

Machine Learning

Support Vector Machines

Fabio Vandin

December 2nd, 2022

Duality

We now present (Hard-)SVM in a different way which is very useful for *kernels*.

We want to solve

data is linearly separable

$$\mathbf{w}_0 = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } \forall i : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$$

Duality

We now present (Hard-)SVM in a different way which is very useful for *kernels*.

We want to solve

$$\mathbf{w}_0 = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } \forall i : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$$

One can prove (details in the book!) that \mathbf{w} that minimizes the function above is equivalent to find α that solves the *dual problem*:

$$\mathcal{S} = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\} \quad \max_{\alpha \in \mathbb{R}^m : \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \right)$$

Duality

We now present (Hard-)SVM in a different way which is very useful for *kernels*.

We want to solve

$$\mathbf{w}_0 = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } \forall i : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$$

One can prove (details in the book!) that \mathbf{w} that minimizes the function above is equivalent to find α that solves the *dual problem*:

$$\max_{\alpha \in \mathbb{R}^m : \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_j, \mathbf{x}_i \rangle \right)$$

Note:

- solution is the vector α which defines the support vectors = $\{\mathbf{x}_i : \alpha_i \neq 0\}$
- \mathbf{w}_0 can be derived from α (see previous slides!)

Duality

We now present (Hard-)SVM in a different way which is very useful for *kernels*.

We want to solve

$$\mathbf{w}_0 = \min_{\mathbf{w}} \frac{1}{2} \|\mathbf{w}\|^2 \text{ subject to } \forall i : y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$$

One can prove (details in the book!) that \mathbf{w} that minimizes the function above is equivalent to find α that solves the *dual problem*:

$$\max_{\alpha \in \mathbb{R}^m : \alpha \geq 0} \left(\sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i=1}^m \sum_{j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_j, \mathbf{x}_i \rangle \right)$$

Note:

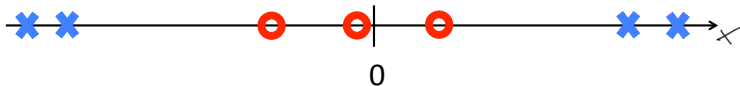
- solution is the vector α which defines the support vectors = $\{\mathbf{x}_i : \alpha_i \neq 0\}$
- \mathbf{w}_0 can be derived from α (see previous slides!)
- dual problem requires only to compute inner products $\langle \mathbf{x}_j, \mathbf{x}_i \rangle$, does not need to consider \mathbf{x}_i by itself

SVM is a powerful algorithm, but still limited to linear models...
and linear models cannot always be used (directly)!

$$\vec{x} = [x]$$

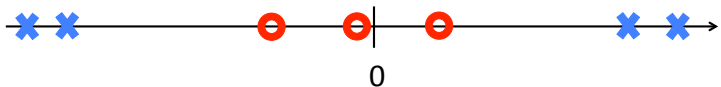
$$\mathcal{X} = \mathbb{R}$$

Example



SVM is a powerful algorithm, but still limited to linear models...
and linear models cannot always be used (directly)!

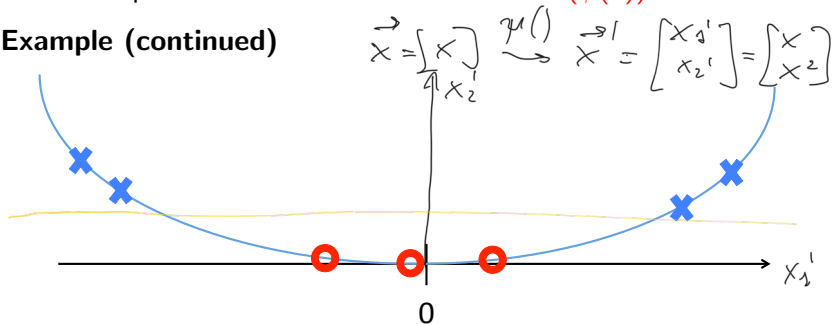
Example



We can:

- apply a nonlinear transformation $\psi()$ to each point in training set S first: $S' = ((\psi(\mathbf{x}_1), y_1), \dots, (\psi(\mathbf{x}_m), y_m))$;
- learn a linear predictor \hat{h} in the transformed space using S' ;
- make prediction for a new instance \mathbf{x} as $\hat{h}(\psi(\mathbf{x}))$

Example (continued)



Kernel Trick for SVM

What if we want to apply a nonlinear transformation before using SVM?

