

# Kernel Methods, SVM

Milan Straka

 November 07, 2022



Charles University in Prague  
Faculty of Mathematics and Physics  
Institute of Formal and Applied Linguistics



unless otherwise stated

# Kernel Linear Regression

Consider linear regression with linear, quadratic and cubic features (for simplicity we consider  $x_i$ ,  $x_i x_j$  and  $x_i x_j x_k$  for any indices) computed by feature mapping  $\varphi : \mathbb{R}^D \rightarrow \mathbb{R}^{1+D+D^2+D^3}$ :

$$\varphi(\mathbf{x}) = \begin{bmatrix} 1 \\ x_1 \\ x_2 \\ \dots \\ x_1^2 \\ x_1 x_2 \\ \dots \\ x_2 x_1 \\ \dots \\ x_1^3 \\ x_1^2 x_2 \\ \dots \end{bmatrix}.$$

The SGD update of a linear regression using a minibatch of examples  $\mathbb{B}$  is then

$$\mathbf{w} \leftarrow \mathbf{w} - \frac{\alpha}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} (\varphi(\mathbf{x}_i)^T \mathbf{w} - t_i) \varphi(\mathbf{x}_i).$$

# Kernel Linear Regression

When the dimensionality of the input is  $D$ , one step of SGD takes  $\mathcal{O}(D^3)$  per example.

Surprisingly, we can do better under some circumstances. We start by noting that we can write the parameters  $\mathbf{w}$  as a linear combination of the  $N$  input feature vectors  $\varphi(\mathbf{x}_i)$ .

By induction, we can start with  $\mathbf{w} = \mathbf{0} = \sum_{i=1}^N 0 \cdot \varphi(\mathbf{x}_i)$ . Assuming  $\mathbf{w} = \sum_{i=1}^N \beta_i \cdot \varphi(\mathbf{x}_i)$ , after an SGD update we get

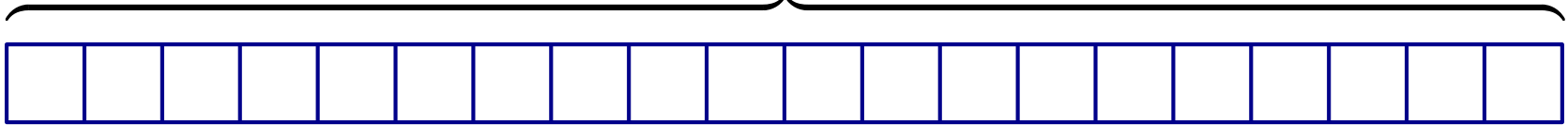
$$\begin{aligned} \mathbf{w} &\leftarrow \mathbf{w} - \frac{\alpha}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} (\varphi(\mathbf{x}_i)^T \mathbf{w} - t_i) \varphi(\mathbf{x}_i) \\ &\leftarrow \sum_{i=1}^N \left( \beta_i - [i \in \mathbb{B}] \cdot \frac{\alpha}{|\mathbb{B}|} (\varphi(\mathbf{x}_i)^T \mathbf{w} - t_i) \right) \varphi(\mathbf{x}_i). \end{aligned}$$

Every  $\beta_i$  for  $i \in \mathbb{B}$  changes to  $\beta_i - \frac{\alpha}{|\mathbb{B}|} (\varphi(\mathbf{x}_i)^T \mathbf{w} - t_i)$ , so after substituting for  $\mathbf{w}$  we get

$$\beta_i \leftarrow \beta_i - \frac{\alpha}{|\mathbb{B}|} \left( \left( \sum_{j=1}^N \beta_j \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j) \right) - t_i \right).$$

$$\mathcal{O}(D^3)$$

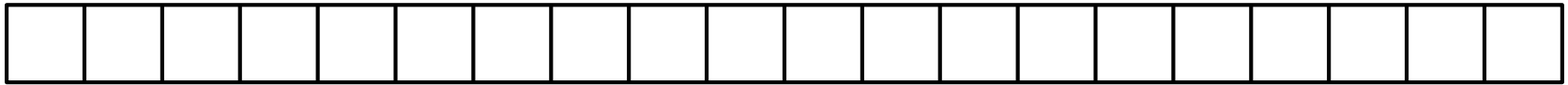
$w$



$\varphi(\mathbf{x}_1)$



$\varphi(\mathbf{x}_2)$



$\varphi(\mathbf{x}_3)$



$\vdots$



$\varphi(\mathbf{x}_N)$



$\beta$

$\beta_1$

$\beta_2$

$\beta_3$

$\vdots$

$\beta_N$



We can formulate an alternative linear regression algorithm (a so-called **dual formulation**):

**Input:** Dataset ( $\mathbf{X} = \{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_N\} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{t} \in \mathbb{R}^N$ ), learning rate  $\alpha \in \mathbb{R}^+$ .

