

Machine Learning

Lectures 7 & 8: Support Vector Machines + Kernels

Prof. Dr. Stephan Günnemann

Data Mining and Analytics
Technical University of Munich

4.12.2017 & 11.12.2017

Reading material

Reading material

- Bishop: chapters 7.1.0, 7.1.1, 7.1.2

Acknowledgements

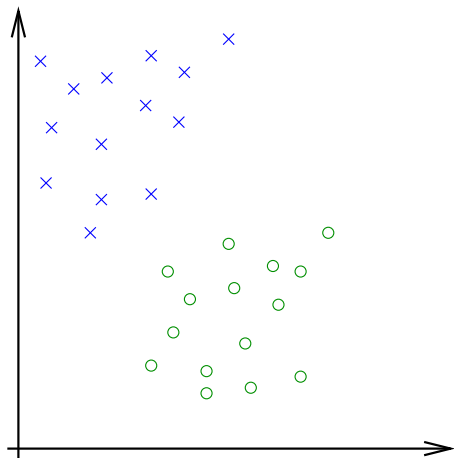
- Slides are based on an older version by S. Urban

Section 1

Linear Classification with a Decision Hyperplane

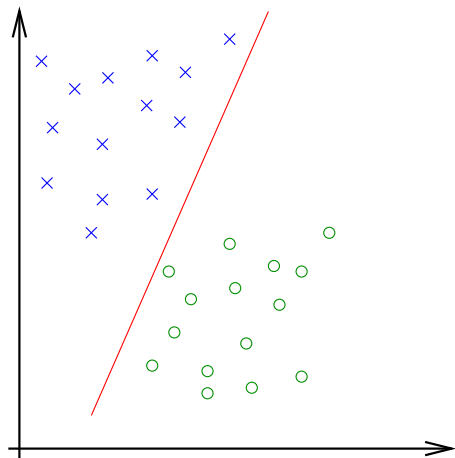
Recall: Linear classification with a hyperplane

- Binary classification problem:
 - × are in class 1
 - are in class -1



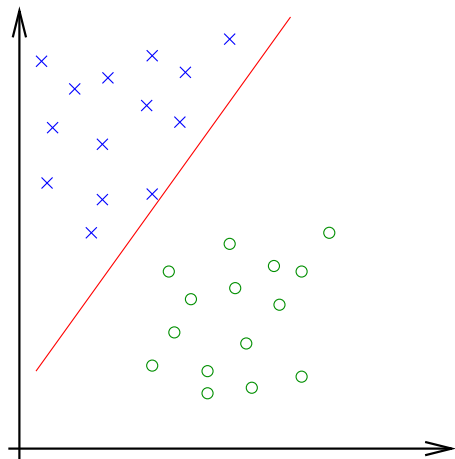
Recall: Linear classification with a hyperplane

- Binary classification problem:
 - × are in class 1
 - are in class -1
- Linear classifier:
Find the optimal decision hyperplane.



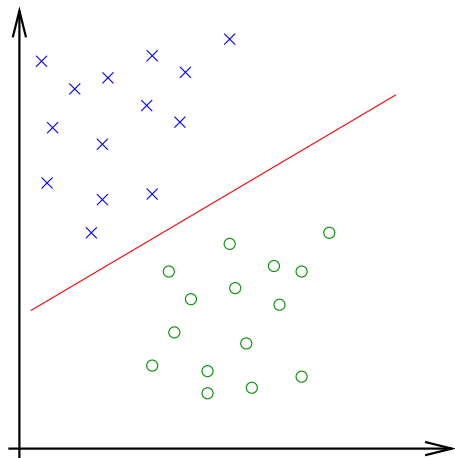
Recall: Linear classification with a hyperplane

- Binary classification problem:
 - × are in class 1
 - are in class -1
- Linear classifier:
Find the optimal decision hyperplane.



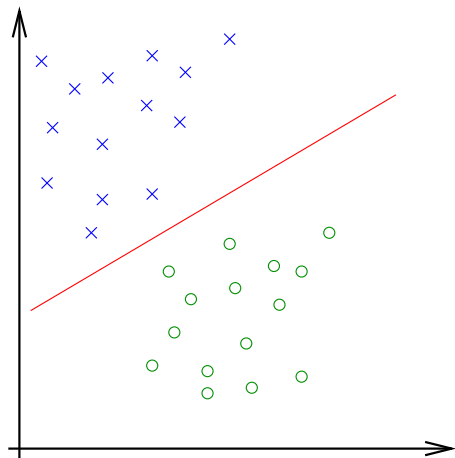
Recall: Linear classification with a hyperplane

- Binary classification problem:
 - × are in class 1
 - are in class -1
- Linear classifier:
Find the optimal decision hyperplane.
- There are infinitely many solutions. Problem is not well-defined.



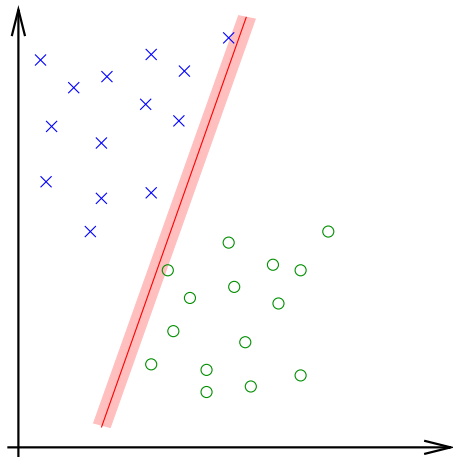
Recall: Linear classification with a hyperplane

- Binary classification problem:
 - × are in class 1
 - are in class -1
- Linear classifier:
Find the optimal decision hyperplane.
- There are infinitely many solutions. Problem is not well-defined.
- A better objective:
Find a separating hyperplane,
such that subsequent points
are classified correctly.



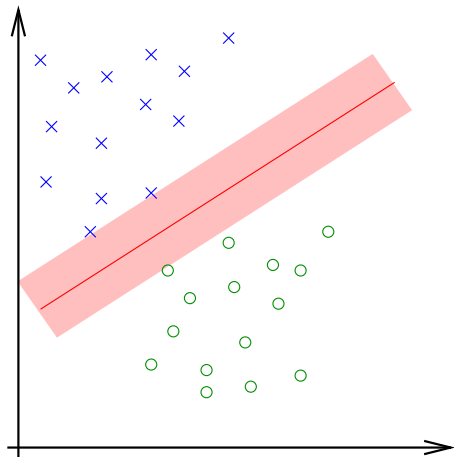
Maximum margin classifier

- Intuitively, a wide margin around the dividing line makes it more likely that new samples will fall on the right side of the boundary.