Let $\psi()$ be the nonlinear transformation

i) Given training set S : obtain training S'

$$S = \{(\vec{x}_1, y_1), \dots, (\vec{x}_m, y_m)\}$$

↓

$$S' = \{(\psi(\vec{x}_1), y_1), \dots, (\psi(\vec{x}_m), y_m)\}$$

ii) Learn a model with SVM using S' , let h_{svm} be the model we learn

iii) Given $\vec{x} \in \mathcal{X}$, the prediction is $h_{\text{svm}}(\psi(\vec{x}))$

Kernel Trick for SVM

What if we want to apply a nonlinear transformation before using SVM?

Let $\psi()$ be the nonlinear transformation

Considering the dual formulation \Rightarrow we only need to be able to compute $\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle$ for some \mathbf{x}, \mathbf{x}' .

Definition

A *kernel function* is a function of the type:

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

where $\psi(\mathbf{x})$ is a transformation of \mathbf{x} .

Kernel Trick for SVM

What if we want to apply a nonlinear transformation before using SVM?

Let $\psi()$ be the nonlinear transformation

Considering the dual formulation \Rightarrow we only need to be able to compute $\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle$ for some \mathbf{x}, \mathbf{x}' .

Definition

A *kernel function* is a function of the type:

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

where $\psi(\mathbf{x})$ is a transformation of \mathbf{x} .

Intuition: we can think of K_{ψ} as specifying *similarity* between instances and of ψ as mapping the domain set \mathcal{X} into a space where these similarities are realized as dot products.

Kernel Trick for SVM(2)

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

It seems that to compute $K_{\psi}(\mathbf{x}, \mathbf{x}')$ requires to be able to compute $\psi(\mathbf{x})$...

Not always... sometimes we can compute $K_{\psi}(\mathbf{x}, \mathbf{x}')$ without computing $\psi(\mathbf{x})$!

Kernel: Example

Consider $\mathbf{x} \in \mathbb{R}^d$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix} \in \mathbb{R}^d$$

$$\psi(\mathbf{x}) = (1, \overbrace{x_1, x_2, \dots, x_d}^d, \overbrace{x_1 x_1, x_1 x_2, x_1 x_3, \dots, x_d x_d}^{d^2})^T$$

What is the dimension of $\psi(\vec{x})$? $1 + d + d^2$

Kernel: Example

Consider $\mathbf{x} \in \mathbb{R}^d$

$$\psi(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, x_1x_3, \dots, x_dx_d)^T$$

The dimension of $\psi(\mathbf{x})$ is $1 + d + d^2$.

$$K_{\psi}(\vec{x}, \vec{x}') = \left\langle \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \vdots \\ x_d \\ x_1x_1 \\ \vdots \\ x_dx_d \end{bmatrix}, \begin{bmatrix} 1 \\ x'_1 \\ x'_2 \\ \vdots \\ x'_d \\ x'_1x'_1 \\ \vdots \\ x'_dx'_d \end{bmatrix} \right\rangle = 1 + (x_1x'_1 + x_2x'_2 + \dots + x_dx'_d) + x_1x_1x'_1x'_1 + x_1x_2x'_1x'_2 + \dots + x_dx_dx'_dx'_d$$

$$\left(\sum_{i=1}^d x_i x'_i \right) \left(\sum_{j=1}^d x_j x'_j \right)$$

$$\langle \vec{x}, \vec{x}' \rangle^2 = \langle \vec{x}, \vec{x}' \rangle \cdot \langle \vec{x}, \vec{x}' \rangle$$

$$= 1 + \underbrace{\sum_{i=1}^d x_i x'_i}_{\langle \vec{x}, \vec{x}' \rangle} + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j$$

Kernel: Example

Consider $\mathbf{x} \in \mathbb{R}^d$

$$\psi(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, x_1x_3, \dots, x_dx_d)^T$$

The dimension of $\psi(\mathbf{x})$ is $1 + d + d^2$.

$$\langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle = 1 + \sum_{i=1}^d x_i x'_i + \sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j$$

Note that

$$\sum_{i=1}^d \sum_{j=1}^d x_i x_j x'_i x'_j = \left(\sum_{i=1}^d x_i x'_i \right) \left(\sum_{j=1}^d x_j x'_j \right) = (\langle \mathbf{x}, \mathbf{x}' \rangle)^2$$

therefore

$$K_\psi(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle = 1 + \langle \mathbf{x}, \mathbf{x}' \rangle + (\langle \mathbf{x}, \mathbf{x}' \rangle)^2$$

We have:

$$\psi(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, x_1x_3, \dots, x_dx_d)^T$$

$$K_\psi(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle = 1 + \langle \mathbf{x}, \mathbf{x}' \rangle + (\langle \mathbf{x}, \mathbf{x}' \rangle)^2$$

Computation of $\psi(\vec{x})$: how many operations? $\Theta(d^2)$
(starting from \vec{x})

Computation of $K_\psi(\vec{x}, \vec{x}')$: how many operations?

If I first compute $\psi(\vec{x})$ and $\psi(\vec{x}')$, and then their \langle, \rangle : $\Theta(d^2)$

If I use $K_\psi(\vec{x}, \vec{x}') = 1 + \underbrace{\langle \vec{x}, \vec{x}' \rangle}_{\Theta(d)} + \underbrace{(\langle \vec{x}, \vec{x}' \rangle)^2}_{\Theta(1)} : \Theta(d)$

We have:

$$\psi(\mathbf{x}) = (1, x_1, x_2, \dots, x_d, x_1x_1, x_1x_2, x_1x_3, \dots, x_dx_d)^T$$

$$K_\psi(\mathbf{x}, \mathbf{x}') = \langle \psi(\mathbf{x}), \psi(\mathbf{x}') \rangle = 1 + \langle \mathbf{x}, \mathbf{x}' \rangle + (\langle \mathbf{x}, \mathbf{x}' \rangle)^2$$

Observation

Computing $\psi(\mathbf{x})$ requires $\Theta(d^2)$ time; computing $K_\psi(\mathbf{x}, \mathbf{x}')$ from the last formula requires $\Theta(d)$ time

When $K_\psi(\mathbf{x}, \mathbf{x}')$ is efficiently computable, we don't need to explicitly compute $\psi(\mathbf{x})$

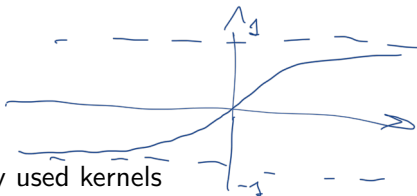
\Rightarrow *kernel trick*

Some Kernels

The following are the most commonly used kernels

- linear kernel: $\psi(\mathbf{x}) = \mathbf{x}$ (standard SVM)

Some Kernels



The following are the most commonly used kernels

- linear kernel: $\psi(\mathbf{x}) = \mathbf{x}$
- sigmoid: $K(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \langle \mathbf{x}, \mathbf{x}' \rangle + \zeta)$ (for $\gamma, \zeta > 0$)

hyperbolic tangent

Some Kernels

The following are the most commonly used kernels

- linear kernel: $\psi(\mathbf{x}) = \mathbf{x}$
- sigmoid: $K(\mathbf{x}, \mathbf{x}') = \tanh(\gamma \langle \mathbf{x}, \mathbf{x}' \rangle + \zeta)$ (for $\gamma, \zeta > 0$)
- degree- Q polynomial kernel
- Gaussian-radial basis function (RBF) kernel