- $\boldsymbol{\beta} \leftarrow \mathbf{0}$
- compute all values  $K_{i,j} = \varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$
- until convergence (or patience runs out), process a minibatch of examples  $\mathbb{B}$ :
  - simultaneously for all  $i \in \mathbb{B}$  (the  $\beta_j$  on the right side must not be modified during the batch update):
    - $\beta_i \leftarrow \beta_i - \frac{\alpha}{|\mathbb{B}|} \left( \left( \sum_{j=1}^N \beta_j K_{i,j} \right) - t_i \right)$

The predictions for input  $\mathbf{z}$  are then performed by computing

$$y(\mathbf{z}) = \varphi(\mathbf{z})^T \mathbf{w} = \sum_{i=1}^N \beta_i \varphi(\mathbf{z})^T \varphi(\mathbf{x}_i).$$

# Bias in Kernel Linear Regression

Until now we did not consider *bias*. Unlike the usual formulation, where we can “hide” it in the weights, we usually handle it manually in the dual formulation.

Specifically, if we want to include bias in kernel linear regression, we modify the predictions to

$$y(\mathbf{z}) = \varphi(\mathbf{z})^T \mathbf{w} + b = \left( \sum_{i=1}^N \beta_i \varphi(\mathbf{z})^T \varphi(\mathbf{x}_i) \right) + b$$

and update the bias  $b$  separately.

To update the bias, we use SGD, and the resulting update is

$$b \leftarrow b - \frac{\alpha}{|\mathbb{B}|} \sum_{i \in \mathbb{B}} \left( \left( \sum_{j=1}^N \beta_j K_{i,j} \right) + b - t_i \right).$$

# Kernel Trick

A single SGD update of a dual-formulation kernel linear regression takes  $\mathcal{O}(N)$  per example, if we pre-compute all the  $\mathcal{O}(N^2)$  dot products  $\varphi(\mathbf{x}_i)^T \varphi(\mathbf{x}_j)$ . Furthermore, inference requires evaluating  $\mathcal{O}(N)$  dot products  $\varphi(\mathbf{z})^T \varphi(\mathbf{x}_i)$ .

Therefore, we need to compute the dot product  $\varphi(\mathbf{x})^T \varphi(\mathbf{z})$  quickly.

Using the previously-defined  $\varphi$ , we get

$$\begin{aligned} \varphi(\mathbf{x})^T \varphi(\mathbf{z}) &= 1 + \sum_i \mathbf{x}_i z_i + \sum_{i,j} \mathbf{x}_i \mathbf{x}_j z_i z_j + \sum_{i,j,k} \mathbf{x}_i \mathbf{x}_j \mathbf{x}_k z_i z_j z_k \\ &= 1 + \sum_i \mathbf{x}_i z_i + \left( \sum_i \mathbf{x}_i z_i \right)^2 + \left( \sum_i \mathbf{x}_i z_i \right)^3 \\ &= 1 + \mathbf{x}^T \mathbf{z} + (\mathbf{x}^T \mathbf{z})^2 + (\mathbf{x}^T \mathbf{z})^3. \end{aligned}$$

Therefore, we can compute the dot product  $\varphi(\mathbf{x})^T \varphi(\mathbf{z})$  in  $\mathcal{O}(D)$  instead of  $\mathcal{O}(D^3)$ .

We define a **kernel** corresponding to a feature map  $\varphi$  as a function

$$K(\mathbf{x}, \mathbf{z}) \stackrel{\text{def}}{=} \varphi(\mathbf{x})^T \varphi(\mathbf{z}).$$

There exist quite a lot of kernels, but the most commonly used are the following:

- **Polynomial kernel of degree  $d$** , also called *homogenous polynomial kernel*,

$$K(\mathbf{x}, \mathbf{z}) = (\gamma \mathbf{x}^T \mathbf{z})^d,$$

corresponds to a feature map returning all combinations of exactly  $d$  input features.