Maximum margin classifier

- Intuitively, a wide margin around the dividing line makes it more likely that new samples will fall on the right side of the boundary.
- Actual rigorous motivation comes from Statistical Learning Theory ¹
- Objective:
Find a hyperplane that separates both classes with the maximum margin.



¹V. Vapnik - "Statistical Learning Theory", 1995

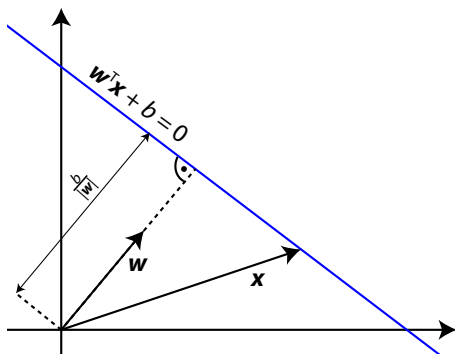
Hyperplanes in Hessian normal form

A hyperplane is defined by a set of all points $x \in \mathbb{R}^D$ that satisfy

$$w^T x + b = 0.$$

Explanation:

- w is the **unit normal vector** and b is the **offset**.
- $w^T x$ is the length of the projection of x on w in units (multiples) of $\|w\|$.



Hyperplanes in Hessian normal form

For a point x_1 that lies on the far side of the hyperplane we have

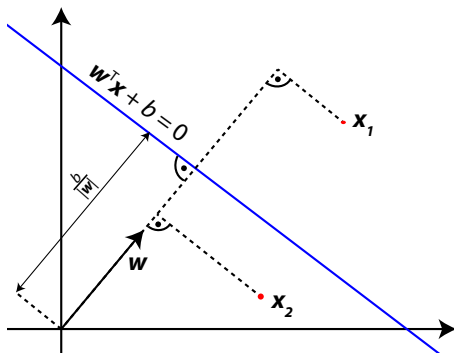
$$w^T x_1 + b > 0$$

because the projection of x_1 on w is *larger* than for points x that are on the hyperplane.

Conversely for a point x_2 that lies on the near side of the hyperplane we have

$$w^T x_2 + b < 0$$

because the projection of x_2 on w is *smaller* than for points x that are on the hyperplane.



Linear classifier

We can use this to build a linear classifier by assigning all x with

$$\mathbf{w}^T \mathbf{x} + b > 0$$

to class **blue** and all x with

$$\mathbf{w}^T \mathbf{x} + b < 0$$

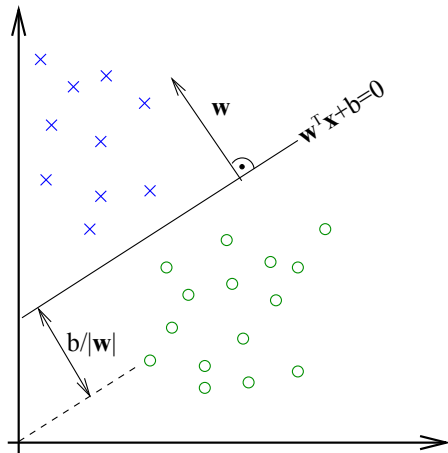
to class **green**.

Thus the class of x is given by

$$h(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b)$$

with

$$\text{sgn}(z) = \begin{cases} -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ +1 & \text{if } z > 0 \end{cases}.$$



Linear classifier with margin

We add two more hyperplanes that are parallel to the original hyperplane and require that no training points must lie between those hyperplanes.

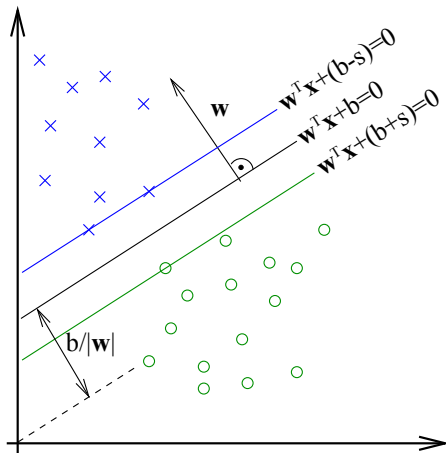
Thus we now require

$$\mathbf{w}^T \mathbf{x} + (b - s) > 0$$

for all \mathbf{x} from class **blue** and

$$\mathbf{w}^T \mathbf{x} + (b + s) < 0$$

for all \mathbf{x} from class **green**.



Size of the margin

Signed distance from the origin to the hyperplane is given by

$$d = -\frac{b}{\|w\|}.$$

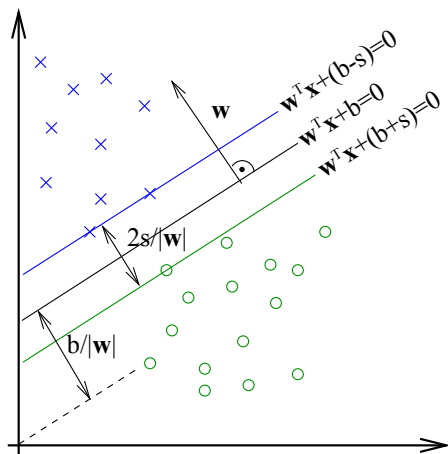
Thus we have

$$d_{blue} = -\frac{b-s}{\|w\|}$$

$$d_{green} = -\frac{b+s}{\|w\|}$$

and the margin is

$$m = d_{blue} - d_{green} = \frac{2s}{\|w\|}.$$

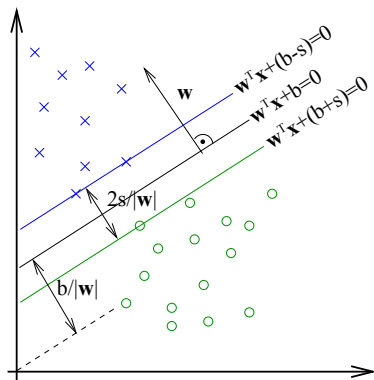


Redundancy of parameter s

The size of the margin,

$$m = \frac{2s}{\|w\|}$$

only depends on the ratio, so w.l.o.g. we can set $s = 1$ and get



Redundancy of parameter s

The size of the margin,

$$m = \frac{2s}{\|w\|}$$

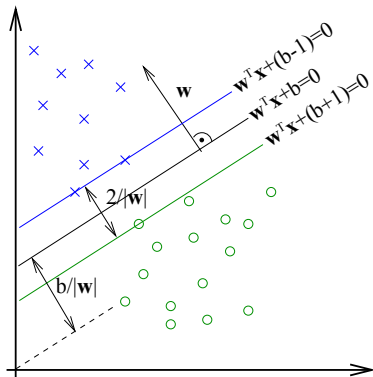
only depends on the ratio, so w.l.o.g. we can set $s = 1$ and get

$$m = \frac{2}{\|w\|}$$

Although the distance from the origin to the black plane,

$$d = -\frac{b}{\|w\|},$$

also depends on two parameters we *cannot* set $b = 1$ as this would link the distance d to the size of the margin m .



Set of constraints

Let \mathbf{x}_i be the i th sample, and $y_i \in \{-1, 1\}$ the class assigned to \mathbf{x}_i .

The constraints

$$\mathbf{w}^T \mathbf{x}_i + b \geq +1 \quad \text{for } y_i = +1,$$

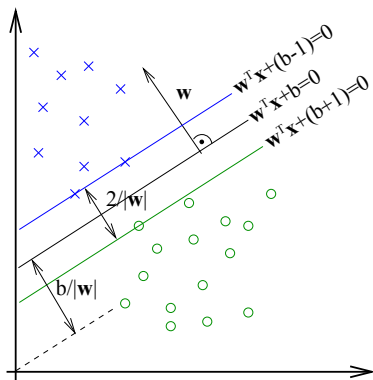
$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 \quad \text{for } y_i = -1$$

can be condensed into

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \text{for all } i.$$

If these constraints are fulfilled the margin is

$$m = \frac{2}{\|\mathbf{w}\|} = \frac{2}{\sqrt{\mathbf{w}^T \mathbf{w}}}.$$



Optimization problem

Let \mathbf{x}_i be the i th data point, $i = 1, \dots, N$, and $y_i \in \{-1, 1\}$ the class assigned to \mathbf{x}_i .

To find the separating hyperplane with the maximum margin we need to find $\{\mathbf{w}, b\}$ that

$$\begin{aligned} &\text{minimize} && f_0(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ &\text{subject to} && f_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \quad \text{for } i = 1, \dots, N. \end{aligned}$$

This is a **constrained convex optimization problem** (more specifically, quadratic programming problem).

We don't yet have the tools to solve this kind of problem, so we make a detour to develop them.

We go from $\|\mathbf{w}\|$ to $\|\mathbf{w}\|^2 = \mathbf{w}^T \mathbf{w}$ as square root is a monotonic function that doesn't change the location of the optimum.

Section 2

Constrained Optimization

Optimization with inequality constraints

Constrained optimization problem

Given $f_0 : \mathbb{R}^D \rightarrow \mathbb{R}$ and $f_i : \mathbb{R}^D \rightarrow \mathbb{R}$,

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, M.\end{array}$$

Feasibility

A point $\mathbf{x} \in \mathbb{R}^D$ is called **feasible** if and only if it satisfies all constraints $f_i(\mathbf{x}) \leq 0, i = 1, \dots, M$ of the optimization problem.

Minimum and minimizer

We call the optimal value the **minimum** p^* , and the point where the minimum is obtained the **minimizer** \mathbf{x}^* . Thus $p^* = f(\mathbf{x}^*)$.

Minimization with a single inequality constraint

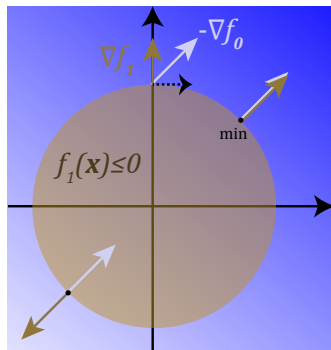
$$\begin{aligned} &\text{minimize} && f_0(\mathbf{x}) = -(x_1 + x_2) \\ &\text{subject to} && f_1(\mathbf{x}) = x_1^2 + x_2^2 - 1 \leq 0 \end{aligned}$$

- If the minimum is in the interior of the circle ($f_1(\mathbf{x}^*) < 0$), we have reached it when $\nabla f_0(\mathbf{x}^*) = 0$.
- If the minimum is on the circle ($f_1(\mathbf{x}^*) = 0$), we have reached it when the negative gradient $-\nabla f_0(\mathbf{x}^*)$ has no component that we could follow without changing the value of f_1 .

Thus at the minimizer \mathbf{x}^* we have

$$-\nabla f_0(\mathbf{x}^*) = \alpha \nabla f_1(\mathbf{x}^*)$$

with $\alpha \geq 0$ to ensure that $-\nabla f_0$ and ∇f_1 do not point in opposite directions.



$f_0(\mathbf{x})$ is color coded.

Multiple inequality constraints

Given $f_0 : \mathbb{R}^D \rightarrow \mathbb{R}$ and $f_i : \mathbb{R}^D \rightarrow \mathbb{R}$,

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, M.\end{array}$$

For multiple constraints f_1, \dots, f_p we have reached the minimum when the negative gradient $-\nabla f_0$ has no component that we could follow without changing the value of **any** constraint.

This is the case when $-\nabla f_0$ is a linear combination of the ∇f_i 's,

$$-\nabla f_0(\mathbf{x}^*) = \sum_{i=1}^M \alpha_i \nabla f_i(\mathbf{x}^*), \quad \alpha_i \geq 0.$$

Lagrangian

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, M\end{array}$$

Definition (Lagrangian)

We define the **Lagrangian** $L : \mathbb{R}^D \times \mathbb{R}^M \rightarrow \mathbb{R}$ associated with the above problem as

$$L(\mathbf{x}, \boldsymbol{\alpha}) = f_0(\mathbf{x}) + \sum_{i=1}^M \alpha_i f_i(\mathbf{x})$$

We refer to $\alpha_i \geq 0$ as the **Lagrange multiplier** associated with the inequality constraint $f_i(\mathbf{x}) \leq 0$.

