

## The Karush-Kuhn-Tucker (KKT) conditions

For several lectures we have been alluding to the **Karush-Kuhn-Tucker (KKT) conditions**. We are finally in a position to provide an intuitive account of what these are and where they come from. Simply put, the KKT conditions are a set of sufficient (and at most times necessary) conditions for an  $\mathbf{x}^*$  to be the solution of a given convex optimization problem. The conditions involve the existence of Lagrange multipliers satisfying certain natural properties, and they play a fundamental role in both the theory and practice of convex optimization.

We will start here by considering a general convex program with **inequality** constraints only. Specifically, we will consider the convex program given by

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} && f(\mathbf{x}) && (1) \\ & \text{subject to} && g_m(\mathbf{x}) \leq 0, && m = 1, \dots, M. \end{aligned}$$

Note that our initial restriction to inequality constraints only is merely to make the exposition easier – after we have this established, we will show how to include equality constraints (which must always be affine in convex programming).

Our analysis of the KKT conditions will exploit concepts from Lagrange duality. Recall that the Lagrangian for (1) is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = f(\mathbf{x}) + \sum_{m=1}^M \lambda_m g_m(\mathbf{x}),$$

the Lagrange dual function is given by

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}),$$

and the associated dual problem is given by

$$\underset{\boldsymbol{\lambda} \in \mathbb{R}^M}{\text{maximize}} \quad d(\boldsymbol{\lambda}) \quad \text{subject to} \quad \boldsymbol{\lambda} \geq \mathbf{0}. \quad (2)$$

Here and throughout this lecture, we assume that the functions  $f(\mathbf{x}), g_1(\mathbf{x}), \dots, g_M(\mathbf{x})$  are convex and differentiable. The KKT conditions for the optimization program in (1) are as follows.

### **KKT (inequality only)**

The KKT conditions for  $\mathbf{x} \in \mathbb{R}^N$  and  $\boldsymbol{\lambda} \in \mathbb{R}^M$  are

$$g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M, \quad (\text{K1})$$

$$\lambda_m \geq 0, \quad m = 1, \dots, M, \quad (\text{K2})$$

$$\lambda_m g_m(\mathbf{x}) = 0, \quad m = 1, \dots, M, \quad (\text{K3})$$

$$\nabla f(\mathbf{x}) + \sum_{m=1}^M \lambda_m \nabla g_m(\mathbf{x}) = \mathbf{0}. \quad (\text{K4})$$

Below we will argue that under certain conditions, a necessary and sufficient condition for  $\mathbf{x}^*$  to be a solution to (1) is for there to exist some  $\boldsymbol{\lambda}^*$  such that  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  satisfy the KKT conditions. Moreover, under the same conditions, a necessary and sufficient condition for  $\boldsymbol{\lambda}^*$  to be a solution to (2) is for there to exist some  $\mathbf{x}^*$  such that  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  satisfy the KKT conditions.

## Sufficiency

We start by establishing that these are sufficient.

Suppose that  $f(\mathbf{x}), g_1(\mathbf{x}), \dots, g_M(\mathbf{x})$  are convex and differentiable. If the KKT conditions hold for some  $\mathbf{x}^* \in \mathbb{R}^N$  and some  $\boldsymbol{\lambda}^* \in \mathbb{R}^M$ , then  $\mathbf{x}^*$  is a solution to (1) and  $\boldsymbol{\lambda}^*$  is a solution to (2).

Suppose that  $\mathbf{x}'$  and  $\boldsymbol{\lambda}'$  satisfy the KKT conditions. We will not use the notation  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  just yet to avoid any confusion, since we have not yet proven that they are solutions. We begin by noting that K4 is equivalent to the condition that

$$\nabla_{\mathbf{x}} \mathcal{L}(\mathbf{x}', \boldsymbol{\lambda}') = \mathbf{0}.$$

Since  $f$  and  $g_1, \dots, g_M$  are convex, this implies that  $\mathbf{x}'$  is a minimizer of  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}')$ , i.e.,

$$\mathcal{L}(\mathbf{x}', \boldsymbol{\lambda}') \leq \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}')$$

for all  $\mathbf{x} \in \mathbb{R}^N$ . This implies that

$$d(\boldsymbol{\lambda}') = \mathcal{L}(\mathbf{x}', \boldsymbol{\lambda}'),$$

and thus

$$\begin{aligned} d(\boldsymbol{\lambda}') &= f(\mathbf{x}') + \sum_{m=1}^M \lambda'_m g_m(\mathbf{x}') \\ &= f(\mathbf{x}'), \end{aligned}$$

where the last equality follows from (K3).