Using  $(a_1 + \dots + a_k)^d = \sum_{n_i \geq 0, \sum n_i = d} \binom{d}{n_1, \dots, n_k} a_1^{n_1} \dots a_k^{n_k}$ , we can verify that

$$\varphi(\mathbf{x}) = \left( \sqrt{\gamma^d \binom{d}{n_1, \dots, n_D}} x_1^{n_1} \dots x_D^{n_D} \right)_{n_i \geq 0, \sum n_i = d}.$$

For example, for  $d = 2$ ,  $\varphi(x_1, x_2) = \gamma(x_1^2, \sqrt{2}x_1x_2, x_2^2)$ .



- **Polynomial kernel of degree at most  $d$** , also called *nonhomogenous polynomial kernel*,

$$K(\mathbf{x}, \mathbf{z}) = (\gamma \mathbf{x}^T \mathbf{z} + 1)^d,$$

corresponds to a feature map generating all combinations of up to  $d$  input features.

Given that  $(\gamma \mathbf{x}^T \mathbf{z} + 1)^d = \sum_i \binom{d}{i} (\gamma \mathbf{x}^T \mathbf{z})^i$ , it is not difficult to derive that

$$\varphi(\mathbf{x}) = \left( \sqrt{\gamma^{d-n_{D+1}} \binom{d}{n_1, \dots, n_{D+1}}} x_1^{n_1} \cdots x_D^{n_D} \right)_{n_i \geq 0, \sum_{i=1}^{D+1} n_i = d}.$$

For example, for  $d = 2$ ,  $\varphi(x_1, x_2) = (1, \sqrt{2\gamma}x_1, \sqrt{2\gamma}x_2, \gamma x_1^2, \sqrt{2\gamma}x_1x_2, \gamma x_2^2)$ .

- **Gaussian Radial basis function (RBF) kernel**

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

corresponds to a dot product in an infinite-dimensional space; it is a combination of polynomial kernels of all degrees. Assuming  $\gamma = 1$  for simplicity, we get

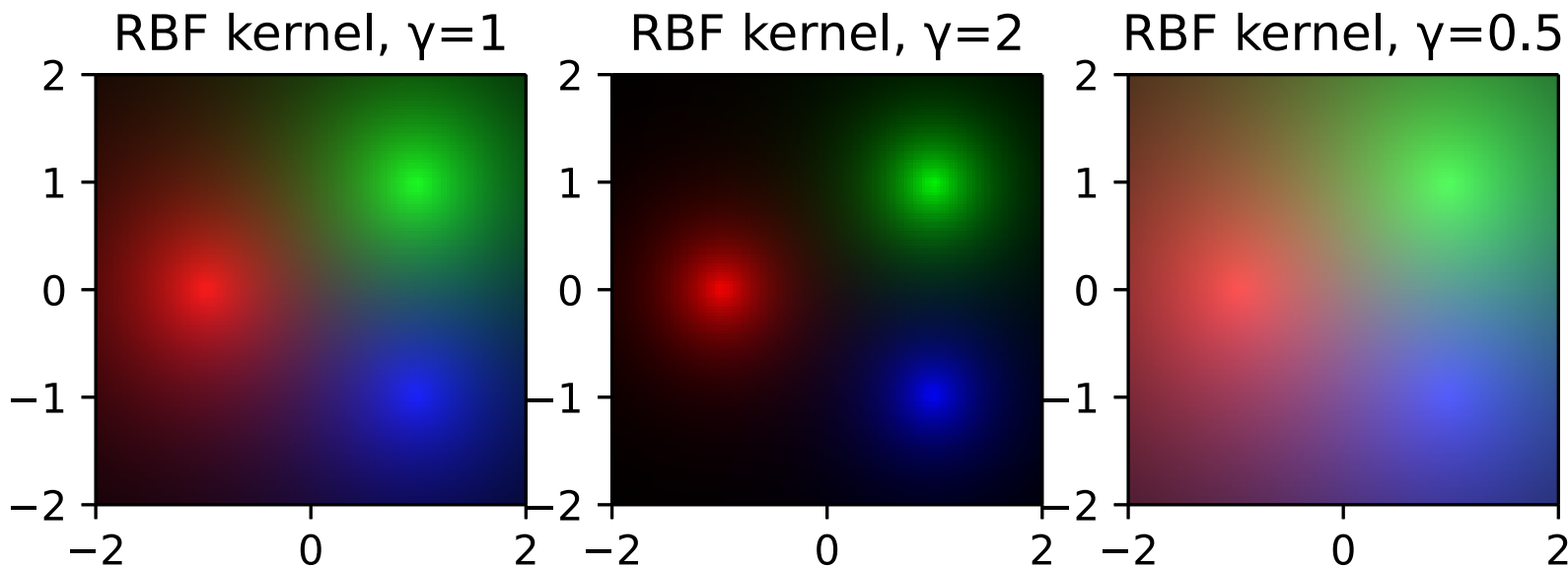
$$e^{-\|\mathbf{x} - \mathbf{z}\|^2} = e^{-\|\mathbf{x}\|^2 + 2\mathbf{x}^T \mathbf{z} - \|\mathbf{z}\|^2} = \sum_{d=0}^{\infty} \frac{(2\mathbf{x}^T \mathbf{z})^d}{d!} e^{-\|\mathbf{x}\|^2 - \|\mathbf{z}\|^2} = \sum_{d=0}^{\infty} \frac{2^d e^{-\|\mathbf{x}\|^2 - \|\mathbf{z}\|^2}}{d!} \left(\mathbf{x}^T \mathbf{z}\right)^d,$$