Calculating the gradient of L w.r.t. \mathbf{x} shows that our optimality criterion at \mathbf{x}^* is recovered,

$$\nabla_{\mathbf{x}^*} L(\mathbf{x}^*, \boldsymbol{\alpha}) = \nabla f_0(\mathbf{x}^*) + \sum_{i=1}^M \alpha_i \nabla f_i(\mathbf{x}^*) = 0.$$

Interpretation of the Lagrangian

$$\begin{array}{ll} \text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0, \quad i = 1, \dots, M \end{array} \quad L(\mathbf{x}, \boldsymbol{\alpha}) = f_0(\mathbf{x}) + \sum_{i=1}^M \alpha_i f_i(\mathbf{x})$$

For every choice of $\boldsymbol{\alpha}$, the corresponding $\min_{\mathbf{x} \in \mathbb{R}^D} L(\mathbf{x}, \boldsymbol{\alpha})$ is a **lower bound** on the optimal value of the constrained problem.

Watch `lagrangian.avi` for a demonstration.

Lagrange dual function

Definition (Lagrange dual function)

The **Lagrange dual function** $g : \mathbb{R}^M \rightarrow \mathbb{R}$ is the minimum of the Lagrangian over \mathbf{x} given $\boldsymbol{\alpha}$,

$$g(\boldsymbol{\alpha}) = \min_{\mathbf{x} \in \mathbb{R}^D} L(\mathbf{x}, \boldsymbol{\alpha}) = \min_{\mathbf{x} \in \mathbb{R}^D} \left(f_0(\mathbf{x}) + \sum_{i=1}^M \alpha_i f_i(\mathbf{x}) \right).$$

It is concave in $\boldsymbol{\alpha}$ since it is the point-wise minimum of a family of affine functions of $\boldsymbol{\alpha}$.

Lagrange dual problem

For each $\alpha \geq 0$ the Lagrange dual function $g(\alpha)$ gives us a lower bound on the optimal value p^* of the original optimization problem.

What is the best (highest) lower bound?

Lagrange dual problem

$$\begin{aligned} & \text{maximize} && g(\alpha) = \min_x L(x, \alpha) \\ & \text{subject to} && \alpha_i \geq 0, \quad i = 1, \dots, m \end{aligned}$$

The maximum d^* of the Lagrange dual problem is the best lower bound on p^* that we can achieve by using the Lagrangian.

Dualities

Weak duality (always)

Since $g(\alpha)$ is a lower bound of p^* we have **weak duality**,

$$d^* \leq p^* .$$

The difference between the solution of the original and the dual problem,

$$p^* - d^* \geq 0 ,$$

is called the **duality gap**.

Strong duality (under certain conditions)

Under certain conditions we have

$$d^* = p^* ,$$

i.e. the solution to the Lagrange dual problem is a solution of the original constrained optimization problem.

We then say that **strong duality** holds.

Dual solution

Let \mathbf{x}^* be a minimizer of the primal problem and $\boldsymbol{\alpha}^*$ a maximizer of the dual problem. If strong duality holds, then

$$L(\mathbf{x}^*, \boldsymbol{\alpha}^*) = g(\boldsymbol{\alpha}^*) = \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}^*).$$

Proof.

Let $\tilde{\mathbf{x}}$ be a minimizer of $L(\mathbf{x}, \boldsymbol{\alpha}^*)$ over \mathbf{x} . We have

$$d^* = g(\boldsymbol{\alpha}^*) = \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}^*) = L(\tilde{\mathbf{x}}, \boldsymbol{\alpha}^*) \quad (\text{minimizer})$$

$$= f_0(\tilde{\mathbf{x}}) + \sum_{i=1}^M \alpha_i^* f_i(\tilde{\mathbf{x}}) \quad (\text{Lagrangian})$$

$$\leq f_0(\mathbf{x}^*) + \sum_{i=1}^M \alpha_i^* f_i(\mathbf{x}^*) \quad (\tilde{\mathbf{x}} \text{ minimizer over } \mathbf{x})$$

$$\leq f_0(\mathbf{x}^*) = p^* \quad (\boldsymbol{\alpha} \geq \mathbf{0}; \mathbf{x}^* \text{ feasible: } f_i(\mathbf{x}^*) \leq 0)$$

and since $d^* = p^*$, we have $g(\boldsymbol{\alpha}^*) = L(\tilde{\mathbf{x}}, \boldsymbol{\alpha}^*) = L(\mathbf{x}^*, \boldsymbol{\alpha}^*)$. □

Corollary

Let \mathbf{x}^ be a minimizer of the primal problem and $\boldsymbol{\alpha}^*$ a maximizer of the dual problem. If strong duality holds and $L(\mathbf{x}, \boldsymbol{\alpha})$ is convex in \mathbf{x} , then*

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}^*).$$

Constraint qualifications for convex problems

Consider the constrained optimization problem,

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0.\end{array}$$

Slater's constraint qualification

If f_0, f_1, \dots, f_M are **convex** and there exists an $\mathbf{x} \in \mathbb{R}^D$ such that

$$f_i(\mathbf{x}) < 0, \quad i = 1, \dots, M$$

or the constraints are affine, that is

$$f_i(\mathbf{x}) = \mathbf{w}_i^T \mathbf{x} + b_i \leq 0,$$

the duality gap is zero.

For a proof see Boyd p. 234.

Recipe for solving constrained optimization problems

The constrained optimization problem,

$$\begin{aligned} & \text{minimize} && f_0(\mathbf{x}) \\ & \text{subject to} && f_i(\mathbf{x}) \leq 0 \quad i = 1, \dots, M, \end{aligned}$$

with f_0 **convex** and f_1, \dots, f_M **convex** can be solved as follows:

1. Calculate the **Lagrangian**

$$L(\mathbf{x}, \boldsymbol{\alpha}) = f_0(\mathbf{x}) + \sum_{i=1}^M \alpha_i f_i(\mathbf{x}).$$

2. Obtain the **Lagrange dual function** $g(\boldsymbol{\alpha})$ by solving

$$\mathbf{x}^* = \arg \min_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}) \stackrel{\text{convex}}{\iff} \nabla_{\mathbf{x}} L(\mathbf{x}, \boldsymbol{\alpha}) = \mathbf{0}.$$

3. Solve the dual problem

$$\begin{aligned} & \text{maximize} && g(\boldsymbol{\alpha}) = L(\mathbf{x}^*, \boldsymbol{\alpha}) \\ & \text{subject to} && \alpha_i \geq 0 \quad i = 1, \dots, M. \end{aligned}$$

The solution of this problem is also the solution of the original problem if Slater's condition is satisfied.

