

Machine Learning

Lecture 9: SVM and Kernels

Prof. Dr. Stephan Günnemann
Aleksandar Bojchevski

Data Analytics and Machine Learning Group
Technical University of Munich

Winter term 2020/2021

Roadmap

1. Support Vector Machines (SVM)
2. Soft Margin Support Vector Machines
3. Kernels

Section 1

Support Vector Machines (SVM)

Linear classifier

$w \rightarrow$ rotation

$b/|w| \rightarrow$ up/down

A linear classifier assigns all x with

$$w^T x + b > 0 \quad y = +1$$

to class blue and all x with

$$w^T x + b < 0 \quad y = -1$$

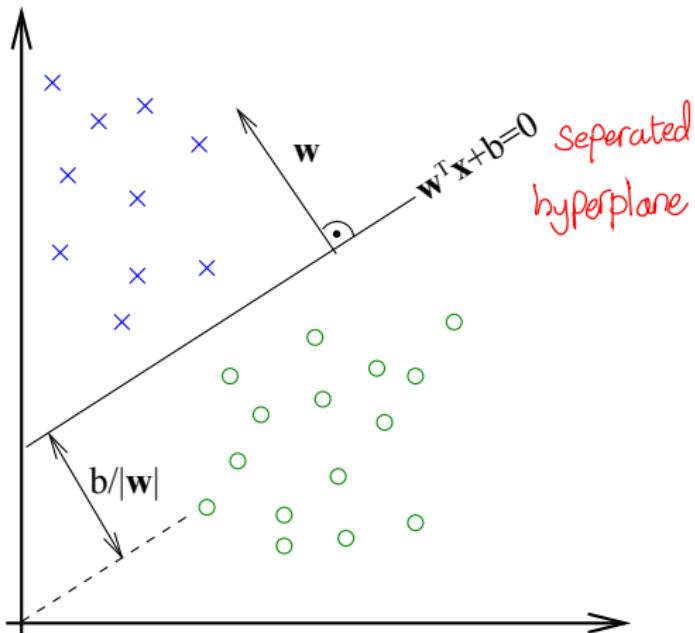
to class green.

Thus the class of x is given by

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

with

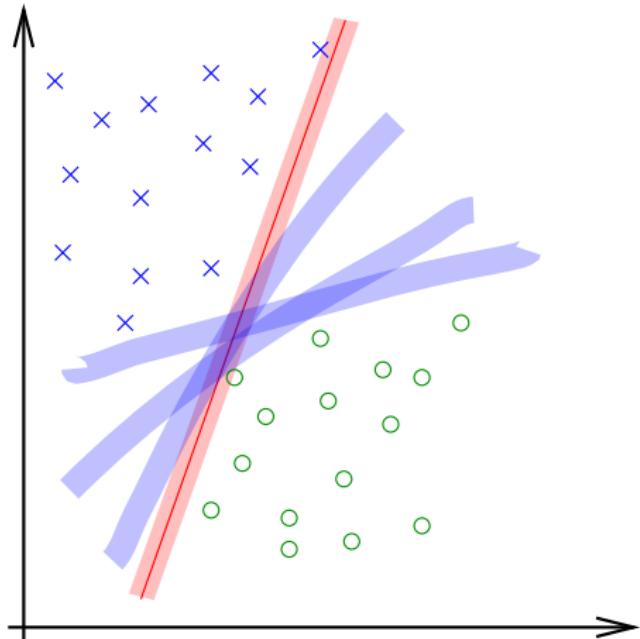
$$\text{sign}(z) = \begin{cases} -1 & \text{if } z < 0 \\ 0 & \text{if } z = 0 \\ +1 & \text{if } z > 0 \end{cases}$$



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.

Linearly separable data guarantee
at least one hyperplane



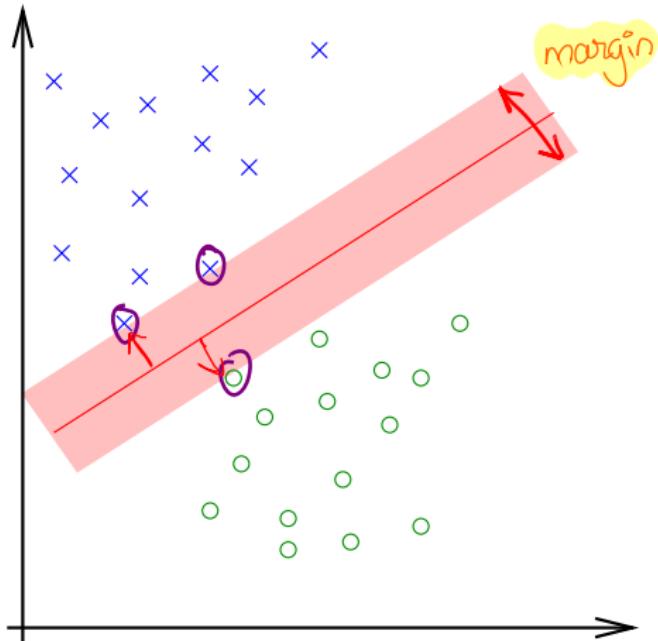
Maximum margin classifier

- 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.

i.e low generalization error



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

VC dim relates margin & error

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.