Recall that for *any* feasible  $\mathbf{x}$  and  $\boldsymbol{\lambda}$  we always have  $d(\boldsymbol{\lambda}) \leq f(\mathbf{x})$ . Since (K1) and (K2) ensure that  $\mathbf{x}'$  and  $\boldsymbol{\lambda}'$  are both feasible, the

above result implies that  $\mathbf{x}'$  is a solution to (1) and  $\boldsymbol{\lambda}'$  is a solution to (2): for any feasible  $\mathbf{x}$  we have  $f(\mathbf{x}) \geq d(\boldsymbol{\lambda}') = f(\mathbf{x}')$  and for any  $\boldsymbol{\lambda} \geq \mathbf{0}$  we have  $d(\boldsymbol{\lambda}) \leq f(\mathbf{x}') = d(\boldsymbol{\lambda}')$ .

Note that along the way we have also shown that the existence of  $\mathbf{x}, \boldsymbol{\lambda}$  satisfying the KKT conditions also implies strong duality.

## Necessity

We have just shown that for any convex problem of the form (1), if we can find a  $\mathbf{x}^*, \boldsymbol{\lambda}^*$  satisfying the KKT conditions, then  $\mathbf{x}^*$  must be a solution to (1). However, this on its own does not ensure that for any solution  $\mathbf{x}^*$  we must necessarily be able to find a  $\boldsymbol{\lambda}^*$  such that  $\mathbf{x}^*, \boldsymbol{\lambda}^*$  obey the KKT conditions. Indeed, we do not necessarily even know if a  $\mathbf{x}, \boldsymbol{\lambda}$  satisfying the KKT conditions exists.

Here we show that if we make the additional assumption that *strong duality* holds (for example, if our constraints satisfy Slater's condition), then the KKT conditions are also necessary.

Suppose that  $f(\mathbf{x}), g_1(\mathbf{x}), \dots, g_M(\mathbf{x})$  are convex and differentiable. Let  $\mathbf{x}^*$  be a solution to (1) and  $\boldsymbol{\lambda}^*$  be a solution to (2). If strong duality holds (e.g., Slater's condition holds) then  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  obey the KKT conditions.

We trivially have (K1) and (K2) simply because  $\mathbf{x}^*$  and  $\boldsymbol{\lambda}^*$  must be feasible to be solutions to (1) and (2) respectively. We very nearly proved the remaining conditions last time without explicitly saying so. In particular, we previously argued that if we have strong duality

then

$$\begin{aligned} f(\mathbf{x}^*) &= d(\boldsymbol{\lambda}^*) \\ &= \inf_{\mathbf{x} \in \mathbb{R}^N} \left( f(\mathbf{x}) + \sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}) \right) \\ &\leq f(\mathbf{x}^*) + \sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}^*) \\ &\leq f(\mathbf{x}^*). \end{aligned} \tag{3}$$

where the last inequality follows from the facts that we must have  $\lambda_m^* \geq 0$  and  $g_m(\mathbf{x}^*) \leq 0$ . Looking at this entire chain of inequalities, where the first and last term are both  $f(\mathbf{x}^*)$ , means that

$$f(\mathbf{x}^*) = \min_{\mathbf{x} \in \mathbb{R}^N} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*) = \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*).$$

Since  $f(\mathbf{x})$ ,  $g_1(\mathbf{x})$ ,  $\dots$ ,  $g_M(\mathbf{x})$  are convex and differentiable, so is  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$  for any  $\boldsymbol{\lambda}^* \geq \mathbf{0}$ . Thus, if  $\mathbf{x}^*$  is a minimizer of  $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}^*)$  then we must have  $\nabla \mathcal{L}(\mathbf{x}^*, \boldsymbol{\lambda}^*) = \mathbf{0}$ , which is precisely (K4).