KKT conditions

If we have the following optimization problem

$$\begin{array}{ll}\text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq 0,\end{array}$$

with f_0 **convex** and f_1, \dots, f_M **convex** and strong duality holds, the following statement is true:

KKT conditions

\mathbf{x}^* and $\boldsymbol{\alpha}^*$ are the optimal solutions of the constrained optimization problem and the corresponding Lagrange dual problem if and only if they satisfy the **Karush-Kuhn-Tucker (KKT)** conditions:

$f_i(\mathbf{x}^*) \leq 0$	primal feasibility,
$\alpha_i^* \geq 0$	dual feasibility,
$\alpha_i^* f_i(\mathbf{x}^*) = 0$	complementary slackness,
$\nabla_{\mathbf{x}} L(\mathbf{x}^*, \boldsymbol{\alpha}^*) = 0$	\mathbf{x}^* minimizes Lagrangian,

for $i = 1, \dots, M$. (Proof in the tutorial)

Summary

- Lagrange formalism reformulates constrained optimization problem into a dual problem.
- For convex objective and constraints, solving the dual problem allows to solve the primal problem (duality gap is zero).
- KKT conditions characterize the solution and sometimes allow to obtain it without minimization.
- Convexity is necessary to guarantee optimality.

Section 3

Back to Support Vector Machines

Recall: Optimization problem

Let \mathbf{x}_i be the i th data point, $i = 1, \dots, N$, and $y_i \in \{-1, 1\}$ the class assigned to \mathbf{x}_i .

To find the separating hyperplane with the maximum margin we need to find $\{\mathbf{w}, b\}$ that

$$\begin{aligned} &\text{minimize} && f_0(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \\ &\text{subject to} && f_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \quad \text{for } i = 1, \dots, N. \end{aligned}$$

This is a **constrained convex optimization problem**.

Primal problem

We apply our recipe for solving the constrained optimization problem.

1. Calculate the Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\alpha}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1].$$

2. Minimize $L(\mathbf{w}, b, \boldsymbol{\alpha})$ w.r.t. \mathbf{w} and b .

$$\nabla_{\mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\alpha}) = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i \stackrel{!}{=} 0$$

$$\frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i \stackrel{!}{=} 0$$

Thus the weights are a linear combination of the training samples,

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i.$$

Dual problem

Substituting both relations back into $L(\mathbf{w}, b, \boldsymbol{\alpha})$ gives the Lagrange dual function $g(\boldsymbol{\alpha})$.

Thus we have reformulated our original problem as

$$\begin{aligned} \text{maximize} \quad & g(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad \text{for } i = 1, \dots, N. \end{aligned}$$

3. Solve this problem.

Solving the dual problem

We can rewrite the dual function $g(\alpha)$ in vector form

$$g(\alpha) = \frac{1}{2} \alpha^T Q \alpha + \alpha^T \mathbf{1}_N$$

where Q is a symmetric negative (semi-)definite matrix, and the constraints on α are linear.

This is an instance of a [quadratic programming](#) problem. There exist efficient algorithms for its solution, such as [Sequential minimal optimization](#) (SMO) ².

A number of implementations, such as LIBSVM ³ are available and are widely used in practice.

²<http://cs229.stanford.edu/materials/smo.pdf>

³C.-C. Chang and C.-J. Lin. *LIBSVM : a library for support vector machines*, 2011

Recovering \mathbf{w} and b from the dual solution α^*

Having obtained the optimal α^* using our favorite QP solver, we can compute the parameters defining the separating hyperplane.

Recall, that from the optimality condition, the weights \mathbf{w} are a linear combination of the training samples,

$$\mathbf{w} = \sum_{i=1}^N \alpha_i^* y_i \mathbf{x}_i$$

From the complementary slackness condition $\alpha_i^* f_i(\mathbf{x}) = 0$ we can easily recover the bias.

When we take any vector \mathbf{x}_i for which $\alpha_i \neq 0$. The corresponding constraint $f_i(\mathbf{w}, b)$ must be zero and thus we have

$$\mathbf{w}^T \mathbf{x}_i + b = y_i .$$

Solving this for b yields the bias

$$b = y_i - \mathbf{w}^T \mathbf{x}_i$$

We can also average the b over all support vectors to get a more stable solution.

Support vectors

Let's look closer at the complimentary slackness condition

$$\alpha_i f_i(\mathbf{x}^*) = 0.$$

In our case this means

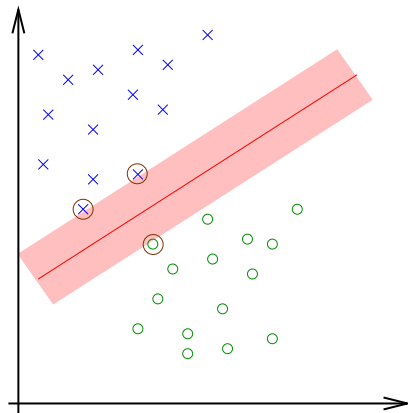
$$\alpha_i [y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1] = 0 \quad \text{for all } i.$$

Hence a training sample \mathbf{x}_i can only contribute to the weight vector ($\alpha_i \neq 0$) if it lies **on the margin**, that is

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

A training sample \mathbf{x}_i with $\alpha_i \neq 0$ is called a **support vector**.

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$



Classifying

The class of \mathbf{x} is given by

$$h(\mathbf{x}) = \text{sgn}(\mathbf{w}^T \mathbf{x} + b).$$

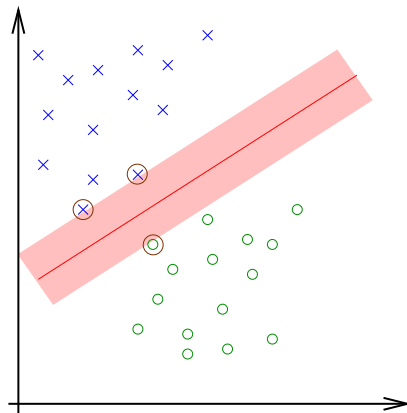
Substituting

$$\mathbf{w} = \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i$$

gives

$$h(\mathbf{x}) = \text{sgn}\left(\sum_{i=1}^N \alpha_i y_i \mathbf{x}_i^T \mathbf{x} + b\right).$$

Since the solution is **sparse** (most α_i s are zero) we only need to remember the few training samples \mathbf{x}_i with $\alpha_i \neq 0$.