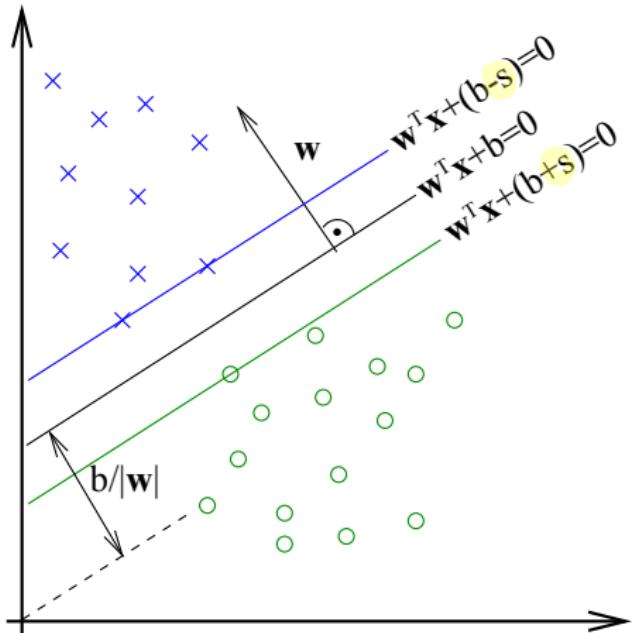
Thus we now require

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

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

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

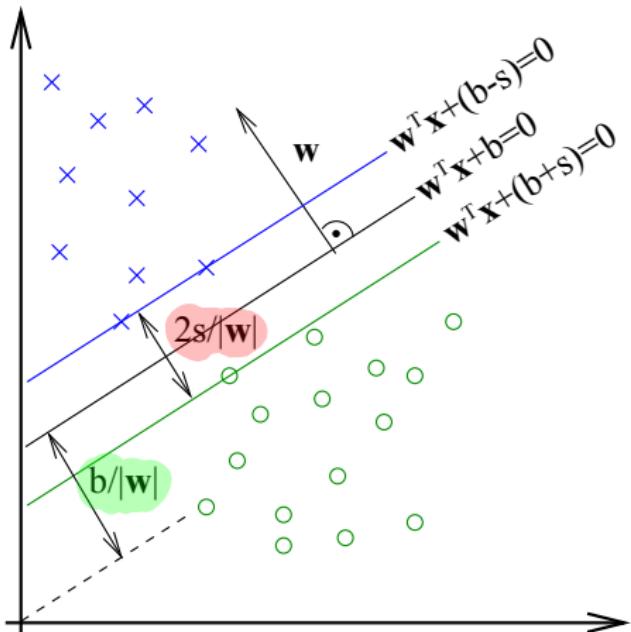
Thus we have

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

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

and the margin is

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

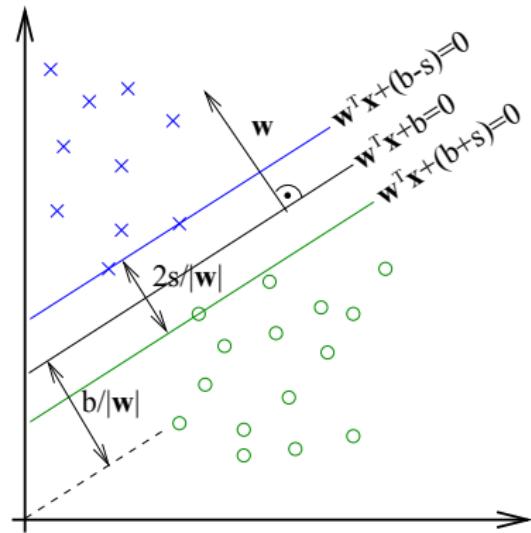


Redundancy of parameter s

The size of the margin,

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

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,

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

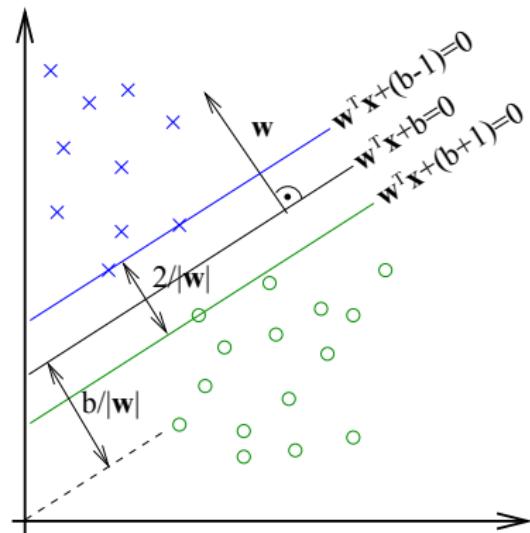
only depends on the ratio, so w.l.o.g. we can set $s = 1$ and get

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

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 .



