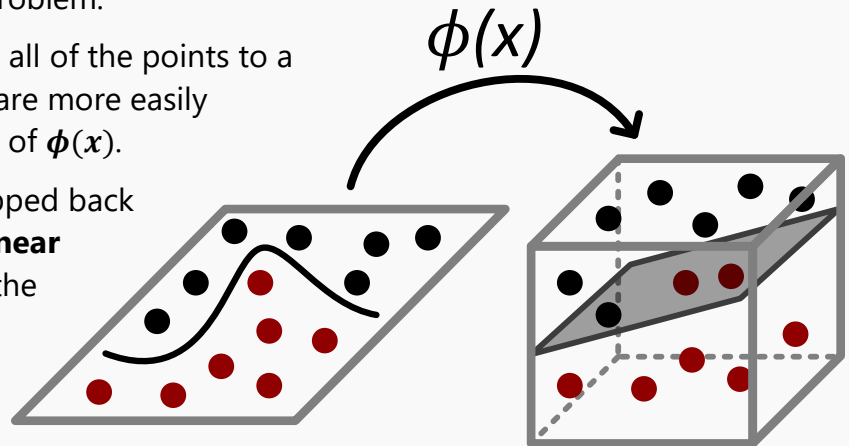


kernels for support vector machines 🏠 (SVM)

Kernels allow for **Support Vector Machines** to convert from **linear models** to **nonlinear models** by utilizing a slightly different optimization problem:

The **Kernel “trick”** allows an **SVM** to map all of the points to a high dimensional space where the points are more easily separated; achieved through the mapping of $\phi(x)$.

Then when the **decision boundary** is mapped back to the original dimensional space, a **nonlinear** boundary results that perfectly separates the data. The above illustration can be slightly misleading because the feature $\phi(x)$ is almost never seen in the process.



An illustrative example of Kernels

Assuming a housing problem with the following features in the dataset:

$$x_i = [x_{i1}, x_{i2}, x_{i3}, x_{i4}, \dots] \rightarrow x_{i1} = \text{List Price}, x_{i2} = \text{Estimate}, x_{i3} = \text{Time on Market}, x_{i4} = \text{Avg Price}$$

For simplicity, the two features below will be of focus:

$$x_i = [x_{i1}, x_{i2}] \rightarrow x_{i1} = \text{House's List Price}, x_{i2} = \text{Zillow Estimate}$$

The **inner product** of features $[x_{i1}, x_{i2}]$ is calculated as a metric to understand **distance in space** (2D)

$$\text{Inner Product } (x_i, x_k) = x_i \times x_k = x_{i1}x_{k1} + x_{i2}x_{k2}$$

In attempts to translate the above **2 dimensional** example into a **3 dimensional** space, the points are mapped through the function $\phi(x)$ as follows:

$$x_i = [x_{i1}, x_{i2}]$$

$$\phi(x_i) = \phi([x_{i1}, x_{i2}]) = [x_{i1}^2, x_{i2}^2, x_{i1}x_{i2}] \quad (2D \rightarrow 3D \text{ for } i)$$

$$\phi(x_k) = \phi([x_{k1}, x_{k2}]) = [x_{k1}^2, x_{k2}^2, x_{k1}x_{k2}] \quad (2D \rightarrow 3D \text{ for } k)$$

x is above being mapped from a **2D** to a **3D** space and the **inner product** can be computed in **3D**:

$$\text{Inner Product } (\phi(x_i), \phi(x_k)) = \phi(x_i) \times \phi(x_k) = x_{i1}^2, x_{i2}^2 + x_{k1}^2, x_{k2}^2 + x_{i1}x_{i2}x_{k1}x_{k2}$$

The above multiplies each dimension **component-wise** and computes the summation.

Therefore, one way to map to a high dimensional space with **SVM** is:

- Decide on the term $\phi(x)$
- Replace all values of x_i in the optimization problem with $\phi(x_i)$
- Solve for the new optimization problem

There is a special property of the SVM optimization problem:

- The only terms involving \mathbf{x}_i are **Inner Product** $(\mathbf{x}_i, \mathbf{x}_k)$
- \mathbf{x}_i will never occur alone or in other ways, they will only occur in inner products

$$\max_{\alpha} \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i, k=1}^m \alpha_i \alpha_j y_i y_k \quad \textcircled{x_i^T x_k} \leftarrow \text{Inner Product}$$

In the objective above, it is noted that the only occurrence of x is inside the **inner product**. Therefore, rather than determining what ϕ is for the purpose of mapping the data to a high dimensional space, the **inner product** in the original space can be replaced with an **inner product** in a different space.