which is a combination of polynomial kernels; therefore, the feature map corresponding to the RBF kernel is

$$\varphi(\mathbf{x}) = \left( e^{-\gamma \|\mathbf{x}\|^2} \sqrt{\frac{(2\gamma)^d}{d!} \binom{d}{n_1, \dots, n_D}} x_1^{n_1} \cdots x_D^{n_D} \right)_{d \in \{0, 1, 2, \dots\}, n_i \geq 0, \sum_{i=1}^D n_i = d}.$$

Note that the RBF kernel is a function of distance – it "weights" more similar examples more strongly. We could interpret it as an extended version of the k-nearest neighbor algorithm, one which considers all examples, each weighted by similarity.

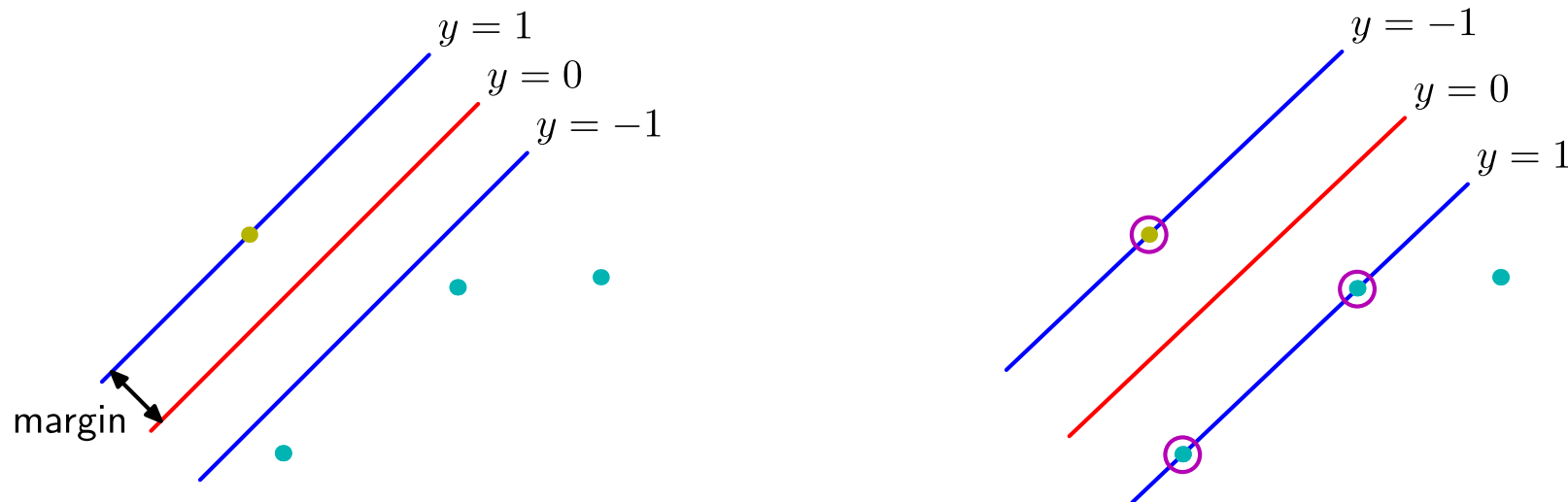
For illustration, we plot RBF kernel values of three points  $(0, -1)$ ,  $(1, 1)$  and  $(1, -1)$  with different values of  $\gamma$ :



# Support Vector Machines

Let us return to a binary classification task. The perceptron algorithm guaranteed finding some separating hyperplane if it existed (but it could find quite a bad one).

We now consider finding the one with **maximum margin**.



**Figure 7.1** The margin is defined as the perpendicular distance between the decision boundary and the closest of the data points, as shown on the left figure. Maximizing the margin leads to a particular choice of decision boundary, as shown on the right. The location of this boundary is determined by a subset of the data points, known as support vectors, which are indicated by the circles.

*Figure 7.1 of Pattern Recognition and Machine Learning.*

# Support Vector Machines

Assume we have a dataset  $\mathbf{X} \in \mathbb{R}^{N \times D}$ ,  $\mathbf{t} \in \{-1, 1\}^N$ , a feature map  $\varphi$  and a model

$$y(\mathbf{x}) \stackrel{\text{def}}{=} \varphi(\mathbf{x})^T \mathbf{w} + b.$$

We already know that the distance of a point  $\mathbf{x}_i$  to the decision boundary is

$$\frac{|y(\mathbf{x}_i)|}{\|\mathbf{w}\|} \stackrel{\text{assuming } y \text{ classifies all } \mathbf{x}_i \text{ correctly}}{=} \frac{t_i y(\mathbf{x}_i)}{\|\mathbf{w}\|}.$$

We therefore want to maximize

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_i \left[ t_i (\varphi(\mathbf{x}_i)^T \mathbf{w} + b) \right].$$

However, this problem is difficult to optimize directly.

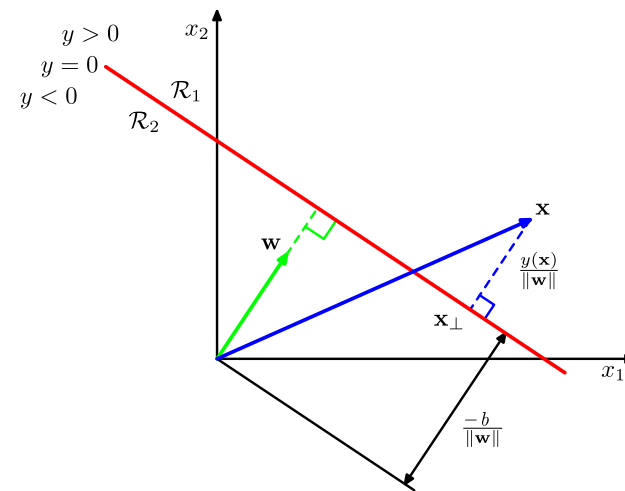


Figure 4.1 of Pattern Recognition and Machine Learning.

Because the model is invariant to multiplying  $\mathbf{w}$  and  $b$  by a constant, we can decide that for the points closest to the decision boundary, it will hold that

$$t_i y(\mathbf{x}_i) = 1.$$

Then for all the points we will have  $t_i y(\mathbf{x}_i) \geq 1$ , and we can simplify

$$\arg \max_{\mathbf{w}, b} \frac{1}{\|\mathbf{w}\|} \min_i \left[ t_i (\varphi(\mathbf{x}_i)^T \mathbf{w} + b) \right].$$

to

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{given that } t_i y(\mathbf{x}_i) \geq 1.$$

# Constrained Optimization – Inequality Constraints

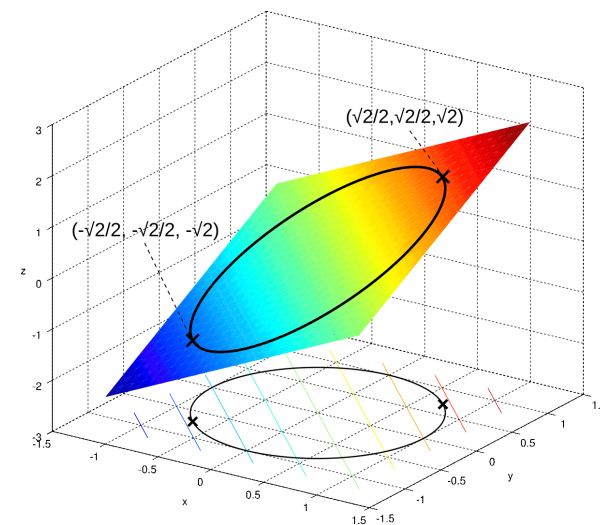
Given a function  $f(\mathbf{x})$ , we can find its minimum with respect to a vector  $\mathbf{x} \in \mathbb{R}^D$ , by investigating the critical points  $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ .

We can even incorporate constraints of form  $g(\mathbf{x}) = 0$  by forming a Lagrangian

$$\mathcal{L}(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

and again investigating the critical points  $\nabla_{\mathbf{x}, \lambda} \mathcal{L}(\mathbf{x}, \lambda) = 0$ .

We now describe how to include inequality constraints  $g(\mathbf{x}) \geq 0$ .



<https://upload.wikimedia.org/wikipedia/commons/e/ed/Lagrange>

# Constrained Optimization – Inequality Constraints

Our goal is to find a minimum of  $f(\mathbf{x})$  subject to a constraint  $g(\mathbf{x}) \geq 0$ .

We start by again forming a Lagrangian  $f(\mathbf{x}) - \lambda g(\mathbf{x})$ .

The optimum can either be attained for  $g(\mathbf{x}) > 0$ , when the constraint is said to be **inactive**, or for  $g(\mathbf{x}) = 0$ , when the constraint is said to be **active**. In the inactive case, the minimum is again a critical point of the Lagrangian with the condition  $\lambda = 0$ .

When the minimum is on a boundary, it corresponds to a critical point with  $\lambda \neq 0$  – but note that this time the sign of the multiplier matters, because the minimum is attained only when the gradient of  $f(\mathbf{x})$  is oriented **into** the region  $g(\mathbf{x}) \geq 0$ . We therefore require  $\nabla f(\mathbf{x}) = \lambda \nabla g(\mathbf{x})$  for  $\lambda > 0$ .

In both cases,  $\lambda g(\mathbf{x}) = 0$ .

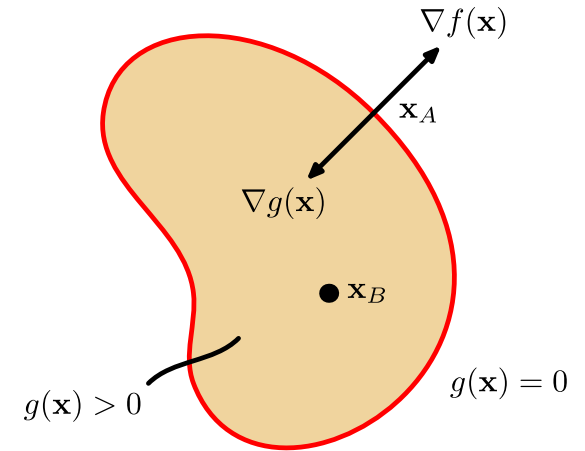
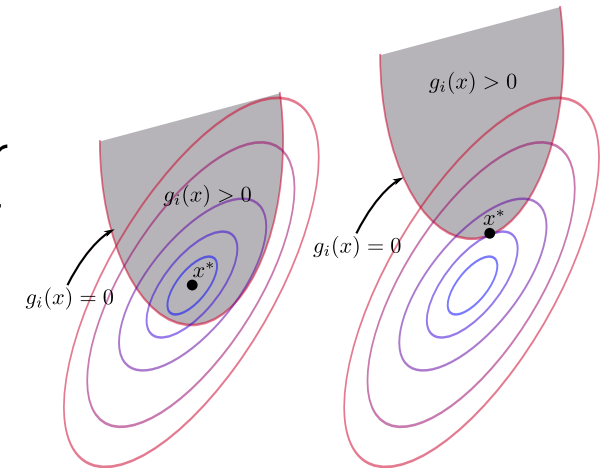


Figure E.3 of Pattern Recognition and Machine Learning.



<https://upload.wikimedia.org/wikipedia/commons/5/5d/Inequ>



# Minimization – Inequality Constraint

Let  $f(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$  be a function, which has a minimum in  $\mathbf{x}$  subject to an inequality constraint  $g(\mathbf{x}) \geq 0$ . Assume that both  $f$  and  $g$  have continuous partial derivatives and that  $\nabla_{\mathbf{x}}g(\mathbf{x}) \neq 0$ .

Then there exists a  $\lambda \in \mathbb{R}$ , such that the **Lagrangian function**

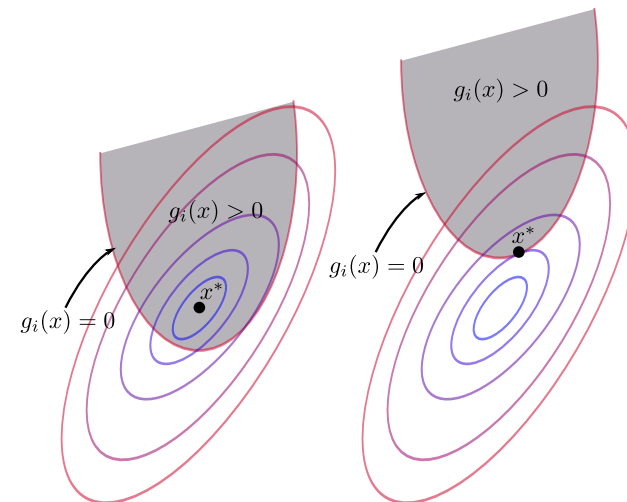
$$\mathcal{L}(\mathbf{x}, \lambda) \stackrel{\text{def}}{=} f(\mathbf{x}) - \lambda g(\mathbf{x})$$

has zero gradient in  $\mathbf{x}$  and the following conditions hold:

$$g(\mathbf{x}) \geq 0,$$

$$\lambda \geq 0,$$

$$\lambda g(\mathbf{x}) = 0.$$



<https://upload.wikimedia.org/wikipedia/commons/5/5d/Inequality>

These conditions are known as **Karush-Kuhn-Tucker (KKT)** conditions.

# Minimization – Inequality Constraint

It is easy to verify that if  $\mathcal{L}(\mathbf{x}, \lambda)$  has a minimum in  $\mathbf{x}$  and  $\lambda$  fulfills the KKT conditions  $g(\mathbf{x}) \geq 0$ ,  $\lambda \geq 0$ ,  $\lambda g(\mathbf{x}) = 0$ , the Lagrangian  $\mathcal{L}$  has a **maximum** in  $\lambda$  subject to  $\lambda \geq 0$ :

- if  $g(\mathbf{x}) = 0$ , then  $\mathcal{L}$  does not change when changing  $\lambda$ ,
- if  $g(\mathbf{x}) > 0$ , then  $\lambda = 0$  from the KKT conditions, which is a maximum of  $\mathcal{L}$ .

On the other hand, if  $\mathcal{L}(\mathbf{x}, \lambda)$  has a minimum in  $\mathbf{x}$ , and a maximum in  $\lambda$  subject to  $\lambda \geq 0$ , all the KKT conditions must hold:

- if  $g(\mathbf{x}) < 0$ , then increasing  $\lambda$  would increase  $\mathcal{L}$ ,
- if  $g(\mathbf{x}) > 0$ , then decreasing  $\lambda$  increases  $\mathcal{L}$ , so  $\lambda = 0$ .

## Maximizing Given $f(\mathbf{x})$

If we instead want to find the constrained maximum of  $f(\mathbf{x})$ , we can search for the minimum of  $-f(\mathbf{x})$ , which results in the Lagrangian  $f(\mathbf{x}) + \lambda g(\mathbf{x})$ , which we *minimize* with respect to  $\lambda$ .

# Necessary and Sufficient KKT Conditions

The KKT conditions are necessary conditions for a minimum (resp. a maximum).

However, it can be proven that in the following settings, the conditions are also **sufficient**:

- if the objective to optimize is a *convex* function (resp. *concave* for maximization) with respect to  $\mathbf{x}$ ;
- the inequality constraints are continuously differentiable convex functions;
- the equality constraints are affine functions (linear functions with an offset).

Therefore, if the above holds and if we find  $\mathbf{x}$  and  $\lambda$  such that:

- $\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}) = 0$ ,
- either
  - $g(\mathbf{x}) \geq 0$ ,  $\lambda \geq 0$ ,  $\lambda g(\mathbf{x}) = 0$ ,
  - or  $\lambda \geq 0$  and  $\mathcal{L}$  has a *maximum* in  $\lambda$ ,

then  $\mathbf{x}$  is a minimum of the function  $f(\mathbf{x})$  subject to an inequality constraint  $g(\mathbf{x}) \geq 0$ .

It is easy to verify that these conditions hold for the SVM optimization problem.

In order to solve the constrained problem of

$$\arg \min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{given that } t_i y(\mathbf{x}_i) \geq 1,$$

we write the Lagrangian with multipliers  $\mathbf{a} = (a_1, \dots, a_N)$  as

$$\mathcal{L} = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_i a_i (t_i y(\mathbf{x}_i) - 1).$$

Setting the derivatives with respect to  $\mathbf{w}$  and  $b$  to zero, we get

$$\begin{aligned} \mathbf{w} &= \sum_i a_i t_i \varphi(\mathbf{x}_i), \\ 0 &= \sum_i a_i t_i. \end{aligned}$$

Substituting these to the Lagrangian, we want to maximize

$$\mathcal{L} = \sum_i a_i - \frac{1}{2} \sum_i \sum_j a_i a_j t_i t_j K(\mathbf{x}_i, \mathbf{x}_j)$$

with respect to  $a_i$  subject to the constraints  $a_i \geq 0$  and  $\sum_i a_i t_i = 0$ , using the kernel  $K(\mathbf{x}, \mathbf{z}) = \varphi(\mathbf{x})^T \varphi(\mathbf{z})$ .