Hyperplane is only determined by direction of w !

Set of constraints

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

The constraints → shall hold

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

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

can be condensed into

enforcing the
correct classification
of data
points

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

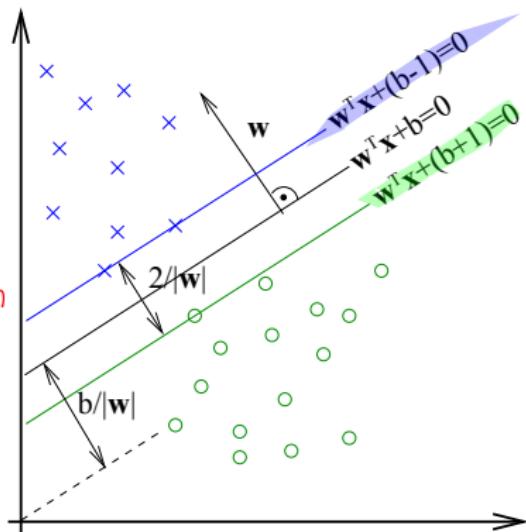
* If these constraints are fulfilled the margin is

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

Solving :

$$\max_{\mathbf{w}, b} \frac{2}{\|\mathbf{w}\|}$$

s.t constraints



SVM's Optimization problem

Hard Margin SVM

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

$$\text{minimize } f_0(\mathbf{w}, b) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{subject to } f_i(\mathbf{w}, b) = y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 \geq 0 \quad \text{for } i = 1, \dots, N.$$

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

* Check extra notes

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.

Optimization with inequality constraints

Constrained optimization problem

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

minimize $_{\theta}$ $f_0(\theta)$

$$\Theta = \{\omega, b\}$$

subject to $f_i(\theta) \leq 0$ for $i = 1, \dots, M$.

Feasibility

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

Minimum and minimizer

We call the optimal value the **minimum** p^* , and the point where the minimum is obtained the **minimizer** θ^* . Thus $p^* = f_0(\theta^*)$.

Lagrangian

Another way of solving it

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

Definition (Lagrangian)

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

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

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

Lagrange dual function

Definition (Lagrange dual function)

The **Lagrange dual function** $g : \mathbb{R}^M \rightarrow \mathbb{R}$ maps α to the minimum of the Lagrangian over θ (possibly $-\infty$ for some values of α),

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



It is concave in α since it is the point-wise minimum of a family of affine functions of α .

θ is completely unconstrained

Interpretation of the Lagrangian

θ^* shall be feasible

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

$$P^* = f_0(\theta^*)$$

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

$$g(\alpha)$$

as long as it
satisfies $f_i(\theta) \leq 0$

* For every choice of α , the corresponding unconstrained
 $g(\alpha) = \min_{\theta \in \mathbb{R}^d} L(\theta, \alpha)$ is a lower bound on the optimal value of the
constrained problem:

$$\min_{\substack{\theta \in \mathbb{R}^d \\ f_i(\theta) \leq 0}} f_0(\theta) = f_0(\theta^*) \geq f_0(\theta^*) + \sum_{i=1}^M \underbrace{\alpha_i}_{\geq 0} \underbrace{f_i(\theta^*)}_{\leq 0} = L(\theta^*, \alpha) \geq \underbrace{\min_{\theta \in \mathbb{R}^d} L(\theta, \alpha)}_{g(\alpha)}$$

Hence, $\forall \alpha \quad f_0(\theta^*) \geq g(\alpha)$.

≥ 0

allowing θ to vary freely;
 $g(\alpha)$ will only get smaller

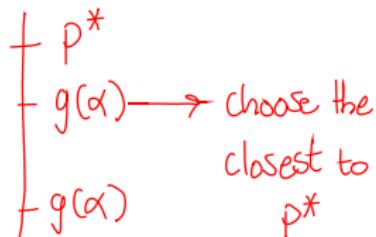
Lagrange dual problem

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

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



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

Note: since we are maximizing g we are not interested in dual multipliers α such that $g(\alpha) = -\infty$, so the condition $g(\alpha) \neq -\infty$ is usually added as an additional constraint to the dual problem → we call α **feasible** if and only if all $\alpha_i \geq 0$ and $L(\theta, \alpha)$ is bounded from below for $\theta \in \mathbb{R}^d$.

Duality

Weak duality (always)

Since for all $\alpha \geq 0$ it holds that $g(\alpha) \leq p^*$ we have **weak duality**,

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

The difference $p^* - d^* \geq 0$ between the solution of the original and the dual problem is called the **duality gap**.

Strong duality (under certain conditions)

Under certain conditions we have

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

* i.e. the **maximum to the Lagrange dual problem** is the **minimum of the original (primal) constrained optimization problem** (i.e. $f_0(\theta^*) = g(\alpha^*)$).

i.e. Solving either is the same

SVM's Primal problem ✘

$$\begin{array}{ll} \min_{\omega, b} & \omega^T \omega \\ \text{s.t.} & y_i (\omega^T \mathbf{x}_i + b) - 1 \geq 0 \end{array}$$

We apply a recipe for solving the constrained optimization problem.