Section 4

Soft Margin Support Vector Machines

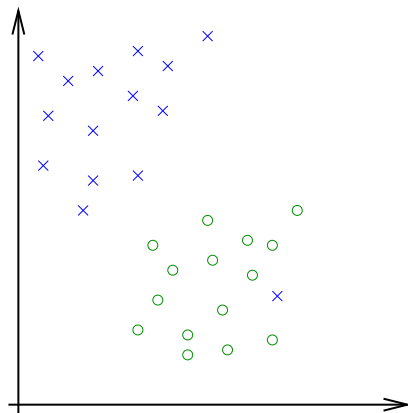
Dealing with noisy data

What if the data is not linearly separable due to some noise?

With our current version of SVM, a single outlier makes the constraint unsatisfiable.

The corresponding Lagrange multiplier α_i would go to infinity and destroy the solution.

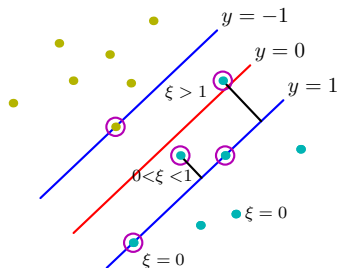
How to make our model more robust?



Slack variables

Idea: Relax the constraints as much as necessary but punish the relaxation of a constraint.

We introduce a **slack variable** $\xi_i \geq 0$ for every training sample x_i that gives the distance of how far the margin is violated by this training sample in units of $\|w\|$.



Slack variables

Idea: Relax the constraints as much as necessary but punish the relaxation of a constraint.

We introduce a **slack variable** $\xi_i \geq 0$ for every training sample x_i that gives the distance of how far the margin is violated by this training sample in units of $\|w\|$.

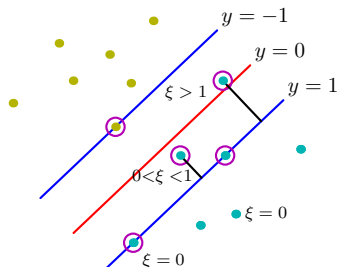
Hence our relaxed constraints are

$$w^T x_i + b \geq +1 - \xi_i \quad \text{for } y_i = +1,$$

$$w^T x_i + b \leq -1 + \xi_i \quad \text{for } y_i = -1.$$

Again, they can be condensed into

$$y_i(w^T x_i + b) \geq 1 - \xi_i \quad \text{for all } i.$$



Slack variables

Idea: Relax the constraints as much as necessary but punish the relaxation of a constraint.

We introduce a **slack variable** $\xi_i \geq 0$ for every training sample \mathbf{x}_i that gives the distance of how far the margin is violated by this training sample in units of $\|\mathbf{w}\|$.

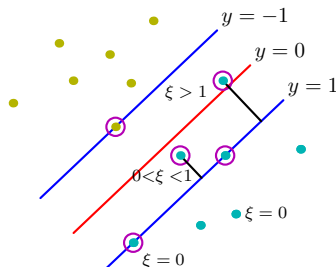
Hence our relaxed constraints are

$$\mathbf{w}^T \mathbf{x}_i + b \geq +1 - \xi_i \quad \text{for } y_i = +1,$$

$$\mathbf{w}^T \mathbf{x}_i + b \leq -1 + \xi_i \quad \text{for } y_i = -1.$$

Again, they can be condensed into

$$y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \text{for all } i.$$



The new cost function is,

$$f_0(\mathbf{w}, b, \boldsymbol{\xi}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i.$$

The factor $C > 0$ determines how heavy a violation is punished.

$C \rightarrow \infty$ recovers hard-margin SVM.

Optimization problem with slack variables

Let \mathbf{x}_i be the i th data point, $i = 1, \dots, N$, and $y_i \in \{-1, 1\}$ the class assigned to \mathbf{x}_i . Let $C > 0$ be a constant.

To find the hyperplane that separates most of the data with maximum margin we

$$\begin{aligned} \text{minimize} \quad & f_0(\mathbf{w}, b, \boldsymbol{\xi}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i \geq 0 & i = 1, \dots, N, \\ & \xi_i \geq 0 & i = 1, \dots, N. \end{aligned}$$

Here we used the 1-norm for the penalty term $\sum_i \xi_i$. Another choice is to use the 2-norm penalty, $\sum_i \xi_i^2$.

The penalty that performs better in practice will depend on the data and the type of noise that has influenced it.

Lagrangian with slack variables

1. Calculate the Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ - \sum_{i=1}^N \alpha_i [y_i (\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i] - \sum_{i=1}^N \mu_i \xi_i.$$

2. Minimize $L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu})$ w.r.t. \mathbf{w} , b and $\boldsymbol{\xi}$.

$$\nabla_{\mathbf{w}} L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\mu}) = \mathbf{w} - \sum_{i=1}^N \alpha_i y_i \mathbf{x}_i = 0; \quad \frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i = 0, \\ \frac{\partial L}{\partial \xi_i} = C - \alpha_i - \mu_i = 0 \quad \text{for } i = 1, \dots, N$$

From $\alpha_i = C - \mu_i$ and dual feasibility $\mu_i \geq 0$, $\alpha_i \geq 0$ we get

$$0 \leq \alpha_i \leq C.$$

Dual problem with slack variables

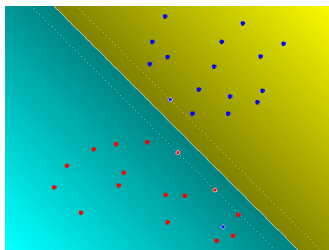
This leads to the dual problem:

$$\begin{aligned} \text{maximize} \quad & g(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C \quad i = 1, \dots, N. \end{aligned}$$

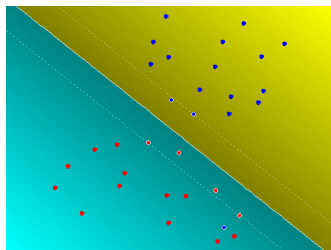
This is nearly the same dual problem as for the case without slack variables.

