

Learning with Kernels

Proseminar Data Mining

Matthieu Bulté

Fakultät für Mathematik

Technische Universität München

Email: matthieu.bulte@tum.de

add to large margin intuition section that larger margin also reduces the space of possible hyperplanes

use jabref for refs

figure out where we need a vector space and where we need hilbert space add new decision function without computation of \mathbf{w} vector, explain why is is desired

Abstract—We present an algorithm for finding the observations relevant to classification of observations in two or more classes: the Support Vector Machine. The kernel trick, a generic method applicable to scalar product-based algorithm is shown to enable to let one incorporate domain knowledge about the classification problem in order to improve classification performance.

Index Terms—support vector machines, machine learning, statistics

I. INTRODUCTION

Machine learning can be seen as the science of automatically finding regularities or patterns in data. The contribution from many fields of science, engineering and mathematics brought a large variety of approaches to solving this problem.

In this work, we will bring our attention to the subject of kernel algorithms with a focus on support vector machines. While most of the focus nowadays goes into neural networks, support vector machine stays a relevant algorithm with useful theoretical, but also practical properties that we will study.

The first section will focus on presenting on a simplified version of the support vector machine: the hard margin classifier. We will show how the algorithm constructs the optimal hyperplane for separating two classes of observations. This will help construct the necessary intuition to understand classical extensions to the algorithm as well as the kernelization of the algorithm.

In the second section, we will study kernel methods and the kernel trick. A method for working with a projection of the data in a higher dimensional space without having to pay the computational cost of applying the projection. We will conclude the section with an overview of different kernels, and discuss how kernels helps us to incorporate domain knowledge to the learned classifier.

We will then conclude this paper with an overview of practical challenges that can be found while using support vector machines on modern data sets, and mention interesting extensions to support vector machines that were added in the 50 years of their existence.

II. MARGIN MAXIMIZATION CLASSIFIER

The hard margin classifier learns the parameters $\mathbf{w} \in \mathbb{R}^n$ and $b \in \mathbb{R}$ of a hyperplane $\mathcal{H} \subset \mathbb{R}^n$, separating two classes of observations $\mathbf{x}_i \in \mathcal{X} \subset \mathbb{R}^n$. To train it, the algorithm takes a set of m observations \mathbf{x}_i together with a vector of labels $y_i = \pm 1$ for each of the observations:

$$(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_m, y_m)$$

The learned parameters are then used to define a decision function f , assigning to points of \mathcal{X} a label ± 1 corresponding to the side of the hyperplane on which the points stands:

$$f(\mathbf{x}) = \text{sgn}(\langle \mathbf{w}, \mathbf{x} \rangle + b) \quad (1)$$

A. Maximizing the margin

One particularity about the learning algorithm used for Support Vector Machines is the loss function the algorithm tries to optimize. Other learning algorithms typically chose a loss function based on the empirical risk, defined as following for a function f and a training set (\mathbf{x}, y) of size m :

$$R_{\text{emp}}[f] = \frac{1}{m} \sum_{i=1}^m \frac{1}{2} |f(\mathbf{x}_i) - y_i| \quad (2)$$

The margin maximization algorithm instead, searches within the space of hyperplanes properly classifying the training set, for the hyperplane with the largest distance from the training points to the hyperplane.

The choice of maximizing the margin has good theoretical roots. Indeed, it is possible to derive an upper bound on the expected misclassification rate of a linear classifier based on the margin of the underlying hyperplane to the training set. It is beyond the scope of this work to go into further details but we will still attempt to provide an intuition on why it might be interesting to optimize it, and for more information the reader can refer to [ref?](#).

The intuition comes from the fact that the observations from the test set, or encountered while using the classifier, have been generated by a similar process as the observations present in the training set. Thus, one can assume that observations are very similar to each other up to some noise. Maximizing the hyperplane margin M is thus the same as increasing the tolerance to noise of the classifier. As shown in figure 1, given a observation x_i from the training set, all new observations

lying in the hypersphere of radius of the margin while be classified similarly to x_i .