1. Calculate the Lagrangian

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

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

Assumption for
Lagrange: $f_i \leq 0$

SVM : $f_i \geq 0$

So, we multiply -1

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

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

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

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

SVM's Dual problem

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

Thus we have reformulated our original problem as

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

$$\begin{aligned} \text{subject to} \quad & \sum_{i=1}^N \alpha_i y_i = 0 \\ & \alpha_i \geq 0, \quad \text{for } i = 1, \dots, N. \end{aligned}$$

3. Solve this problem.

Solving the dual problem

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

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

\min , Sym PSD
 \max , sym NSD

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

equivalent

★ 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 LIBSVM³ 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 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 w are a linear combination of the training samples,

$$w = \sum_{i=1}^N \alpha_i^* y_i x_i$$

* Check extra notes

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

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

$$f_i(\theta^*) = 0$$

$$y_i (w^T x_i + b) - 1 = 0$$

Solving this for b yields the bias

$$b = y_i - w^T x_i$$



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

Support vectors

those that support the hyperplane

From complimentary slackness

$$\alpha_i^* f_i(\theta^*) = 0.$$

only non-ve α_i :

contribute to w

$$w = \sum_{i=1}^N \alpha_i y_i x_i = \sum_{\substack{S \\ \{\alpha_i | \alpha_i > 0\}}} \alpha_i y_i x_i$$

we have

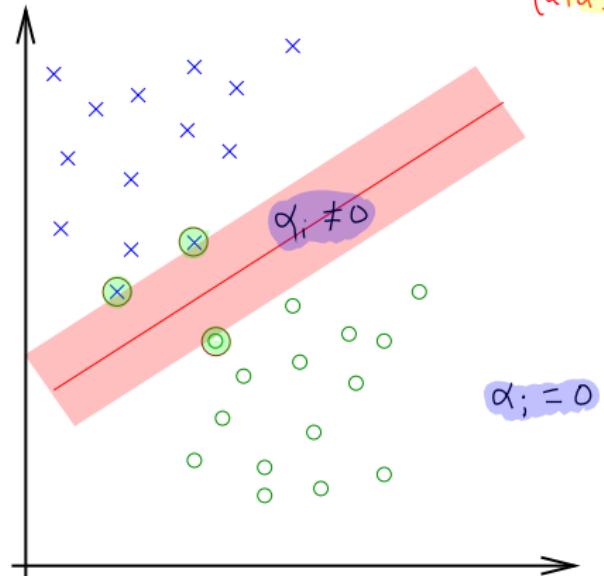
$$\alpha_i [y_i (w^T x_i + b) - 1] = 0 \quad \text{for all } i.$$

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

$$\alpha_i > 0$$

$$y_i (w^T x_i + b) = 1.$$

A training sample x_i with $\alpha_i \neq 0$ is called a **support vector**.



Classifying

The class of x is given by

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

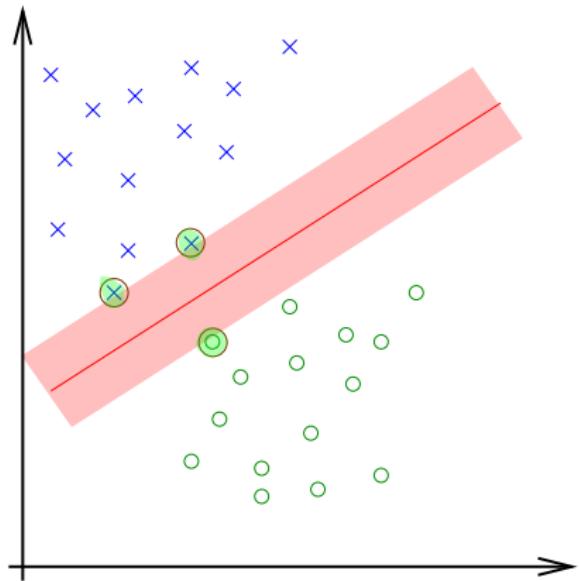
Substituting

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

gives

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

* 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 2

Soft Margin Support Vector Machines

Dealing with noisy data

$$\begin{array}{l} \text{max } \text{margin} \\ \text{min } \|\omega\|_2 \perp \text{hyperplane} \end{array}$$

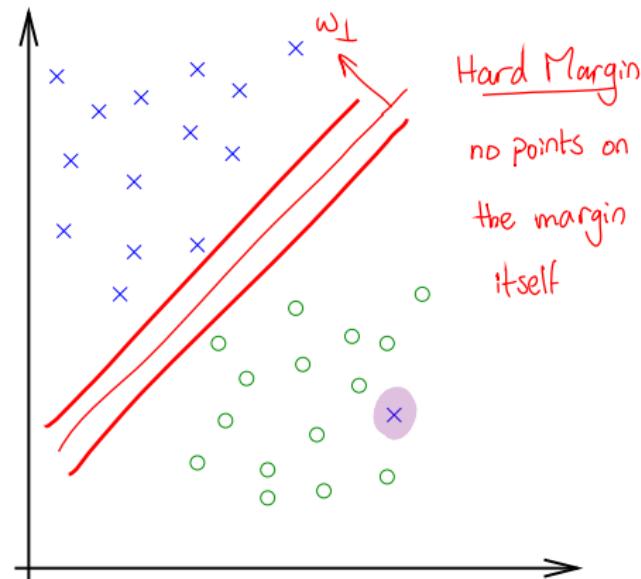
What if the data is not linearly separable due to some noise?