Only the constraint $\alpha_i \leq C$ is new. It ensures that α_i is bounded and cannot go to infinity.

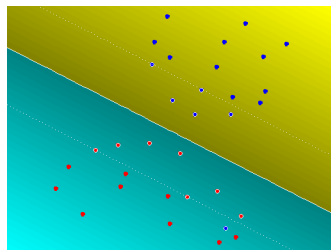
Influence of the penalty C



$C = 100$



$C = 10$



$C = 1$

Hinge loss formulation

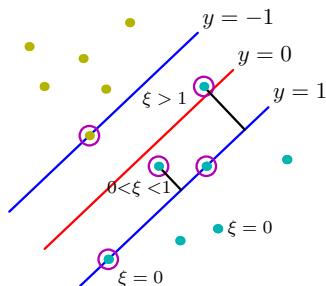
We can have another look at our **constrained** optimization problem.

$$\begin{aligned} \text{minimize} \quad & f_0(\mathbf{w}, b, \boldsymbol{\xi}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \xi_i \\ \text{subject to} \quad & y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i \geq 0, \quad i = 1, \dots, N, \\ & \xi_i \geq 0, \quad i = 1, \dots, N. \end{aligned}$$

Clearly, for the optimal solution the slack variables are

$$\xi_i = \begin{cases} 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b), & \text{if } y_i(\mathbf{w}^T \mathbf{x}_i + b) < 1 \\ 0 & \text{else} \end{cases}$$

since we are minimizing over $\boldsymbol{\xi}$.



Hinge loss formulation

Thus, we can rewrite the objective function as an **unconstrained** optimization problem known as the **hinge loss** formulation

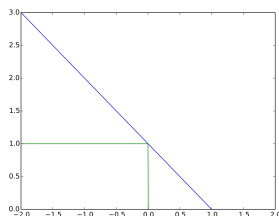
$$\min_{\mathbf{w}, b} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^N \max\{0, 1 - y_i(\mathbf{w}^T \mathbf{x}_i + b)\}$$

The **hinge loss** function

$E_{\text{hinge}}(z) = \max\{0, 1 - z\}$ penalizes the points that lie within the margin.

The name comes from the shape from the function, as can be seen in the figure to the right.

We can optimize this hinge loss objective directly, using standard gradient-based methods.



Hinge loss (**blue**) can be viewed as an approximation to zero-one loss (**green**).

Section 5

Kernels

Feature space

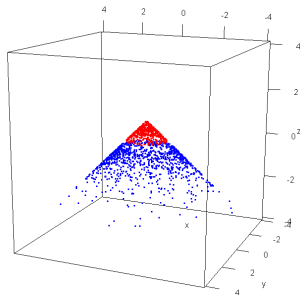
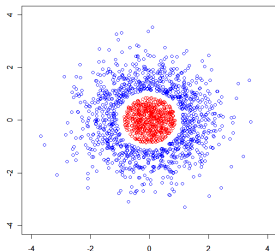
So far we can only construct linear classifiers.

Before, we used **basis functions** $\phi(\cdot)$ to make the models nonlinear

$$\phi : \mathbb{R}^D \rightarrow \mathbb{R}^M \quad \mathbf{x}_i \mapsto \phi(\mathbf{x}_i)$$

For example, with the following mapping the data becomes linearly separable

$$\phi(x, y) = \begin{pmatrix} x \\ y \\ -\sqrt{x^2 + y^2} \end{pmatrix}$$



Kernel trick

In the dual formulation of SVM, the samples \mathbf{x}_i only enter the dual objective as **inner products** $\mathbf{x}_i^T \mathbf{x}_j$

$$g(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{x}_i^T \mathbf{x}_j ,$$

For basis functions this means that

$$g(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

Kernel trick

We can define a **kernel function** $k : \mathbb{R}^D \times \mathbb{R}^D \rightarrow \mathbb{R}$

$$k(\mathbf{x}_i, \mathbf{x}_j) := \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

and rewrite the dual as

$$g(\boldsymbol{\alpha}) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j k(\mathbf{x}_i, \mathbf{x}_j)$$

This operation is referred to as **the kernel trick**.

It can be used not only for SVM. Kernel trick can be used in any model that can be formulated such that it only depends on the inner products $\mathbf{x}_i^T \mathbf{x}_j$. (e.g. linear regression, k-nearest neighbors)

Kernel trick

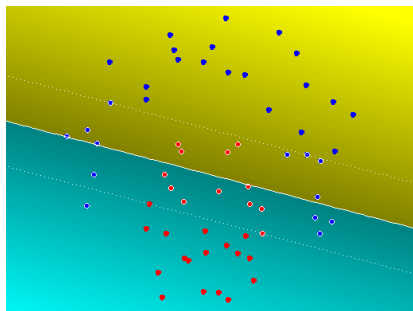
The kernel represent an inner product in the **feature space** spanned by ϕ . Like before, this makes our models non-linear w.r.t. the data space.

What's the point of using kernels if we can simply use the basis functions?

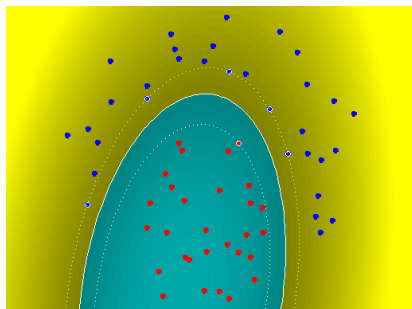
- Some kernels are equivalent to using infinite-dimensional basis functions. While computing these feature transformations would be impossible, directly evaluating the kernels is often easy.
- Kernels can be used to encode similarity between arbitrary non-numerical data, from strings to graphs.
For example, we could define

$$k(\text{lemon}, \text{orange}) = 10 \quad \text{and} \quad k(\text{apple}, \text{orange}) = -5$$

SVM using kernels example



Linear kernel (no kernel)
 $a^T b$



2nd order polynomial kernel
 $(a^T b)^2$

What makes a valid kernel?

A kernel is **valid** if it corresponds to an inner product in some feature space. An equivalent formulation is given by Mercer's theorem.