Moreover, the chain of (in)equalities concluding in (3) also implies that

$$\sum_{m=1}^M \lambda_m^* g_m(\mathbf{x}^*) = 0.$$

Since each  $\lambda_m^* g_m(\mathbf{x}^*)$  is non-positive, the only way the sum can equal zero is if every term is zero, which yields (K3).

## KKT conditions with equality constraints

Now suppose that we have an optimization problem involving equality constraints:

$$\begin{aligned} & \underset{\mathbf{x} \in \mathbb{R}^N}{\text{minimize}} && f(\mathbf{x}) && (4) \\ & \text{subject to} && g_m(\mathbf{x}) \leq 0, && m = 1, \dots, M_I \\ & && h_m(\mathbf{x}) = 0, && m = 1, \dots, M_E. \end{aligned}$$

Using essentially the same arguments as before, we can show that the following KKT conditions on  $\mathbf{x}$ ,  $\boldsymbol{\lambda}$ , and  $\boldsymbol{\nu}$  are sufficient for  $\mathbf{x}$  and  $(\boldsymbol{\lambda}, \boldsymbol{\nu})$  to be solutions of (4) and its dual, respectively. Moreover, if strong duality holds they are also necessary.

### KKT (with equality constraints)

The KKT conditions for  $\mathbf{x} \in \mathbb{R}^N$ ,  $\boldsymbol{\lambda} \in \mathbb{R}^{M_I}$ , and  $\boldsymbol{\nu} \in \mathbb{R}^{M_E}$  are

$$g_m(\mathbf{x}) \leq 0, \quad m = 1, \dots, M_I, \quad (\text{K1})$$

$$h_m(\mathbf{x}) = 0, \quad m = 1, \dots, M_E,$$

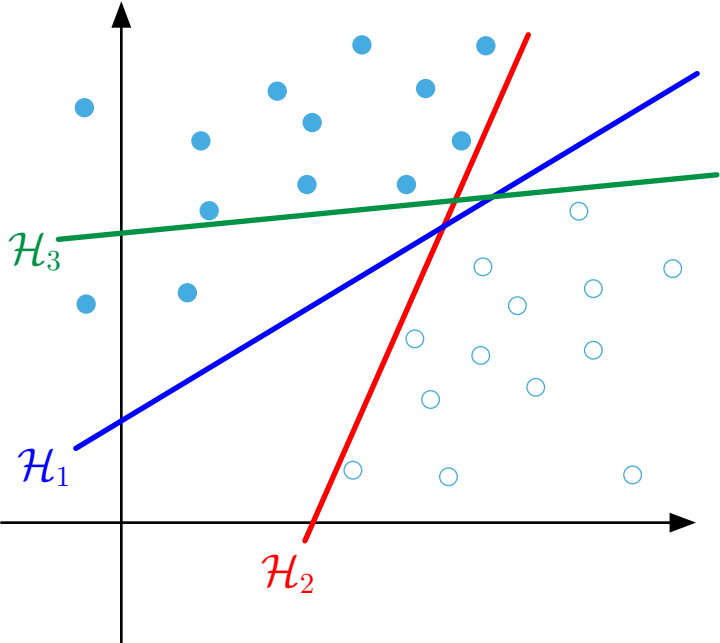
$$\lambda_m \geq 0, \quad m = 1, \dots, M_I, \quad (\text{K2})$$

$$\lambda_m g_m(\mathbf{x}) = 0, \quad m = 1, \dots, M_I, \quad (\text{K3})$$

$$\nabla f(\mathbf{x}) + \sum_{m=1}^{M_I} \lambda_m \nabla g_m(\mathbf{x}) + \sum_{m=1}^{M_E} \nu_m \nabla h_m(\mathbf{x}) = \mathbf{0}. \quad (\text{K4})$$

# Example: Support vector machines

Consider the following fundamental binary classification problem. We are given points  $\mathbf{x}_1, \dots, \mathbf{x}_M \in \mathbb{R}^N$  with class labels  $y_1, \dots, y_M$ , where  $y_m \in \{-1, +1\}$ . We would like to find a hyperplane (i.e., affine functional) which *separates* the classes:



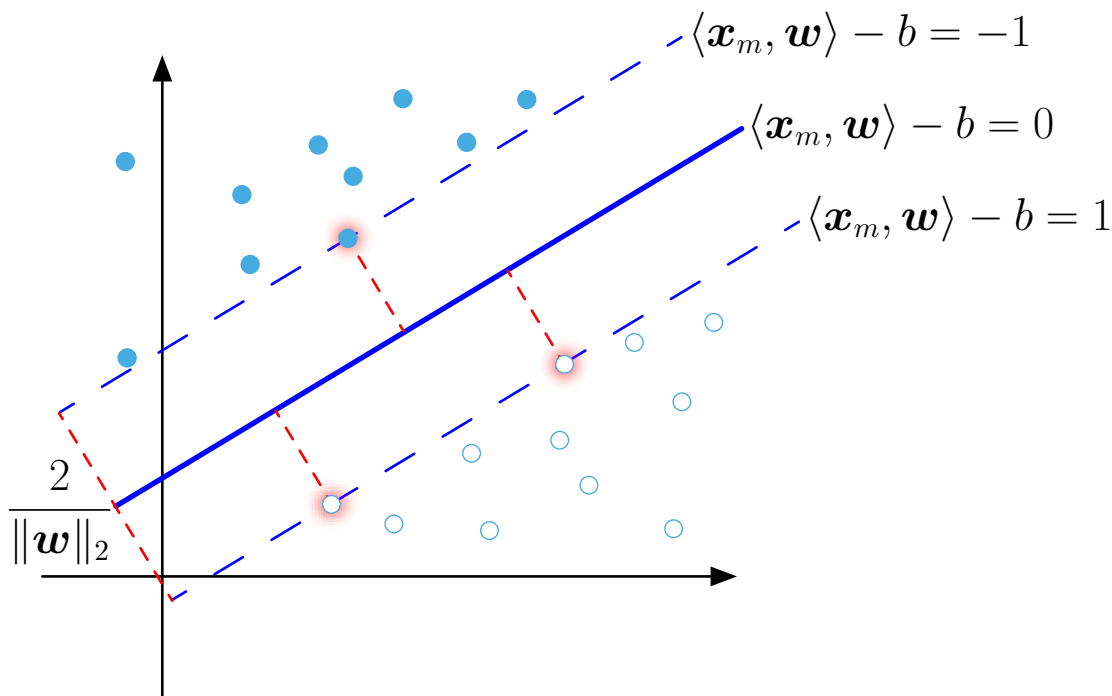
$\mathcal{H}_1$  and  $\mathcal{H}_2$  above both separate the classes in  $\mathbb{R}^2$ , but  $\mathcal{H}_3$  does not. While separating the classes is obviously desirable, we still need a good method to choose from among the many hyperplanes that do separate the classes – and some will perform better than others. Support vector machines (SVMs) take the one with *maximum margin*, i.e., we choose the hyperplane that maximizes the distance to the closest point in either class.

To restate this, we want to find a  $\mathbf{w} \in \mathbb{R}^N$  and  $b \in \mathbb{R}$  such that

$$\begin{aligned} \langle \mathbf{x}_m, \mathbf{w} \rangle - b &\geq 1, & \text{when } y_m = 1, \\ \langle \mathbf{x}_m, \mathbf{w} \rangle - b &\leq -1, & \text{when } y_m = -1. \end{aligned}$$

Of course, it is possible that no separating hyperplane exists; in this case, there will be no feasible points in the program above. It is straightforward, though, to modify this discussion to allow “misclassified” points.

In the formulation above, the distance between the two (parallel) hyperplanes is  $2/\|\mathbf{w}\|_2$ :



Thus maximizing this distance is the same as minimizing  $\|\mathbf{w}\|_2$ .