This idea translates well into the following optimization problem, trying to maximize the margin, while adding a constrain to force the classifier to properly classify the training set

$$\begin{aligned} \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{maximize}} \quad & M := \frac{1}{\|\mathbf{w}\|} \min_i y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \\ \text{subject to} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq M \text{ for all } i = 1 \dots m \end{aligned} \quad (3)$$

Assuming that a solution to this optimization problem exists, it still has the issue that although the solution hyperplane is unique, there is an infinity of parameter vectors $\|\mathbf{w}\|$ resulting in the solution hyperplane, only differing in their length. Uniqueness of \mathbf{w} can be ensured by adding to its length the following constraint:

$$M\|\mathbf{w}\| = 1$$

We can thus reformulate the optimization problem into the following quadratic optimization problem which is known to have a numerical solution:

$$\begin{aligned} \underset{\mathbf{w} \in \mathbb{R}^n, b \in \mathbb{R}}{\text{minimize}} \quad & \|\mathbf{w}\| \\ \text{subject to} \quad & y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) \geq 1 \text{ for all } i = 1 \dots m \end{aligned} \quad (4)$$

B. Support vectors

Although the previous formulation of the problem helped us to better understand the ideas of support vector machines, the resulting quadratic problem will become impractical when later introducing kernels methods. We will introduce in this section an equivalent formulation of the optimization problem (4) that will not only solve computational issues, but will also give us more insights about the resulting solution.

In order to get to these benefits, we introduce the Lagrangian L together with Lagrange multipliers $\alpha_i \geq 0$:

$$L(\mathbf{w}, b, \alpha) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^m \alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] \quad (5)$$

This is called the dual formulation of our original optimization problem, called the primal problem, in which \mathbf{w} and b are called primal variables, and α the dual variable. It can be shown [ref?](#) that the primal optimization problem is equivalent to finding a saddle point of L , minimizing with respect to \mathbf{w} and b , while maximizing it with respect to α .

Because the solution of this dual problem is a saddle point, thus an extremum with respect to all variables. Thus our initial constraints are equivalent to the following constraints

$$\frac{\partial}{\partial b} L(\mathbf{w}, b, \alpha) = 0 \text{ and } \frac{\partial}{\partial \mathbf{w}} L(\mathbf{w}, b, \alpha) = 0 \quad (6)$$

leading to

$$\sum_{i=1}^m \alpha_i y_i = 0 \text{ and } \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \quad (7)$$

Further, by replacing (7) and (8) in the objective function (6), we obtain an optimization problem free of primal variables, and end up only optimizing against the vector α

$$\begin{aligned} \underset{\alpha \in \mathbb{R}^n}{\text{maximize}} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \\ \text{subject to} \quad & \alpha_i \geq 0 \text{ for all } i = 1 \dots m \\ & \sum_{i=1}^m \alpha_i y_i = 0 \end{aligned} \quad (8)$$

Finally, the Karush-Kuhn-Tucker theorem states that the solution α satisfies the following equality for all $i = 1 \dots m$ [ref?](#)

$$\alpha_i [y_i (\langle \mathbf{w}, \mathbf{x}_i \rangle + b) - 1] = 0 \quad (9)$$

The importance of this equality is twofold. First, it allows us to compute b from an observation i for which $\alpha_i \neq 0$. Second and most importantly, this equality implies that every observations not lying on the margin must have a vanishing Lagrangian coefficient. These special observations are called *support vectors*.

Support vectors have a special role. Indeed, looking at (7), one notices that the hyperplane only is only dependent on the support vectors. All other observations vanish from the sum because of their null Lagrange coefficient. This means that all the information extracted from the training set was represented by the support vectors, and all other observations had no role in the training of the algorithm, thus have no role in defining the decision boundary between the two classes of observations.

Furthermore, it can be shown that the amount support vectors give us an upper bound on the out of sample error of the classifier. This means that a low number of support vectors can be an indicator of the capacity of a classifier to generalize to unseen samples. The upper bound is following

$$E[Pr(\text{error})] \leq \frac{E[\text{number of support vectors}]}{\text{number of training observations}} \quad (10)$$

This inequality allows us to quickly derive a simple estimation of the quality of the learned classifier by just looking at the number of support vector.

C. Towards SVMs

We have chose in this work to present a simplified version of support vector machines called the hard margin classifier. This was done in order to focus on building an intuition of the way support vectors machines work. In practice, several considerations have to be made in order to make this classifier usable. These points will be handled in the next paragraphs.