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

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

How to make our model more robust?

Hard Constraint holds true for all data points



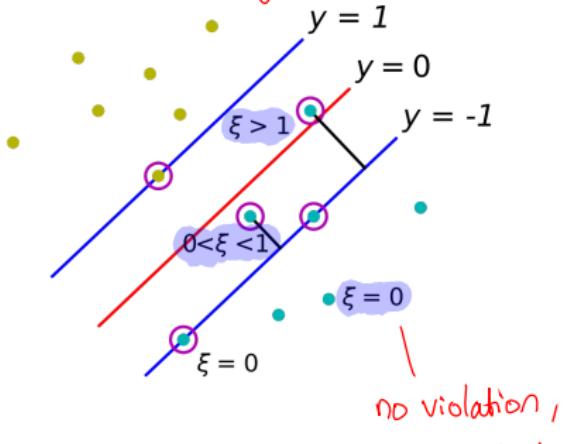
Slack variables

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 $\|w\|$.

the farther away the misclassified point,

the larger ξ_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 x_i that gives the distance of how far the margin is violated by this training sample in units of $\|\mathbf{w}\|$.

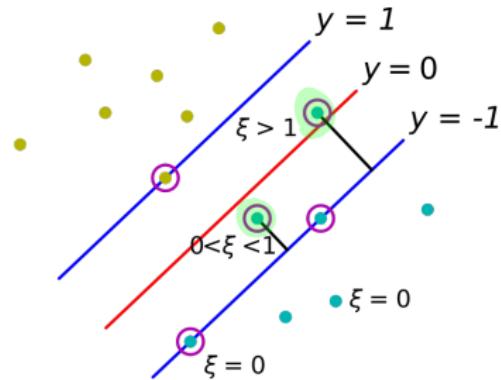
Hence our **relaxed constraints** are

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

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

Again, they can be condensed into

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



let $\underbrace{\mathbf{w}^T \mathbf{x}_i + b}_{-0.3} \leq -1$ Hence, $\xi_i = +0.7$
not more!

Only add as much as necessary
for the equality to hold

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 $\|\mathbf{w}\|$.

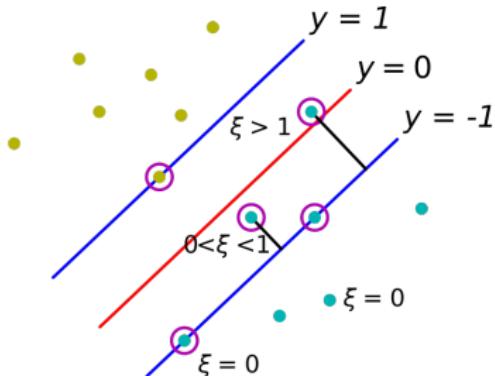
Hence our relaxed constraints are

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

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

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,

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

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

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

Optimization problem with slack variables

Hard Margin
(slide 12)

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

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

optimizing
over
 $\begin{pmatrix} w \\ b \end{pmatrix}, \xi$

$$\text{minimize } f_0(w, b, \xi) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$
$$\text{subject to } y_i(w^T x_i + b) - 1 + \xi_i \geq 0 \quad i = 1, \dots, N,$$

(2 Constraints) $\xi_i \geq 0 \quad i = 1, \dots, N.$

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

penalize those farther

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

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

1 Calculate the Lagrangian

$$\underbrace{f_0(\theta)}$$

$$\alpha = [\alpha, \mu]$$

$$L(w, b, \xi, \alpha, \mu) = \frac{1}{2} w^T w + C \sum_{i=1}^N \xi_i$$

$$\Theta = [w, b, \xi]$$

$$- \sum_{i=1}^N \alpha_i \underbrace{[y_i(w^T x_i + b) - 1 + \xi_i]}_{f_i(\theta)} - \sum_{i=1}^N \mu_i \xi_i.$$

2 constraints

Dual

2 Minimize $L(w, b, \xi, \alpha, \mu)$ w.r.t. w , b and ξ . Unconstrained Θ

$$\nabla_w L(w, b, \xi, \alpha, \mu) = w - \sum_{i=1}^N \alpha_i y_i x_i \stackrel{!}{=} 0; \quad \frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i y_i \stackrel{!}{=} 0,$$

$$\frac{\partial L}{\partial \xi_i} = C - \alpha_i - \mu_i \stackrel{!}{=} 0 \quad \text{for } i = 1, \dots, N$$

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

Condensing constraints into

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

Dual problem with slack variables

This leads to the **dual problem**:

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

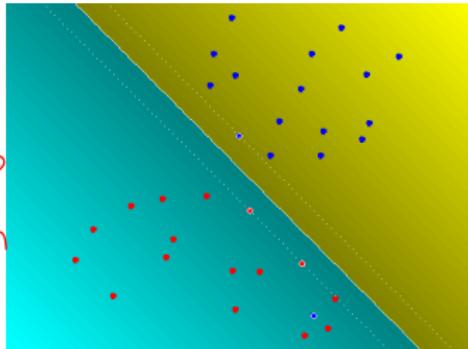

This is nearly the same dual problem as for the case without slack variables.