The solution will fulfill the KKT conditions, meaning that

$$a_i \geq 0, \quad t_i y(\mathbf{x}_i) - 1 \geq 0, \quad a_i (t_i y(\mathbf{x}_i) - 1) = 0.$$

Therefore, either a point  $\mathbf{x}_i$  is on a boundary, or  $a_i = 0$ . Given that the prediction for  $\mathbf{x}$  is  $y(\mathbf{x}) = \sum_i a_i t_i K(\mathbf{x}, \mathbf{x}_i) + b$ , we only need to keep the training points  $\mathbf{x}_i$  that are on the boundary, the so-called **support vectors**. Therefore, even though SVM is a nonparametric model, it needs to store only a subset of the training data.

The dual formulation allows us to use nonlinear kernels.

**Figure 7.2** Example of synthetic data from two classes in two dimensions showing contours of constant  $y(\mathbf{x})$  obtained from a support vector machine having a Gaussian kernel function. Also shown are the decision boundary, the margin boundaries, and the support vectors.

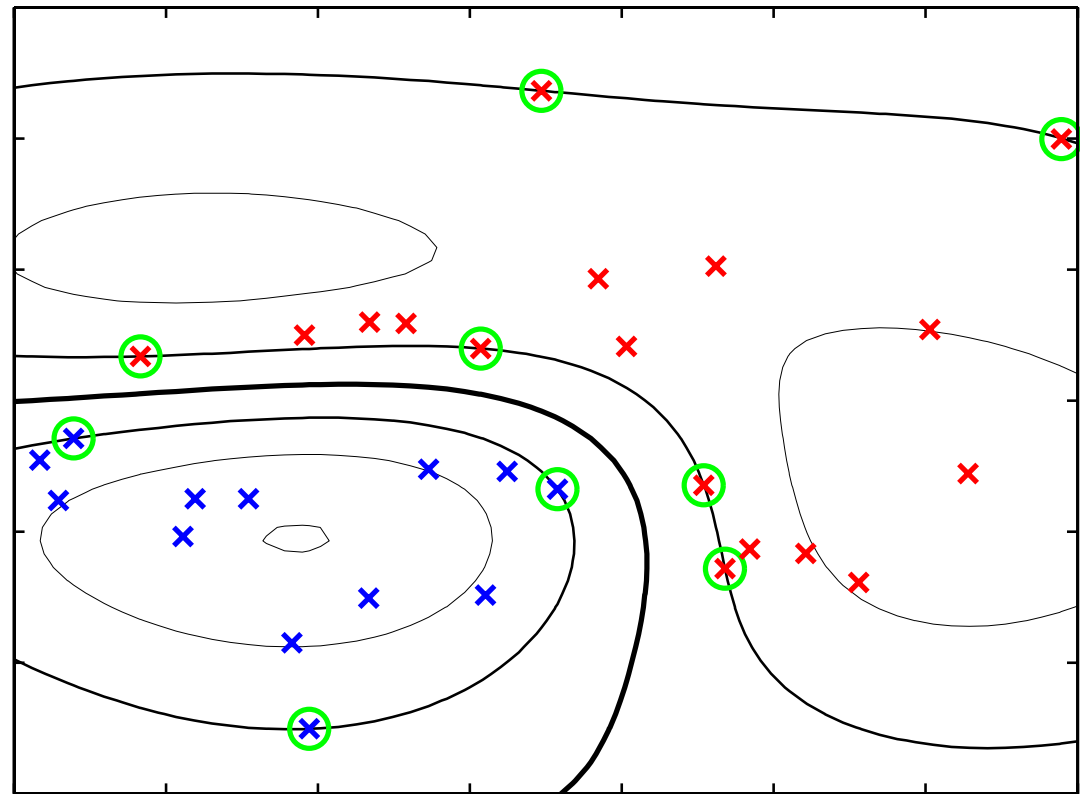


Figure 7.2 of Pattern Recognition and Machine Learning.