This leads to the program

$$\begin{aligned} & \underset{\mathbf{w} \in \mathbb{R}^N, b \in \mathbb{R}}{\text{minimize}} && \frac{1}{2} \|\mathbf{w}\|_2^2 \\ & \text{subject to} && y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1 \leq 0, \quad m = 1, \dots, M. \end{aligned}$$

This is a linearly constrained quadratic program, and is clearly convex. The Lagrangian is

$$\begin{aligned} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) &= \frac{1}{2} \|\mathbf{w}\|_2^2 + \sum_{m=1}^M \lambda_m [y_m(b - \langle \mathbf{x}_m, \mathbf{w} \rangle) + 1] \\ &= \frac{1}{2} \|\mathbf{w}\|_2^2 + b \boldsymbol{\lambda}^\top \mathbf{y} - \boldsymbol{\lambda}^\top \mathbf{X}^\top \mathbf{w} + \boldsymbol{\lambda}^\top \mathbf{1}, \end{aligned}$$

where  $\mathbf{X}$  is the  $N \times M$  matrix

$$\mathbf{X} = \begin{bmatrix} y_1 \mathbf{x}_1 & y_2 \mathbf{x}_2 & \cdots & y_M \mathbf{x}_M \end{bmatrix}.$$

The dual function is

$$d(\boldsymbol{\lambda}) = \inf_{\mathbf{w}, b} \left( \frac{1}{2} \|\mathbf{w}\|_2^2 + b \boldsymbol{\lambda}^\top \mathbf{y} - \boldsymbol{\lambda}^\top \mathbf{X}^\top \mathbf{w} + \boldsymbol{\lambda}^\top \mathbf{1} \right).$$

Since  $b$  is unconstrained above, we see that the presence of  $b \boldsymbol{\lambda}^\top \mathbf{y}$  means that the dual will be  $-\infty$  unless  $\langle \boldsymbol{\lambda}, \mathbf{y} \rangle = 0$ . Minimizing over  $\mathbf{w}$ , we need the gradient equal to zero,

$$\nabla_{\mathbf{w}} \mathcal{L}(\mathbf{w}, b, \boldsymbol{\lambda}) = \mathbf{0}, \quad \Rightarrow \quad \mathbf{w} - \mathbf{X} \boldsymbol{\lambda} = \mathbf{0}.$$

This means that we must have  $\mathbf{w} = \mathbf{X} \boldsymbol{\lambda}$ , which itself is a very handy fact as it gives us a direct passage from the dual solution to the primal

solution. With these substitutions, the dual function is

$$d(\boldsymbol{\lambda}) = \begin{cases} \frac{1}{2}\|\mathbf{X}\boldsymbol{\lambda}\|_2^2 - \boldsymbol{\lambda}^\top \mathbf{X}^\top \mathbf{X}\boldsymbol{\lambda} + \boldsymbol{\lambda}^\top \mathbf{1}, & \langle \boldsymbol{\lambda}, \mathbf{y} \rangle = 0, \\ -\infty, & \text{otherwise.} \end{cases}$$

Thus, the dual SVM program is then

$$\begin{aligned} & \underset{\boldsymbol{\lambda}}{\text{maximize}} && -\frac{1}{2}\|\mathbf{X}\boldsymbol{\lambda}\|_2^2 + \sum_{m=1}^M \lambda_m \\ & \text{subject to} && \langle \boldsymbol{\lambda}, \mathbf{y} \rangle = 0, \quad \boldsymbol{\lambda} \geq \mathbf{0}. \end{aligned}$$

Given the solution  $\boldsymbol{\lambda}^*$  to the dual, we can take  $\mathbf{w}^* = \mathbf{X}\boldsymbol{\lambda}^*$ , and the classifier is

$$\begin{aligned} f(\mathbf{x}) &= \langle \mathbf{x}, \mathbf{w}^* \rangle - b^* \\ &= \langle \mathbf{x}, \mathbf{X}\boldsymbol{\lambda}^* \rangle - b^* \\ &= \sum_{m=1}^M \lambda_m^* y_m \langle \mathbf{x}, \mathbf{x}_m \rangle - b^*. \end{aligned}$$

Notice that the data  $\mathbf{x}_m$  appear only through inner products with  $\mathbf{x}$ .