Mercer's theorem

A kernel is valid if it gives rise to a **positive (semi)-definite** kernel matrix \mathbf{K} for any input data \mathbf{X} .

Kernel matrix (also known as **Gram matrix**) $\mathbf{K} \in \mathbb{R}^{N \times N}$ is defined as

$$\mathbf{K} = \begin{pmatrix} k(\mathbf{x}_1, \mathbf{x}_1) & k(\mathbf{x}_1, \mathbf{x}_2) & \cdots & k(\mathbf{x}_1, \mathbf{x}_N) \\ k(\mathbf{x}_2, \mathbf{x}_1) & k(\mathbf{x}_2, \mathbf{x}_2) & \cdots & k(\mathbf{x}_2, \mathbf{x}_N) \\ \vdots & \vdots & \ddots & \vdots \\ k(\mathbf{x}_N, \mathbf{x}_1) & k(\mathbf{x}_N, \mathbf{x}_2) & \cdots & k(\mathbf{x}_N, \mathbf{x}_N) \end{pmatrix}$$

What happens if we use a non-valid kernel?

Our optimization problem might become non-convex, so we may not get a globally optimal solution.

Examples of kernels

- Polynomial:

$$k(\mathbf{a}, \mathbf{b}) = (\mathbf{a}^T \mathbf{b})^p \quad \text{or} \quad (\mathbf{a}^T \mathbf{b} + 1)^p$$

- Gaussian kernel:

$$k(\mathbf{a}, \mathbf{b}) = \exp \left(-\frac{\|\mathbf{a} - \mathbf{b}\|^2}{2\sigma^2} \right)$$

- Sigmoid:

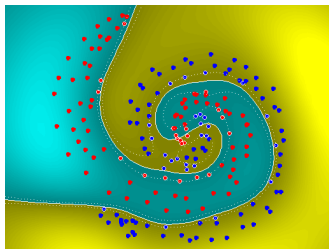
$$k(\mathbf{a}, \mathbf{b}) = \tanh(\kappa \mathbf{a}^T \mathbf{b} - \delta) \quad \text{for } \kappa, \delta > 0$$

In fact, the sigmoid kernel **is not** PSD, but still works well in practice.

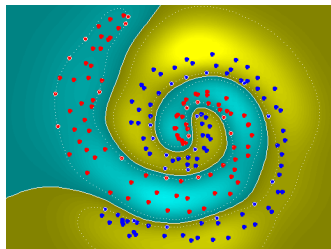
Some kernels introduce additional hyperparameters, that affect the behavior of the algorithm.

Note, that the sigmoid kernel is different from the sigmoid function from *Linear Classification*.

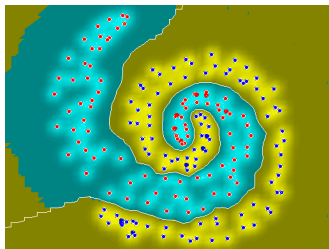
Gaussian kernel ($C=1000$)



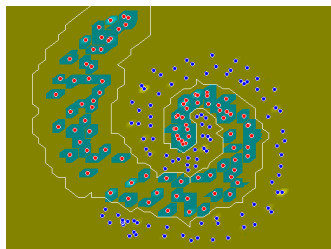
$\sigma = 0.5$



$\sigma = 0.25$



$\sigma = 0.05$



$\sigma = 0.005$

Classifying a new point with kernelized SVM

We denote the set of support vectors (points \mathbf{x}_i for which it holds $0 < \alpha_i < C$ and $\xi_i = 0$) as \mathcal{S} .

From the complementary slackness condition, they must satisfy

$$y_i \left(\sum_{\{j|\mathbf{x}_j \in \mathcal{S}\}} \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) + b \right) = 1$$

Like for the regular SVM, bias can be recovered as

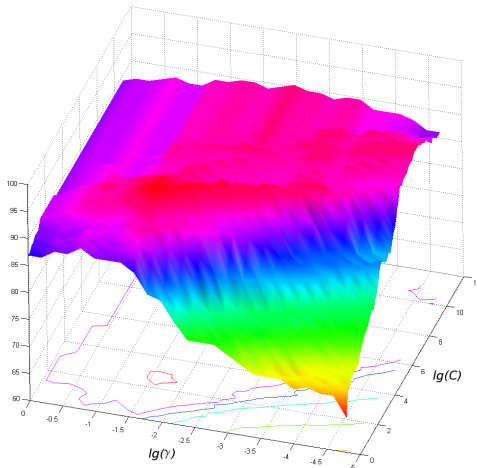
$$b = y_i - \left(\sum_{\{j|\mathbf{x}_j \in \mathcal{S}\}} \alpha_j y_j k(\mathbf{x}_i, \mathbf{x}_j) \right)$$

Thus, a new point \mathbf{x} can be classified as

$$h(\mathbf{x}) = \text{sgn} \left(\sum_{i=1}^N \alpha_i y_i k(\mathbf{x}_i, \mathbf{x}) + b \right)$$

How to choose the hyperparameters?

The best setting of the penalty parameter C and the kernel hyperparameters (e.g., γ or σ) can be determined by performing **cross validation** with **random search** over the parameter space.



Multiple classes

The standard SVM model cannot handle multiclass data.

Two approaches to address this issue are:

One-vs-rest: Train C SVM models for C classes, where each SVM is being trained for classification of one class against all the remaining ones. The winner is then the class, where the distance from the hyperplane is maximal.

One-vs-one: Train $\binom{C}{2}$ classifiers (all possible pairings) and evaluate all. The winner is the class with the majority vote (votes are weighted according to the distance from the margin).

Summary

- Support Vector Machine is a linear binary classifier that, motivated by learning theory, maximizes the margin between the two classes.
- The SVM objective is convex, so a globally optimal solution can be obtained.
- The dual formulation is a quadratic programming problem, that can be solved efficiently, e.g. using standard QP libraries.
- Soft-margin SVM with slack variables can deal with noisy data. Smaller values for the penalty C lead to a larger margin at the cost of misclassifying more samples.
- Linear soft-margin SVM (= hinge loss formulation) can be solved directly using gradient descent.
- We can obtain a nonlinear decision boundary by moving to an implicit high-dimensional feature space by using the kernel trick. This only works in the dual formulation.