2^2 + c \sum_i \xi_i - \sum_i \mu_i \xi_i$ $- \sum_i \lambda_i (y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1 + \xi_i)$ | |

Lagrangian (cont.)

Partial derivatives I:

$$\frac{\partial L(\alpha, \alpha_0, \xi, \lambda, \mu)}{\partial \alpha} = \alpha - \sum_i \lambda_i y_i x_i \stackrel{!}{=} 0.$$

Thus we have:

$$\alpha = \sum_i \lambda_i y_i x_i.$$

Lagrangian (cont.)

Partial derivatives II:

$$\frac{\partial L(\alpha, \alpha_0, \xi, \lambda, \mu)}{\partial \alpha_0} = - \sum_i \lambda_i y_i \stackrel{!}{=} 0$$

Partial derivatives III:

$$\frac{\partial L(\alpha, \alpha_0, \xi, \lambda, \mu)}{\partial \xi_i} = c - \mu_i - \lambda_i \stackrel{!}{=} 0$$

Lagrange Dual

Let us consider the Lagrange function for the dual problem for the hard margin case:

$$\begin{aligned}
 L_D &= \frac{1}{2} \alpha^T \alpha - \sum_i \lambda_i (y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1) \\
 &= \frac{1}{2} \alpha^T \alpha - \underbrace{\left(\sum_i \lambda_i y_i \cdot \mathbf{x}_i \right)^T \alpha}_{\alpha^T} - \underbrace{\sum_i \lambda_i y_i}_{=0} \alpha_0 + \sum_i \lambda_i \\
 &= -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \cdot \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i
 \end{aligned}$$

The Lagrange Dual Problem

The Lagrange dual problem is given by the optimization problem:

$$\begin{aligned}
 &\text{maximize} && -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \cdot \mathbf{x}_i^T \mathbf{x}_j + \sum_i \lambda_i \\
 &\text{subject to} && \lambda \succeq 0 \\
 &&& \sum_i \lambda_i y_i = 0
 \end{aligned}$$

Benefits of the dual representation

- The model can be reformulated using kernels.
- SVMs can be applied efficiently to feature spaces whose dimensionality exceeds the number of samples.

Lagrange Dual Problem (cont.)

For convex optimization problems with differentiable objective and constraint functions, the duality gap is zero, if the KKT conditions are satisfied.

Especially the **complementary slackness** condition is interesting for us:

$$\forall i: \quad \lambda_i f_i(\mathbf{x}) = 0$$

Lagrange Dual Problem (cont.)

Complementary slackness for hard margin SVMs:

$$\forall i: \quad \lambda_i (y_i \cdot (\alpha^T \mathbf{x}_i + \alpha_0) - 1) = 0$$

Implications:

1. If $\lambda_i > 0$, then $y_i (\alpha^T \mathbf{x}_i + \alpha_0) - 1 = 0$, and thus:

$$y_i (\alpha^T \mathbf{x}_i + \alpha_0) = 1.$$

All \mathbf{x}_i with $\lambda_i > 0$ are elements at the boundary of the slab;
these samples are called **support vectors**.

2. We have seen that $\alpha = \sum_i \lambda_i y_i \mathbf{x}_i$, thus the norm vector of the decision boundary is a linear combination of support vectors.

Dual Representation

The decision function can also be rewritten using the duality:

$$f(\mathbf{x}) = \alpha^T \mathbf{x} + \alpha_0 = \sum_{i=1}^m \lambda_i y_i \mathbf{x}_i^T \mathbf{x} + \alpha_0$$

Conclusion:

Feature vectors only appear in inner products, both in the learning and the classification phase.

Feature Transforms

Linear decision boundaries in its current form have serious limitations:

- Non-linearly separable data cannot be classified.
- Noisy data cause problems.
- Formulation deals with vectorial data only.

Possible solution:

- Map data into richer feature space using non-linear feature transform, then use a linear classifier.

Feature Transforms (cont.)

We select a feature transform

$$\phi : \mathbb{R}^d \rightarrow \mathbb{R}^D, \quad D \geq d$$

such that the resulting features

$$\phi(\mathbf{x}_i), \quad i = 1, 2, \dots, m$$

are linearly separable.

Feature Transforms (cont.)

Example

Assume the decision boundary is given by the quadratic function

$$f(\mathbf{x}) = a_0 + a_1 x_1^2 + a_2 x_2^2 + a_3 x_1 x_2 + a_4 x_1 + a_5 x_2.$$

Obviously this is not a linear decision boundary.

By the following mapping, we get features that have a linear decision boundary:

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

Feature Transforms (cont.)

These feature transforms can be easily incorporated into SVMs:

- Decision boundary:

$$f(\mathbf{x}) = \sum_i \lambda_i y_i \cdot \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}) \rangle + \alpha_0$$

- The Lagrange dual problem is given by the **optimization problem**:

$$\text{maximize} \quad -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j \cdot \langle \phi(\mathbf{x}_i), \phi(\mathbf{x}_j) \rangle + \sum_i \lambda_i$$

$$\text{subject to} \quad \lambda \succeq 0, \quad \sum_i \lambda_i y_i = 0$$

Kernel Functions

We now define **kernel functions**:

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

Typical kernel functions are:

- Linear:

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

- Polynomial:

$$k(\mathbf{x}, \mathbf{x}') = (\langle \mathbf{x}, \mathbf{x}' \rangle + 1)^k$$

- Radial basis function:

$$k(\mathbf{x}, \mathbf{x}') = e^{-\gamma \|\mathbf{x} - \mathbf{x}'\|_2^2}$$

- Sigmoid kernel:

$$k(\mathbf{x}, \mathbf{x}') = \tanh(\alpha \langle \mathbf{x}, \mathbf{x}' \rangle + \beta)$$

Lessons Learned


- Lagrangian formulation of the hard and soft margin problems
- Langrange dual representation
- Idea of feature transforms



Next Time in Pattern Recognition



Further Readings

- Bernhard Schölkopf, Alexander J. Smola:
[Learning with Kernels](#),
The MIT Press, Cambridge, 2003.
- Vladimir N. Vapnik:
[The Nature of Statistical Learning Theory](#),
Information Science and Statistics, Springer, Heidelberg, 2000.
- S. Boyd, L. Vandenberghe:
[Convex Optimization](#),
Cambridge University Press, 2004.
 <http://www.stanford.edu/~boyd/cvxbook/>
- Christopher M. Bishop:
[Pattern Recognition and Machine Learning](#),
Springer, New York, 2006

Comprehensive Questions

- What is the Lagrangian of the hard margin SVM?
- What are the KKT optimality conditions for the hard margin SVM?
- How do we apply the KKT conditions to the Lagrange Dual?
- What can we conclude from this reformulated Lagrange Dual?