(slide 20)

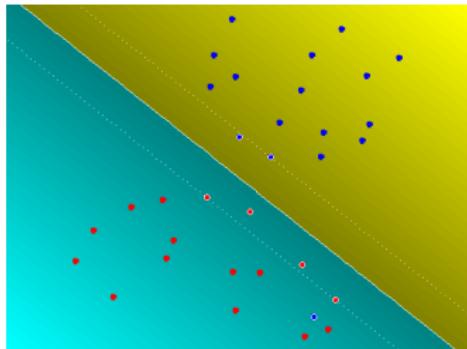
- * Only the constraint $\alpha_i \leq C$ is new. It ensures that α_i is bounded and cannot go to infinity.

Influence of the penalty C

equivalent to
hard margin

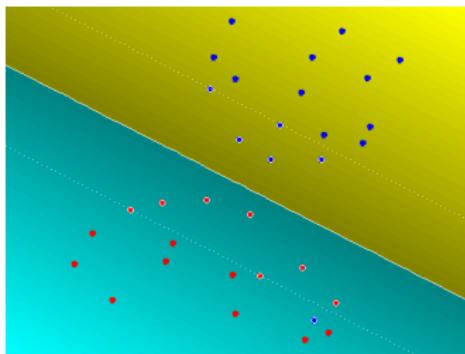


$C = 100$



$C = 10$

margin gets smaller
as well as the
points within



$C = 1$

allowing many points to
be within the margin ($\epsilon_i > 0$)

Hinge loss formulation

Rewriting the primal

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

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

subject to $y_i(\mathbf{w}^T \mathbf{x}_i + b) - 1 + \xi_i \geq 0, \quad i = 1, \dots, N,$

$\xi_i \geq 0, \quad i = 1, \dots, N.$

Observing values
at the optimal soln

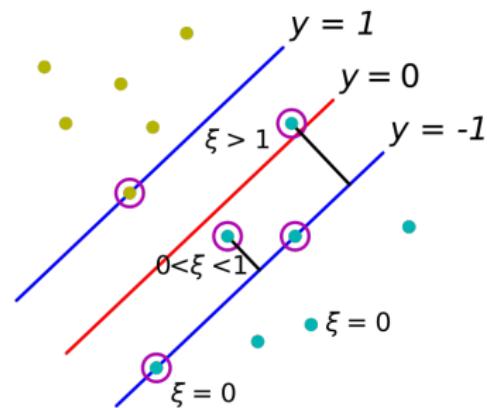
Clearly, for the **optimal solution** the
slack variables are

margin violated
i.e. constraint unsatisfied

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

Correctly classified

since we are minimizing over ξ .



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

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

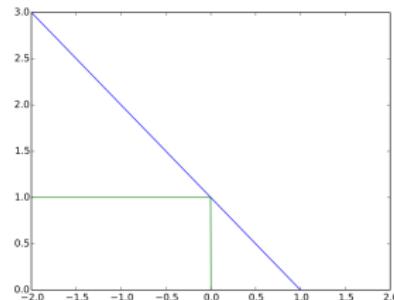


The hinge loss function

$E_{\text{hinge}}(z) = \max\{0, 1 - z\}$ 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 3

Kernels

A general concept

Feature space

So far we can only construct linear classifiers.

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

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

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

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

check

extra

notes

!!

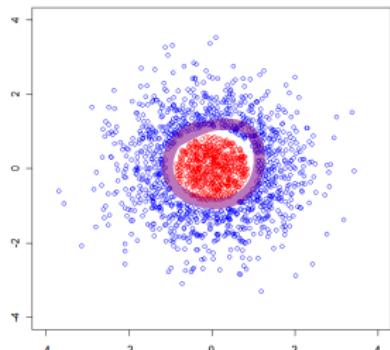


fig #1

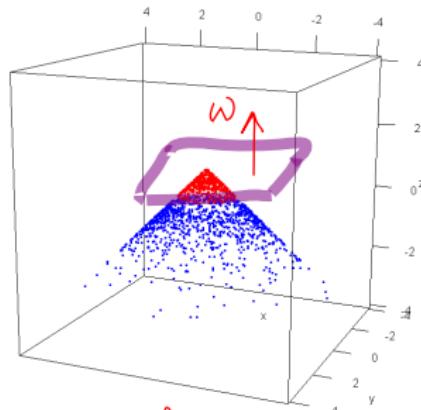


fig #2

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

Kernelized SVM

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

* Applying such transformation, we will get a non-linear DB
in the original input space

Kernel trick directly applying Kernel fn instead of the basis

$K: 2\text{vec} \rightarrow \text{scalar}$

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

$\phi: \text{vec} \rightarrow \text{vec}$

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

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)

* Think of $K(\cdot)$ as measure of similarity bet any 2 objects.
i.e high $K \Rightarrow$ high similarity

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?