The most important addition to the margin maximization algorithm that we have omitted so far is the ability to train

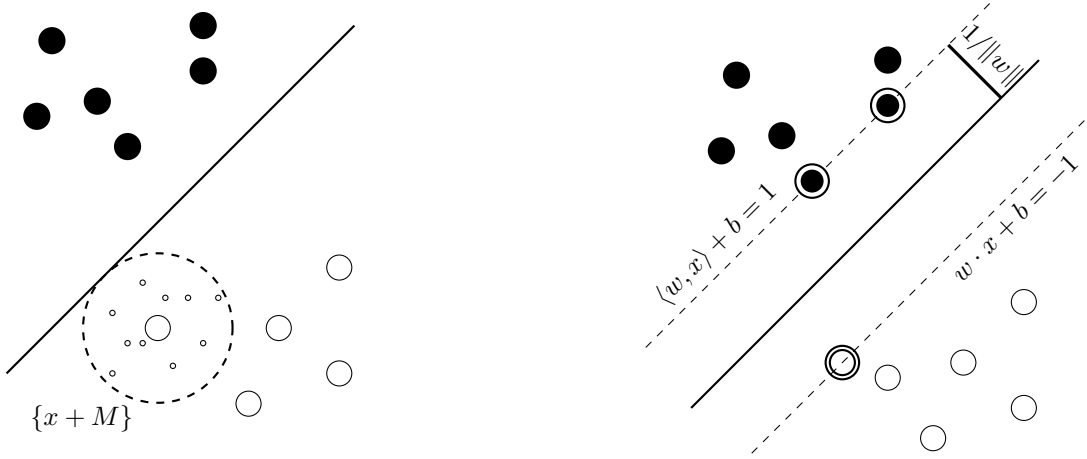


Fig. 1. Left side illustrates the idea behind margin maximization. Smaller points in the $\{x + M\}$ are generated by adding noise to an observation of the training set. The right figure shows a dataset separated by a hyperplane with parameters w, b . The norm of w is determined by the distance from the plane to the support vectors.

on a non linearly separable dataset caused by noise is the observations. Indeed, real world datasets are often generated from measuring physical phenomena. This kind of dataset is subject to noise that can cause an overlapping of the two classes around the decision boundary. Thankfully, these issues can be mitigated by introducing slack variables that will compensate the offset of an observation in the direction of the other class. Figure X [todo or remove?](#) illustrate the concept of these slack variables. A thorough investigation of this method can be found in [ref?](#).

III. KERNEL METHODS

The previous section introduced the large margin classifier. The large margin classifier, as we saw, allows us the find the best separating hyperplane between two linearly seperable classes of a dataset. In this section we introduce the concept of kernels, allowing us to learn non linear decision boudaries, continuing with the Mercer theorem and kernel trick, results from functional analysis allowing us to use the newly introduced ideas of kernels with the previously introduces algorithm, with only light modifications.

A. Non linearly separable data

Very often, the decision boundary one is trying to learn is in it's nature non-linear. One trick that is often used in order to use linear learning algorithms is to add one or several dimensions to the data points by applying a non linear function of the coordinates of the point, to then learn the separating hyperplane in the projected space. This method is known as features engineering, because we add engineered features to our data points.

This approach would be very simple to implement in our current algorithm. Let $\phi : \mathcal{X} \rightarrow \mathcal{V}$ be a non linear transformation from our original space \mathcal{X} to some higher dimensional Hilbert space \mathcal{V} . We can now modify the previously introduced algorithm by replacing simply every scalar product $\langle \cdot, \cdot \rangle$ by the scalar product on the projected points in the \mathcal{V} space, namely $\langle \phi(\cdot), \phi(\cdot) \rangle$.