Degree- Q polynomial kernel

Definition

For given constants $\gamma > 0, \zeta > 0$ and for $Q \in \mathbb{N}$, the *degree- Q polynomial kernel* is

$$K(\mathbf{x}, \mathbf{x}') = (\zeta + \gamma \langle \mathbf{x}, \mathbf{x}' \rangle)^Q$$

time? ~~$\binom{Q}{d}$~~
 ~~$\binom{Q}{Q+d}$~~
 $\binom{Q}{d + \frac{Q}{2}}$

Example

For $Q = 2$:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_d \end{bmatrix}$$

$$\psi(\mathbf{x}) = [\zeta, \sqrt{2\zeta\gamma}x_1, \sqrt{2\zeta\gamma}x_2, \dots, \sqrt{2\zeta\gamma}x_d, \gamma x_1 x_1, \gamma x_1 x_2, \dots, \gamma x_d x_d]^T \in \mathbb{R}^{1+d+d^2}$$

In general: $\psi(\vec{x}) \in \mathbb{R}^{\binom{Q+d}{d}}(d^Q)$

Gaussian-RBF Kernel

Definition

For a given constant $\gamma > 0$ the *Gaussian-RBF kernel* is

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

What is $\psi(\mathbf{x})$? Assume $\gamma = 1$ and $\mathbf{x} = x \in \mathbb{R}$ for simplicity, then

$$\begin{aligned} K(x, x') &= e^{-\|x - x'\|^2} \\ &= e^{-x^2} e^{2xx'} e^{-(x')^2} \\ &= e^{-x^2} \left(\sum_{k=0}^{+\infty} \frac{2^k (x)^k (x')^k}{k!} \right) e^{-(x')^2} \end{aligned}$$

Taylor expansion of $e^{2xx'}$

$$= \sum_{k=0}^{+\infty} \frac{(2xx')^k}{k!}$$

Gaussian-RBF Kernel

Definition

For a given constant $\gamma > 0$ the *Gaussian-RBF kernel* is

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

What is $\psi(\mathbf{x})$? Assume $\gamma = 1$ and $\mathbf{x} = x \in \mathbb{R}$ for simplicity, then

$$\begin{aligned} K(x, x') &= e^{-\|x - x'\|^2} \\ &= e^{-x^2} e^{2xx'} e^{-(x')^2} \\ &= e^{-x^2} \left(\sum_{k=0}^{+\infty} \frac{2^k (x)^k (x')^k}{k!} \right) e^{-(x')^2} = \langle \psi(x), \psi(x') \rangle \end{aligned}$$

$$\Rightarrow \psi(x) = e^{-x^2} \left(1, \sqrt{\frac{2}{1!}} x, \sqrt{\frac{2^2}{2!}} x^2, \sqrt{\frac{2^3}{3!}} x^3, \dots \right)^T$$

Gaussian-RBF Kernel

Definition

For a given constant $\gamma > 0$ the *Gaussian-RBF kernel* is

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

What is $\psi(\mathbf{x})$? Assume $\gamma = 1$ and $\mathbf{x} = x \in \mathbb{R}$ for simplicity, then

$$\begin{aligned} K(x, x') &= e^{-\|x - x'\|^2} \\ &= e^{-x^2} e^{2xx'} e^{-(x')^2} \\ &= e^{-x^2} \left(\sum_{k=0}^{+\infty} \frac{2^k (x)^k (x')^k}{k!} \right) e^{-(x')^2} \end{aligned}$$

$$\Rightarrow \psi(x) = e^{-x^2} \left(1, \sqrt{\frac{2}{1!}}x, \sqrt{\frac{2^2}{2!}}x^2, \sqrt{\frac{2^3}{3!}}x^3, \dots \right)^T$$

$\Rightarrow \psi(x)$ has infinite number of dimensions!

Choice of Kernel

Notes

- polynomial kernel: usually used with $Q \leq 10$
- Gaussian-RBF kernel: usually $\gamma \in [0, 1]$
- many other choices are possible!

Mercer's condition

$K(\mathbf{x}, \mathbf{x}')$ is a valid kernel function if and only if the kernel matrix

$$K = \begin{bmatrix} K(\mathbf{x}_1, \mathbf{x}_1) & K(\mathbf{x}_1, \mathbf{x}_2) \dots & K(\mathbf{x}_1, \mathbf{x}_m) \\ K(\mathbf{x}_2, \mathbf{x}_1) & K(\mathbf{x}_2, \mathbf{x}_2) \dots & K(\mathbf{x}_2, \mathbf{x}_m) \\ \vdots & \vdots & \vdots \\ K(\mathbf{x}_m, \mathbf{x}_1) & K(\mathbf{x}_m, \mathbf{x}_2) \dots & K(\mathbf{x}_m, \mathbf{x}_m) \end{bmatrix}$$

is always symmetric positive semi-definite for any given $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$.

Support Vector Machines for Regression

$$y_i \in \mathbb{R}$$

SVMs can be also used for regression. The function to be minimized will be

$$\ell_2\text{-regulization} \leftarrow \left(\frac{\lambda}{2} \|\mathbf{w}\|^2 + \sum_{i=1}^m V_{\epsilon}(y_i - \langle \mathbf{x}_i, \mathbf{w} \rangle - b) \right)$$

Annotations:

- $\frac{\lambda}{2} \|\mathbf{w}\|^2$ is circled and labeled $\ell_2\text{-regulization}$.
- V_{ϵ} is labeled "loss function".
- y_i is labeled "observed value".
- $\langle \mathbf{x}_i, \mathbf{w} \rangle + b$ is labeled "predicted value (by model \vec{w}, b)".

where

$$V_{\epsilon}(r) = \begin{cases} 0 & \text{if } |r| < \epsilon \\ |r| - \epsilon & \text{otherwise} \end{cases}$$

One can prove that the solution has the form:

$$\mathbf{w} = \sum_{i=1}^m (\alpha_i^* - \alpha_i) \mathbf{x}_i$$

\vec{x}_i in the training set ($\vec{x}_i \in S$)

and that the final model produced in output is

the prediction for $\vec{x} \in \mathbb{R}^d$

$$h(\mathbf{x}) = \sum_{i=1}^m (\alpha_i^* - \alpha_i) \langle \mathbf{x}_i, \mathbf{x} \rangle + b$$

training set points and \vec{x}

where $\alpha_i^*, \alpha_i \geq 0$ and are the solution to a suitable QP.

Definition

Support vector: \mathbf{x}_i such that $\alpha_i^* - \alpha_i \neq 0$

One can define kernels, similarly to SVM for classification.

Exercise 4

Assuming we have the following dataset ($x_i \in \mathbb{R}^2$) and by solving the SVM for classification we get the corresponding optimal dual variables:

i	x_i^T	y_i	α_i^*
1	[0.2 -1.4]	-1	0
2	[-2.1 1.7]	1	0
3	[0.9 1]	1	0.5
4	[-1 -3.1]	-1	0
5	[-0.2 -1]	-1	0.25
6	[-0.2 1.3]	1	0
7	[2.0 -1]	-1	0.25
8	[0.5 2.1]	1	0

Answer to the following:

- (A) What are the support vectors?
- (B) Draw a schematic picture reporting the data points (approximately) and the optimal separating hyperplane, and mark the support vectors. Would it be possible, by moving only two data points, to obtain the SAME separating hyperplane with only 2 support vectors? If so, draw the modified configuration (approximately).

Bibliography [UML]

SVM: Chapter 15

- no sections 15.1.2, 15.2.1, 15.2.2,15.2.3,

Kernels: Chapter 16

- no section 16.3