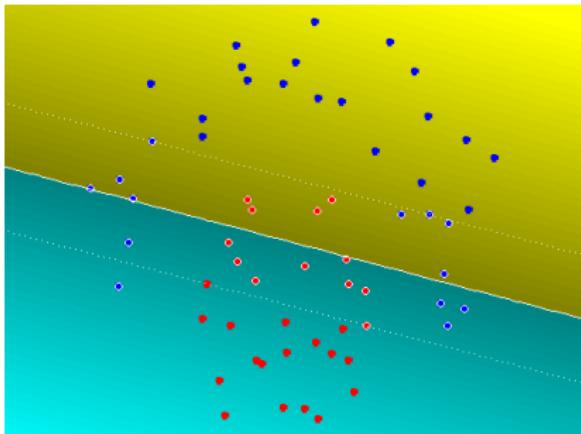
- Some kernels are equivalent to using infinite-dimensional basis functions. While computing these feature transformations would be impossible, directly evaluating the kernels is often easy.
- Kernels can be used to encode similarity between arbitrary non-numerical data, from strings to graphs.

For example, we could define

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

Check extra notes !!

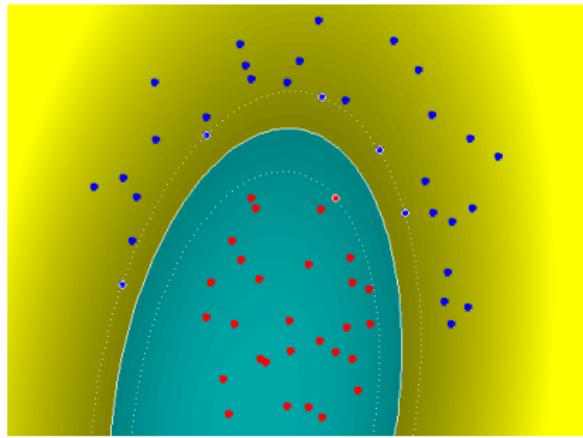
SVM using kernels example



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

$$K(a,b) = \phi(a)^T \phi(b)$$

$$\phi(a) = a$$



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

$$K(a,b)$$

What makes a valid kernel?

$K(a, b) = \phi(a)^T \phi(b)$ shall hold true
for 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 symmetric, positive semidefinite kernel matrix K for any input data X .

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

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

↙ 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.

Using any fn (not necessarily valid K) might work in practice

Convexity is tied to
Sym PSD

Kernel preserving operations

Let $k_1 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $k_2 : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be kernels, with $\mathcal{X} \subseteq \mathbb{R}^N$.
Then the following functions k are kernels as well:

- $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) + k_2(\mathbf{x}_1, \mathbf{x}_2)$
- $k(\mathbf{x}_1, \mathbf{x}_2) = c \cdot k_1(\mathbf{x}_1, \mathbf{x}_2)$, with $c > 0$
- $k(\mathbf{x}_1, \mathbf{x}_2) = k_1(\mathbf{x}_1, \mathbf{x}_2) \cdot k_2(\mathbf{x}_1, \mathbf{x}_2)$
- $k(\mathbf{x}_1, \mathbf{x}_2) = k_3(\phi(\mathbf{x}_1), \phi(\mathbf{x}_2))$, with the kernel k_3 on $\mathcal{X}' \subseteq \mathbb{R}^M$ and $\phi : \mathcal{X} \rightarrow \mathcal{X}'$
- $k(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1 \mathbf{A} \mathbf{x}_2$, with $\mathbf{A} \in \mathbb{R}^N \times \mathbb{R}^N$ symmetric and positive semidefinite

Examples of kernels

adding bias,

K is still valid

- Polynomial:

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

* **Gaussian kernel:** Corresponds to dot product in infinite-dim feature space

Sim \propto dist

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

more control over neighborhood

- Sigmoid:

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

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

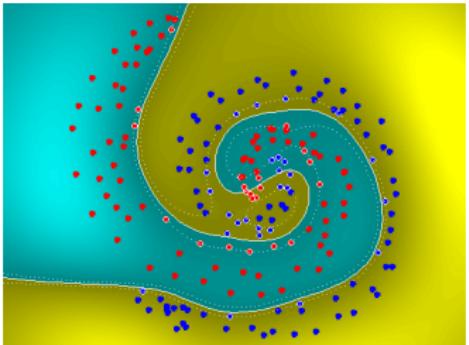
Some kernels introduce additional hyperparameters, that affect the behavior of the algorithm.