In order to better understand kernels and the kernel trick, let's follow the example of polynomial projection of degree d , in which every point is projected to a vector containing every monomial of degree d . For instance, choosing $d = 2$ together with a 2 dimensional input space leads to the projection $\phi : \mathbb{R}^2 \rightarrow \mathbb{R}^4$ defined as

$$\phi(x_1, x_2) \mapsto (x_1^2, x_2^2, x_1x_2, x_2x_1)$$

This representation, while adding more separation capabilities has the problem that the image feature space grows at an exponential rate together with d , making the choice of a larger d prohibitivly expensive in terms of space usage. Though, the margin maximization algorithm only uses the scalar product of the observations in the feature space, and the following equations show that the projection in the high dimensional space is not required to the computation of the scalar product

$$\begin{aligned} \langle \phi(x), \phi(x') \rangle &= [x]_1^2 [x']_1^2 + [x]_2^2 [x']_2^2 + 2[x]_1 [x]_2 [x']_1 [x']_2 \\ &= \langle x, x' \rangle^2 \\ &=: k(x, x') \end{aligned}$$

We call $k(\cdot, \cdot) : \mathbb{R}^2 \times \mathbb{R}^2 \rightarrow \mathbb{R}$ the kernel representation of the scalar product in the space $\phi(\mathbb{R}^2)$. One can generalize this idea to any polynomial degree d as shown in [ref?](#), making the computation of the scalar product in the d -th monomial space as trivial as computing the scalar product in the original space.

B. Kernel trick

We have seen in the previous example, that it is sometimes possible to find a function $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with $k(x, x') = \langle \phi(x), \phi(x') \rangle$ for some projection $\mathcal{X} \rightarrow \mathcal{V}$. With these special functions, it is possible to then, only lightly modifying our original algorithm, to run our learning algorithm in another vector space without the computational cost of projecting our data in this other space. This is called the kernel trick.

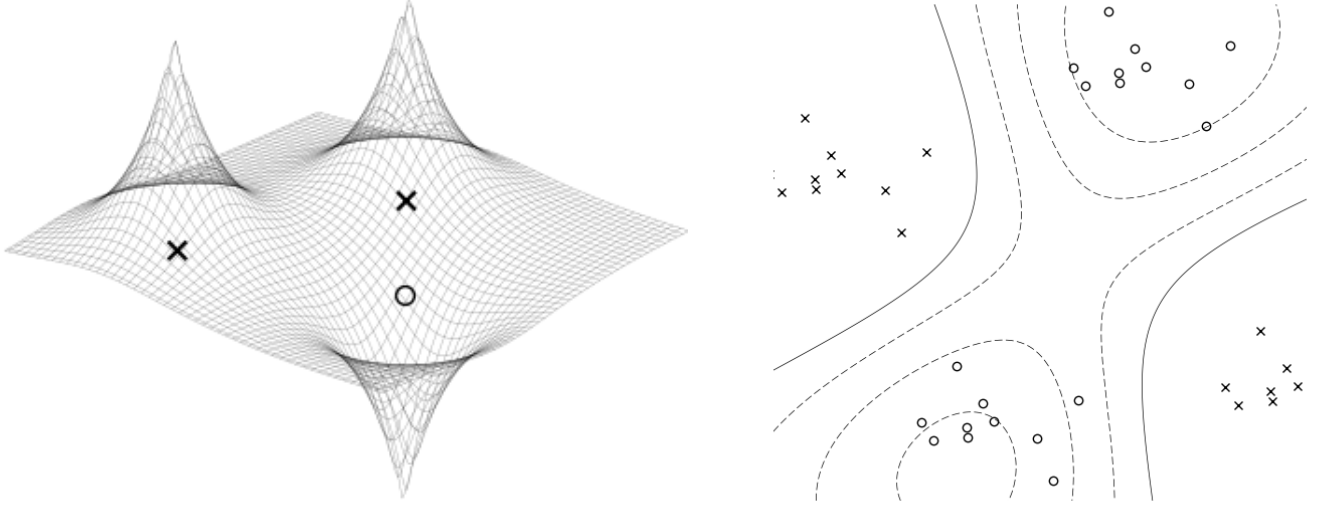


Fig. 2. Left side shows on the vertical axis the influence level of the point on the 2d plane from two positive points drawn with \times markers and of one negative points marked with \circ . On the right pane, the solid line represents the decision boundary and the dashed lines represents the contour levels of the decision function with learned using an RBF kernel.

We will now change our point of view, and try to define the properties defining the class of functions being the representation of a scalar product of the projection of its inputs in another vector space. This means, which properties must hold for k in order for a ϕ to exist with the correspondance property. **better name needed, also a definition would be nice**

Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ with $k(x, x') = \langle \phi(x), \phi(x') \rangle_{\mathcal{V}}$ for some projection $\mathcal{X} \rightarrow \mathcal{V}$, because $\langle \cdot, \cdot \rangle_{\mathcal{V}}$ is a scalar product, the following properties must hold

- *Symmetry*

$$\begin{aligned} \forall y_1, y_2 \in \mathcal{V}. \langle y_1, y_2 \rangle_{\mathcal{V}} &= \overline{\langle y_2, y_1 \rangle_{\mathcal{V}}} & \Rightarrow \\ \forall y_1, y_2 \in \phi(\mathcal{X}). \langle y_1, y_2 \rangle_{\mathcal{V}} &= \overline{\langle y_2, y_1 \rangle_{\mathcal{V}}} & \Rightarrow \\ \forall x_1, x_2 \in \mathcal{X}. \langle \phi(x_1), \phi(x_2) \rangle_{\mathcal{V}} &= \overline{\langle \phi(x_2), \phi(x_1) \rangle_{\mathcal{V}}} & \Rightarrow \\ \forall x_1, x_2 \in \mathcal{X}. k(x_1, x_2) &= \overline{k(x_2, x_1)} \end{aligned}$$

- *Positive definitness* Positive definitness of a kernel is a stronger property than the one required for the positive definiteness in the features space. A kernel is said to be positive definite if for every $\alpha_1, \dots, \alpha_n \in \mathbb{R}$ following inequality holds

$$\sum_{i,j=1}^n \alpha_i \alpha_j k(x_i, x_j) \geq 0$$

We will show by construction, that these properties are sufficient to the proof of the existence of the desired space. The literature contains several examples of vector spaces and projections with the desired property. We will here prove the existence of such a space by construction.

Let $k : \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ be symmetrical and positive definite, and $\mathbb{R}^{\mathcal{X}}$ be the set of functions from a non empty set \mathcal{X} to \mathbb{R} .

We define the reproducing kernel map as follow

$$\begin{aligned} \phi : \mathcal{X} &\rightarrow \mathcal{H} \\ x &\mapsto k(\cdot, x) \end{aligned} \quad (11)$$

We now show how the mapped observations of the training set $\{k(\cdot, x_1), \dots, k(\cdot, x_n)\}$ spans a Hilbert space in which the reproducing property hold with the scalar defined on the span set $\langle k(\cdot, x_i), k(\cdot, x_j) \rangle = k(x_i, x_j)$ and its canonical derivation for the spanned space

$$\begin{aligned} \langle f, g \rangle &= \left\langle \sum_{i=1}^n \alpha_i k(\cdot, x_i), \sum_{j=1}^n \beta_j k(\cdot, x_j) \right\rangle \\ &= \sum_{i,j=1}^n \alpha_i \beta_j \langle k(\cdot, x_i), k(\cdot, x_j) \rangle \\ &= \sum_{i,j=1}^n \alpha_i \beta_j k(x_i, x_j) \end{aligned}$$

Positive definitness, bilinearity and symmetry can all be trivially derived from the definition of this scalar product together with the symmetry and positive definitness of the kernel function. Then attentive reader will notice that the functions $\{k(\cdot, x_1), \dots, k(\cdot, x_n)\}$ must not necessarily be linearly independant, leading to a non unique mapping from f to the coefficients $\alpha_1, \dots, \alpha_n$. To see that this doesn't imply the ill definitness of the defined scalar product, we note that

$$\langle f, g \rangle = \sum_{j=1}^n \beta_j f(x_j)$$

and by symmetry

$$\langle f, g \rangle = \langle g, f \rangle = \sum_{i=1}^n \alpha_i g(x_i)$$

Which shows that the scalar product does not depend on the choice of coefficients for the functions f and g confirming the scalar product is well defined. Thus, the previously defined scalar product, together with the norm $|f| = \sqrt{\langle f, f \rangle}$ define a valid Hilbert space and the correspondence property holds

$$\begin{aligned} \langle \phi(x_i), \phi(x_j) \rangle &= \langle k(\cdot, x_i), k(\cdot, x_j) \rangle \\ &= k(x_i, x_j) \end{aligned}$$

Now that we have presented the kernel trick as well as the existence of a corresponding Hilbert space, the next section will focus on finding and creating useful kernels, as well as presenting some classical kernels.

C. Some useful kernels

We have defined in the last section what kernels are and how we can modify our algorithm to learn non-linear boundaries in our data. What we have omitted so far is how is one supposed to choose develop or choose a kernel when facing a practical problem.

One first way to choose a kernel is by using existing knowledge about the shape the decision boundary should have. Let's recall the equation of the decision boundary defined by the learned parameters w and b using the kernel k

$$\mathcal{H} = \{x \in \mathbb{R}^n \mid k(w, x) + b = 0\}$$

For illustration, if we choose the kernel to be a polynomial kernel of degree 2, meaning $k(x, x') = (\langle x, x' \rangle + 1)^2$, we can refine the definition of the decision boundary as the set of every $x \in \mathbb{R}^2$ for which following holds

$$\begin{aligned} 0 &= (\langle w, x \rangle + 1)^2 + b \\ &= (w_1 x_1 + w_2 x_2 + 1)^2 + b \\ &= (w_1 x_1)^2 + (w_2 x_2)^2 + 2w_1 x_1 w_2 x_2 + w_1 x_1 + w_2 x_2 + 1 + b \\ &= x^T \begin{bmatrix} w_1^2 & w_1 w_2 \\ w_1 w_2 & w_2^2 \end{bmatrix} x + w^T x + 1 \end{aligned}$$

The decision boundary formed by this equation is called a quadric, the generalization of a conic section, which takes a shape from a known set of different kind of solutions. Thus, if one knows a quadric decision boundary is needed, then using a polynomial kernel of degree 2 will learn the proper parameters of the quadric.

We have seen how previous knowledge on the shape of the decision boundary can lead to the choice of a particular kernel. Unfortunately, it is often the case that the best kind of decision boundary can't be determined. Thankfully, there is still a way to use domain knowledge about the problem to choose a kernel that will best help training a performant classifier.

In order to incorporate this other kind of domain knowledge, we have to look at kernels, not as the computations of a scalar product in some feature space, but as a similarity

measurement. Taking this point of view, the kernel trick becomes a way to find the proper Hilbert space in which one's definition of similarity defines a scalar proper product, letting us use the Support Vector Machine machinery as a way to train a classifier based on our definition of similarity.

What do we mean by similarity? Let's first write down a more geometrical definition of the euclidean scalar product in \mathbb{R}^n

$$\langle x, x' \rangle = \|x\| \|x'\| \cos(\text{angle}(x, x'))$$

This equivalent definition of the euclidean scalar product makes it more visible what the underlying notion of similarity of the euclidean scalar product is measured by the angle between the two vectors, scaled by their length. A similar kernel often used is the cosine similarity kernel, defined as

$$k(x, x') = \frac{x^T x'}{\|x\| \|x'\|} = \cos(\text{angle}(x, x'))$$

This notion of similarity will thus compare observations by first normalizing them to then compare the directions of the vectors. This similarity measure is often used in the context of text classification, where each entry of the vector corresponds to the number of occurrences of a word in a text. The normalization will transform count of occurrences in frequencies, thus having a similarity based on the relative importance of each word, independantly of the length of the text.

Cosine similarity is an interesting similarity, but one natural way to reason about similarity of points is their distance from one another. The next kernel we introduce incorporates the notions of distance as a similarity measurement, is called the Radial Basis Function kernel and is defined as followed for some hyperparameter $\sigma \in \mathbb{R}$

$$k(x, x') = \exp\left(-\frac{\|x - x'\|^2}{\sigma^2}\right)$$

This kernel contains the idea that similarity between two observations should decay exponentially with the distance between the two vectors. This kernel isn't only a good model of our notion of similarity, but it also has interesting properties when used with the Support Vector Machines training algorithm.

Let's first recall what was the decision function that was learned by the algorithm

$$f(x) = \text{sgn}\left(\sum_{x_i \text{ is SV}} \alpha_i y_i k(x, x_i) + b\right)$$

We evaluate the sign of a linear combination of the similarity of the current observation to the support vectors. One can see understand each support vector as a point of influence in the direction of the class it belongs to, with it's influence intensity exponentially decaying with the distance. The observation will then be classified by taking a weighted sum of the influence of each support vector and verifying if it is over or

under some threshold b . This intuition can be visualised in figure 2 where we show the influence map built by two points on in a one dimensional space as well as a decision boundary learned using a radial basis function kernel.

IV. ZUSAMMENFASSUNG UND AUSBLICK

blabla