The question arises: How to compute the **inner products** in the new space without computing ϕ 's?

$$x_i = [x_{i1}, x_{i2}]$$

$$\textbf{Inner Product } (\phi(x_i), \phi(x_k)) = \phi(x_i) \times \phi(x_k) = x_{i1}^2, x_{i2}^2 + x_{k1}^2, x_{k2}^2 + x_{i1}x_{i2}x_{k1}x_{k2}$$

It does not matter if ϕ is infinite dimensional or impossible to calculate, as long as the **Kernel** can be calculated between two points, the optimization formula can be run.

A particularly useful Kernel in practice is the **Gaussian Kernel**:

$$K(\mathbf{x}_i, \mathbf{x}_k) = \exp\left(-\frac{\|\mathbf{x}_i - \mathbf{x}_k\|_2^2}{2\sigma^2}\right), \text{ where } \sigma \text{ is selected}$$

Kernel:

the

Gaussian Kernel

x_j x_k

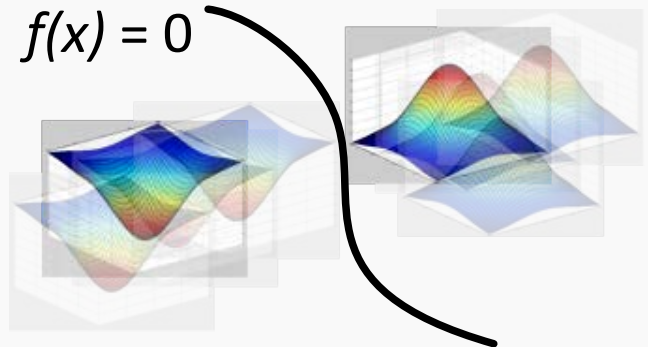
$\phi(x_j)$ $\phi(x_k)$

When using a **Gaussian Kernel**, also referred to as a **Radial Basis Function (RBF) Kernel**, a **normalized, high-dimensional distribution** is essentially being applied to each datapoint. They are subsequently weighted and summed. The distributions can be both **positive (+)** and **negative (-)** depending on the solution to the optimization problem. When a prediction is computed for a new observation in the dataset, a weighted summation of the normalized distributions is computed:

$$f(x) = \sum_{i=1}^n \alpha_i K(x, x_i) + b$$

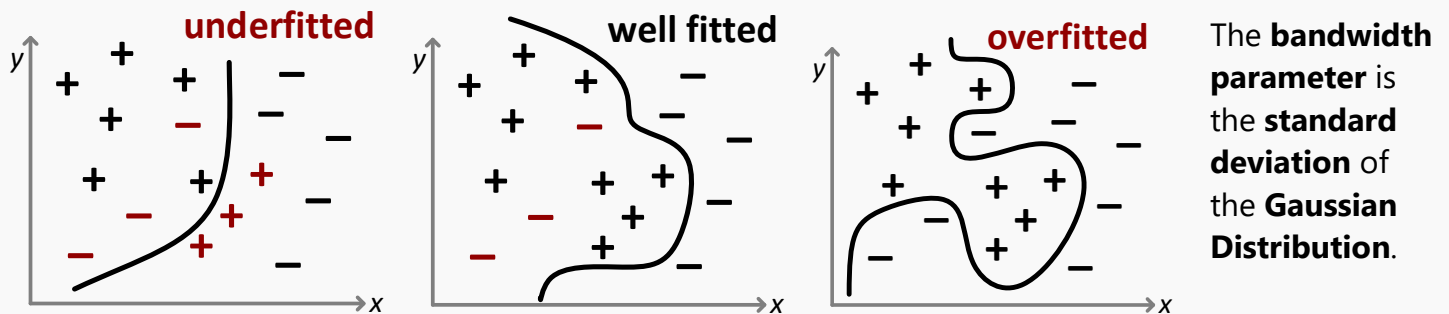
→ the **weights** α_i are determined by the solution to the SVM optimization problem

The **Decision Boundary** is simply constructed where the function $f(x) = 0$ as illustrated →



The examples above illustrated how **Support Vector Machines** are able to compute **complex nonlinear decision boundaries**.

It is important to note that if the **bandwidth parameter** is too small, the model can **overfit** data.



Common Kernels for Support Vector Machines

- Linear Kernel (synonymous to no kernel and using the regular inner product)
- Polynomial (the degree of polynomial is chosen and will overfit the data if too large)
- Gaussian Kernel (the bandwidth is chosen and will overfit the data if too large)

Advantages of Support Vector Machines

- Easily handles nonlinearities (with Kernels)
- Based on quadratic programming with specialized solvers (like coordinate descent with dual-adjustments to variables at a time)
- Reproducible results regardless of user or runs (theoretically)
- Easily handles imbalanced data by reweighting the points

Disadvantages of Support Vector Machines

- The solvers can sometimes be slow
- Does not naturally scale to larger datasets
- The performance is not as reliable in practice (even after adjusting the kernel parameters)
- Support Vector Machine models are uninterpretable

- The choice of Kernel can be difficult