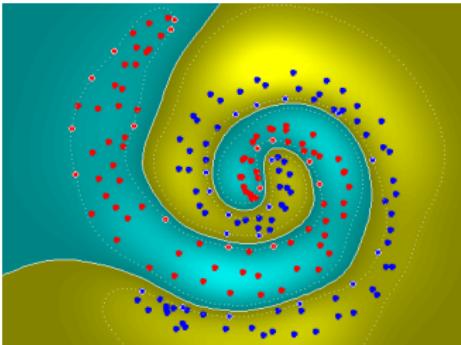
Note, that the sigmoid kernel is different from the sigmoid function from *Linear Classification*.

Gaussian kernel (C=1000)

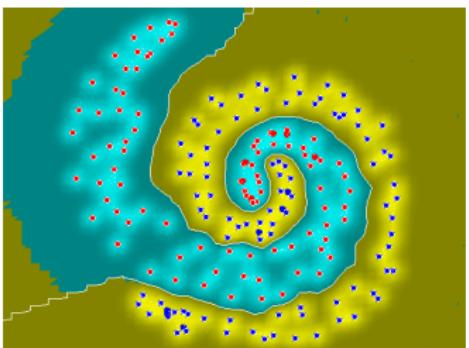
check extra notes !!



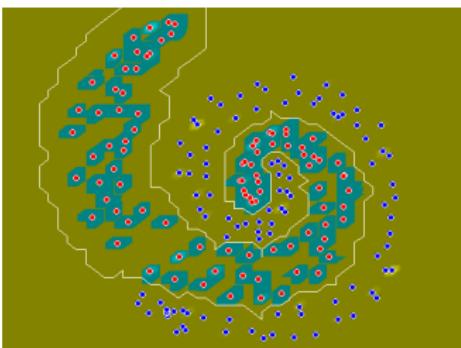
$$\sigma = 0.5$$



$$\sigma = 0.25$$



$$\sigma = 0.05$$



$$\sigma = 0.005$$

Classifying a new point with kernelized SVM

behaving as nearest neighbor classifier

We denote the set of support vectors as \mathcal{S} (points \mathbf{x}_i for which holds $0 < \alpha_i \leq C$). Note: If $0 < \alpha_i < C$ then $\xi_i = 0$; if $\alpha_i = C$ then $\xi_i > 0$.

From the complementary slackness condition, points $\mathbf{x}_i \in \mathcal{S}$ with $\xi_i = 0$ must satisfy

$$\sum_i \alpha_i f_i(\theta) = 0$$

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

Only those \mathbf{x}_i affect DB

Like for the regular SVM, the bias can be recovered as

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

$$w^* = \sum_i \alpha_i y_i \phi(\mathbf{x}_i)$$

as linear combination
of data points

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

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

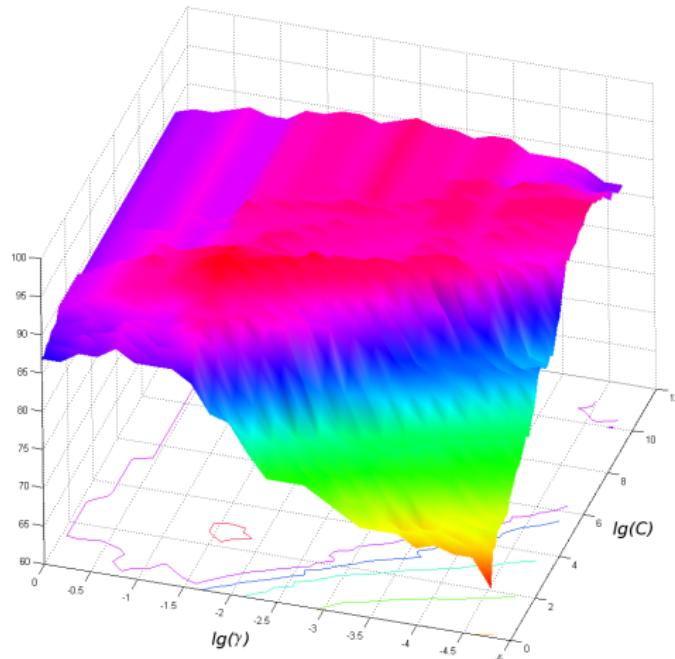
/ \

SV new point

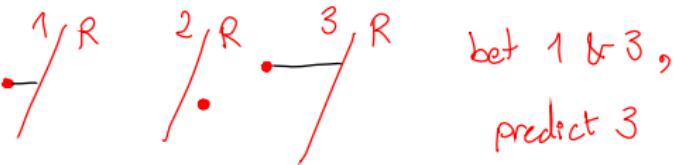
How to choose the hyperparameters?

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.

Choosing those
achieving the
highest acc



Multiple classes



The standard SVM model cannot handle multiclass data.

Two approaches to address this issue are:

scales linearly
with # classes

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.

might not scale

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. Since the parameters are unconstrained
- * 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.

performing dot product in potentially infinite-dim feature space

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