A key realization about the SVM is that the for the dual program, the objective function depends on the data  $\mathbf{x}_m$  only through inner products, as

$$\|\mathbf{X}\boldsymbol{\lambda}\|_2^2 = \sum_{\ell=1}^M \sum_{m=1}^M y_\ell y_m \langle \mathbf{x}_\ell, \mathbf{x}_m \rangle.$$

This means that we can replace  $\langle \mathbf{x}_\ell, \mathbf{x}_m \rangle$  with any “positive kernel function”  $K(\mathbf{x}_\ell, \mathbf{x}_m) : \mathbb{R}^N \otimes \mathbb{R}^N \rightarrow \mathbb{R}$  – a positive kernel just means that the  $M \times M$  matrix  $K(\mathbf{x}_\ell, \mathbf{x}_m)$  is in  $S_+^M$  for all choices of  $\mathbf{x}_1, \dots, \mathbf{x}_M$ .

For example, you might take

$$K(\mathbf{x}_\ell, \mathbf{x}_m) = (1 + \langle \mathbf{x}_\ell, \mathbf{x}_m \rangle)^2 = 1 + 2\langle \mathbf{x}_\ell, \mathbf{x}_m \rangle + \langle \mathbf{x}_\ell, \mathbf{x}_m \rangle^2.$$

This means we have replaced the inner product of two vectors with the inner product between two vectors which have been mapped into a higher-dimensional space:

$$\begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_N \end{bmatrix} \rightarrow \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_N \\ x_1^2 \\ x_2^2 \\ \vdots \\ x_N^2 \\ \sqrt{2}x_1x_2 \\ \vdots \\ \sqrt{2}x_{N-1}x_N \end{bmatrix}.$$

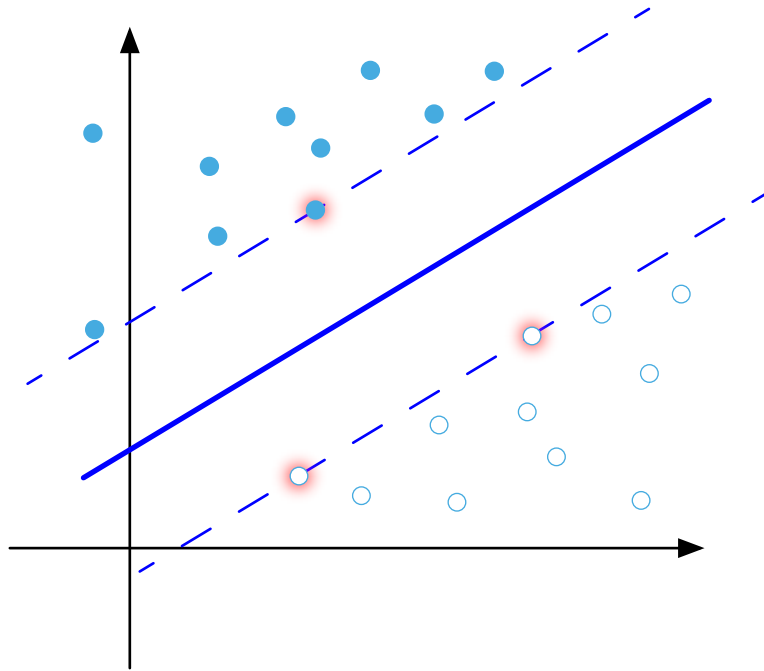
A set of linear constraints on the coordinates on the right, then, corresponds to a second order curve constraint (parabola, ellipse, hyperbola) on the coordinates on the left.

Many kernels are possible. The advantage is that to train and use the classifier, you never have to explicitly move to the higher-dimensional space – you just need to be able to compute  $K(\mathbf{x}_\ell, \mathbf{x}_m)$  for any pair of inputs in  $\mathbb{R}^N$ . A popular choice of kernel is

$$K(\mathbf{x}_\ell, \mathbf{x}_m) = \exp(-\gamma \|\mathbf{x}_\ell - \mathbf{x}_m\|_2^2).$$

This is a perfectly valid positive kernel, and it is straightforward to compute it for any pair of inputs. But it corresponds to mapping the  $\mathbf{x}_m$  into an infinite dimensional space, then finding a separating hyperplane.

Here is an example of a linear classifier in a higher-dimensional space:



that results in a nonlinear classifier in a lower-dimensional space:

