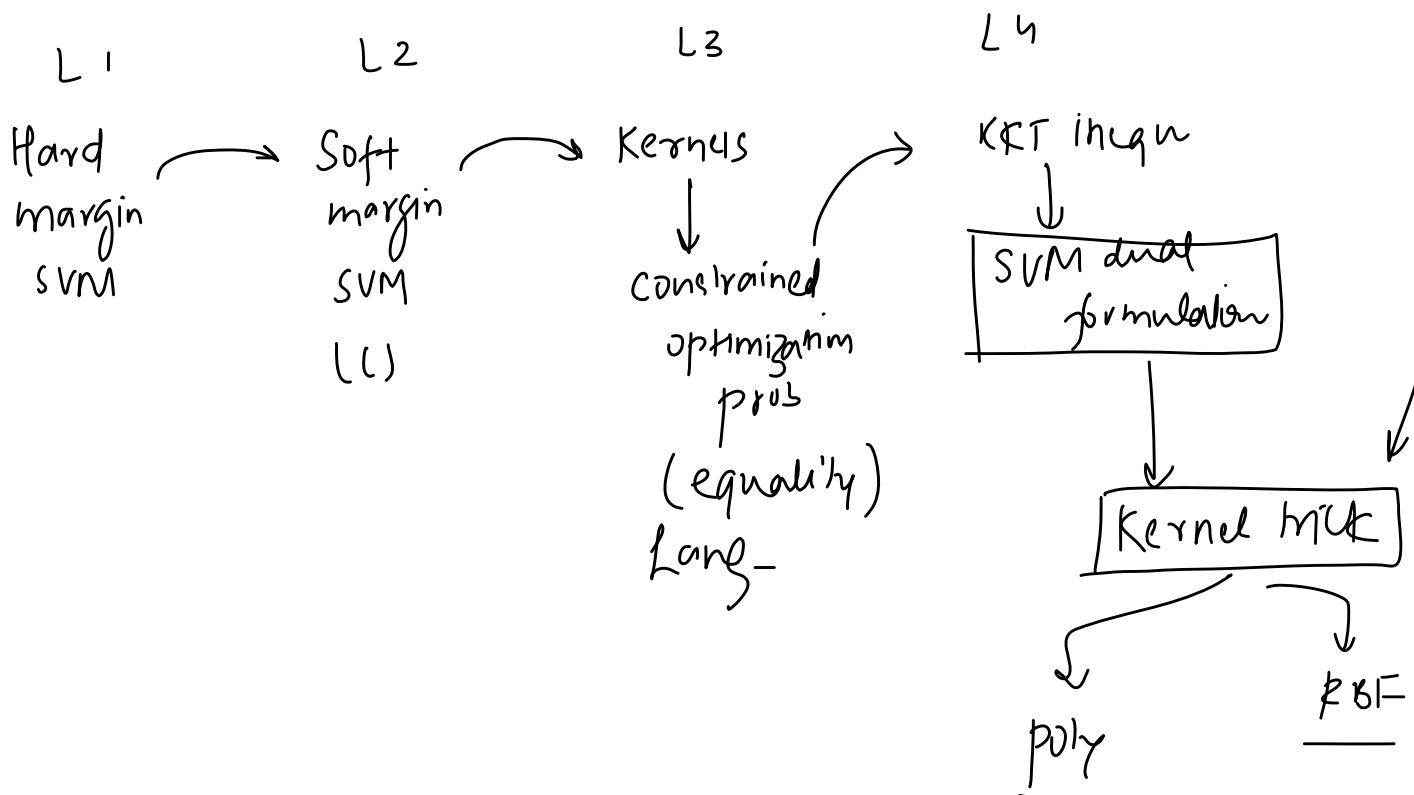


Recap

19 July 2023

17:04

SVM



19 July 2023 19:03

$$\underline{x_i \cdot x_j}$$

$$x_i, x_j$$

$$y_i y_j = 1 \times 1$$

$$X_1 \cdot X_1 \rightarrow X_{11} X_{11} + X_{12} X_{12}$$

25 terms

$$\alpha_i \geq 0$$

$$\sum_{i=1}^n d_i y_i = 0$$

$$x_{11}x_{21} + x_{12}x_{22}$$

$$- \left(S V \right)$$

$$\alpha_i > 0$$

$$\alpha_j = 0$$

$$\left\{ \begin{array}{l} \rightarrow \alpha_i = 0 \text{ for all non support vectors} \\ \rightarrow \alpha_i > 0 \text{ for all SV} \end{array} \right.$$

$$\underbrace{X_1 \quad X_2}$$

$\max_{\alpha_i} \alpha_1 + \alpha_2 + \alpha_3 + \alpha_4 + \alpha_5 - \frac{1}{2} \left[\begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \right] + \begin{array}{c} \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \\ \text{---} \end{array} \end{array}$

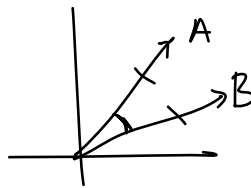
The Similarity Perspective

18 July 2023 08:52

$$\text{cosine similarity} = \frac{A \cdot B}{\|A\| \|B\|}$$

$$\|A\|=1$$

$$\|B\|=1$$



kernel trick

$$\rightarrow \boxed{A \cdot B} \rightarrow \text{dot product of } A \text{ and } B$$

$$\max_{\alpha_i} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{(x_i \cdot x_j)}_{\text{similarity}}$$

x_i and x_j

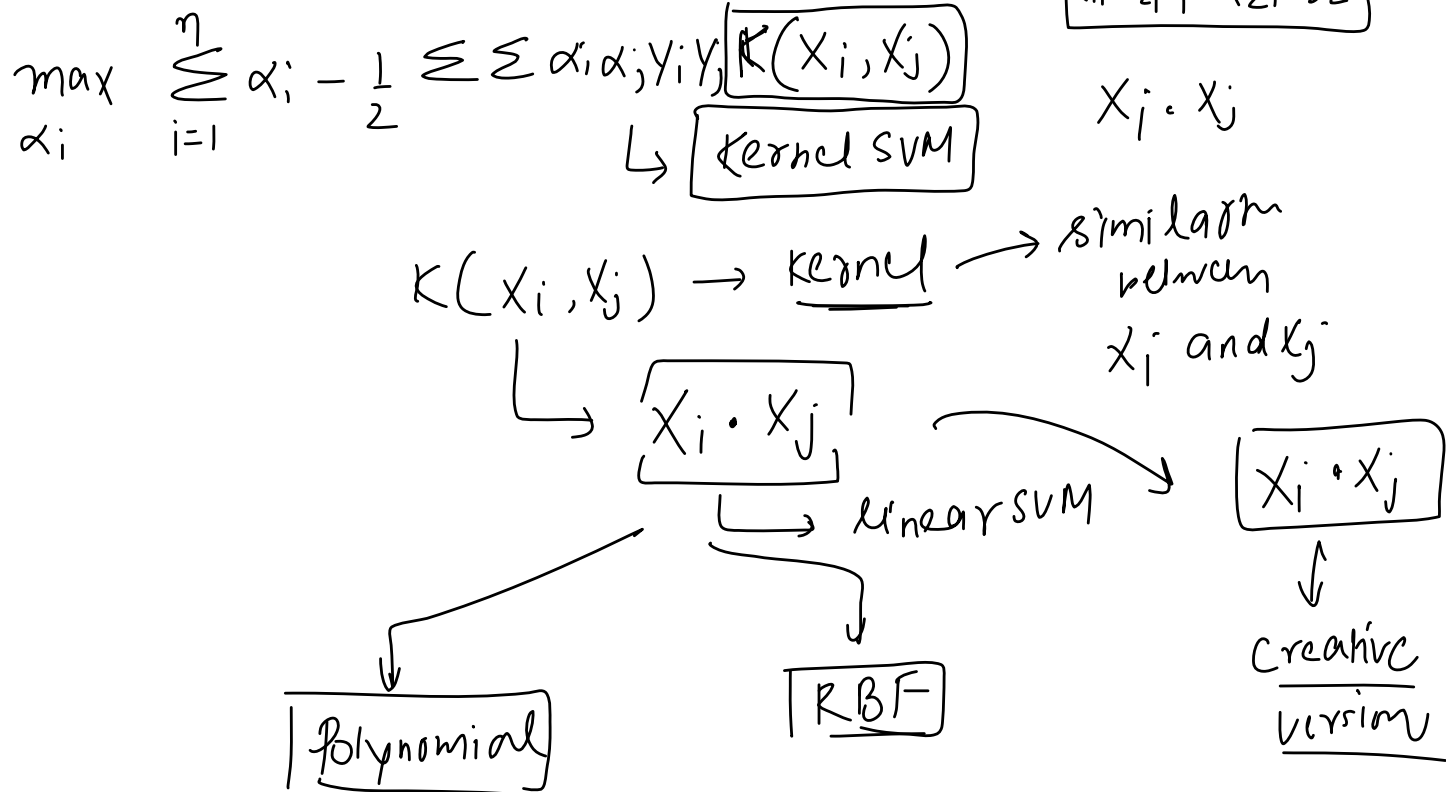
maximizing the similarity of S_V based on mis sign

$$(x_i \cdot x_j) \rightarrow \boxed{\text{sim}(x_i, x_j)}$$

↓
kernels

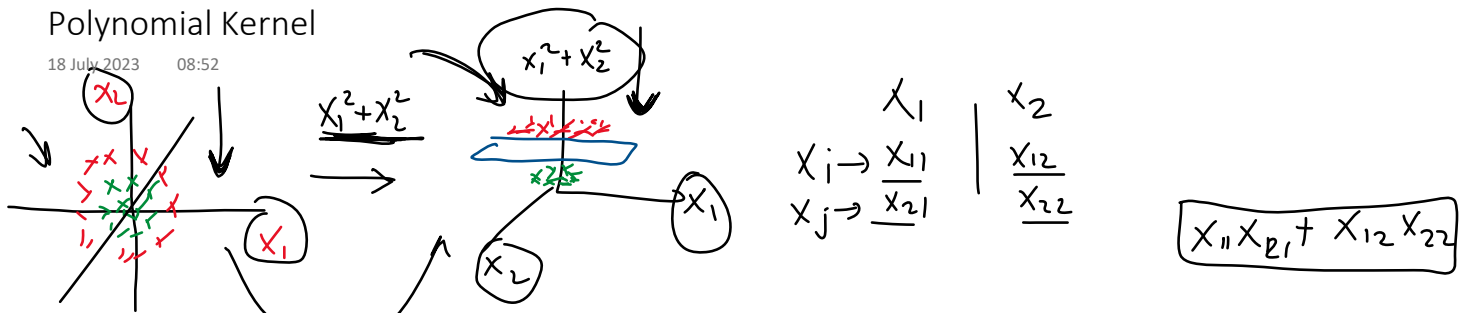
Kernel SVM

19 July 2023 07:41



Polynomial Kernel

18 July 2023 08:52



$$K(x_i, x_j) = \left(\gamma + x_i \cdot x_j \right)^d \quad \gamma=1, d=2, 3$$

$$(1 + x_i^T x_j)^k \quad \gamma=1, d=2$$

$$= (1 + x_i \cdot x_j)^2 = (1 + x_{11}x_{21} + x_{12}x_{22})^2$$

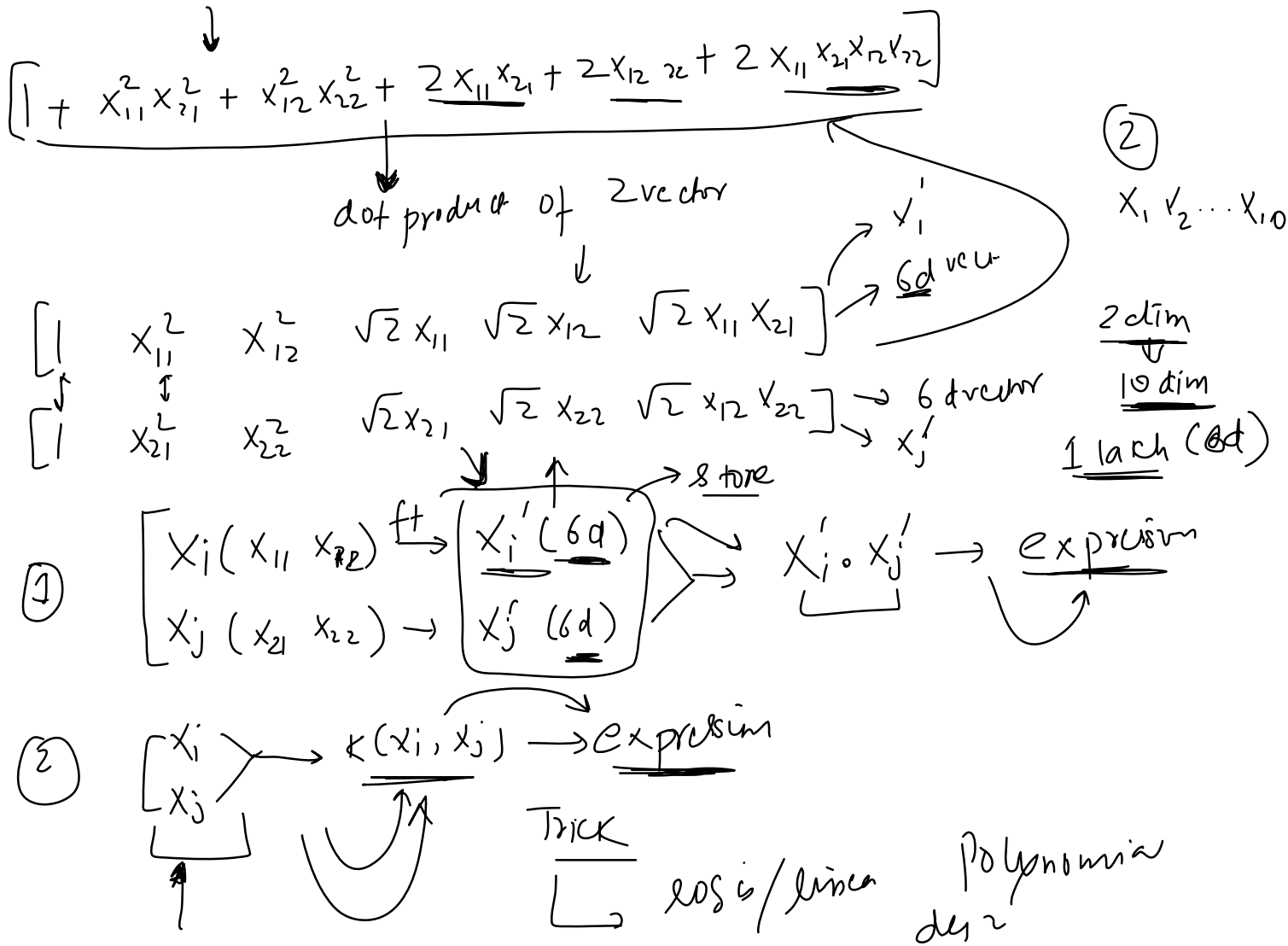
$$= 1 + x_{11}^2 x_{21}^2 + x_{12}^2 x_{22}^2 + 2x_{11}x_{21} + 2x_{12}x_{22} + 2x_{11}x_{21}x_{12}x_{22}$$

kernel trick

↳ polynomial term

The Trick

19 July 2023 07:30

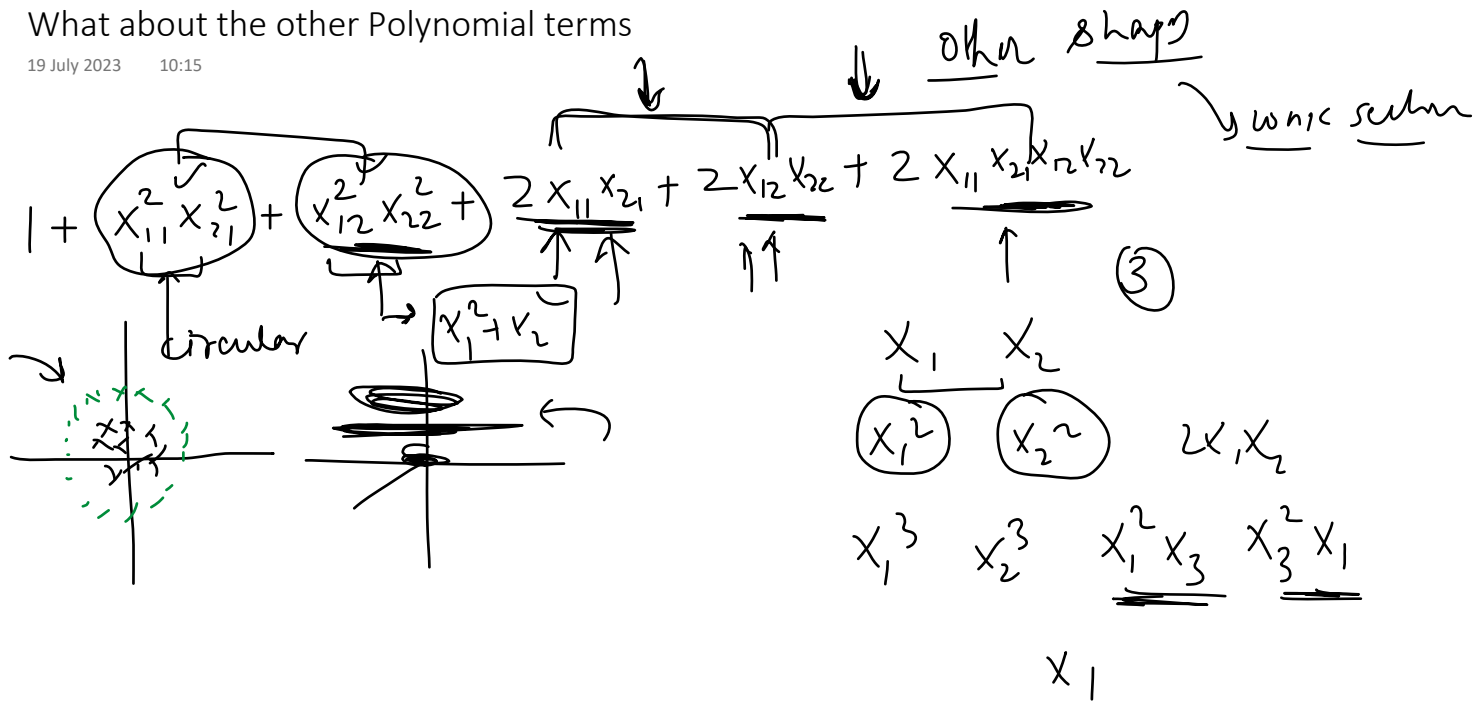


$$x_1 \quad x_2$$

$$x_1 \mid x_1^2 \mid x_2 \mid x_2^2 \mid x_1 x_2$$

What about the other Polynomial terms

19 July 2023 10:15



Radial Basis function { Normal dist }

- popular
- Best out of the box kernel
- powerful

$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

polynomial

x x
x x
xx

$$e^{-\gamma \|x_i - x_j\|^2}$$

hyperparameter

$$\gamma = \frac{1}{2\sigma^2}$$

$$K \propto \frac{1}{\text{dist}}$$

similarity

$\|x_i - x_j\| \rightarrow$ euclidean dist
between x_i and x_j

neural networks

$$K(x_i, x_j) = e^{-\frac{\text{dist}^2}{2\sigma^2}}$$

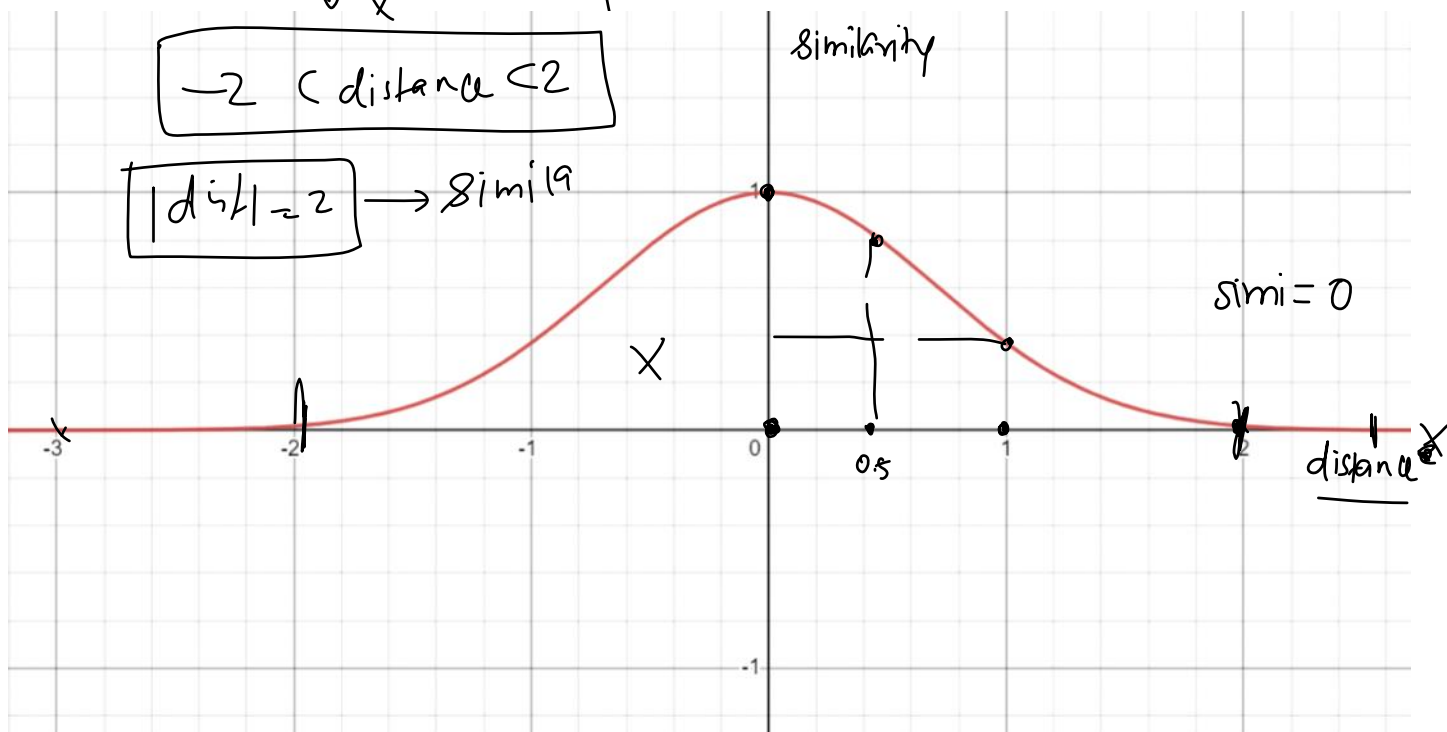
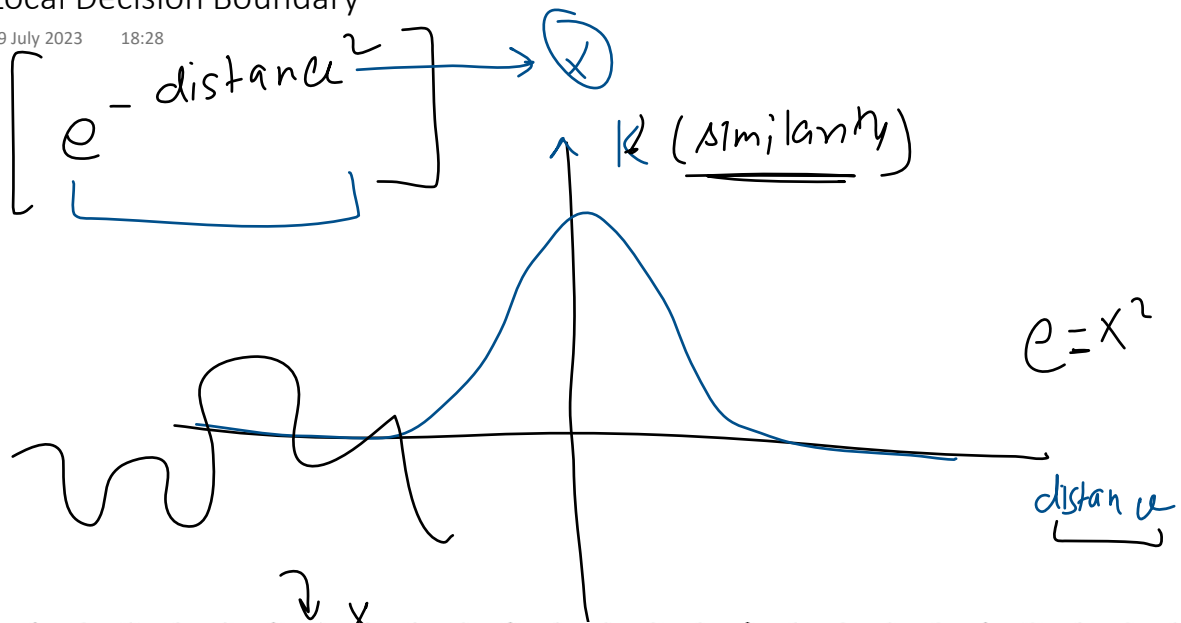
- **Non-linear Transformations:** The RBF kernel enables the use of non-linear transformations, which can map the original feature space to a higher-dimensional space where the data becomes linearly separable. This is particularly useful for problems where the decision boundary is not linear.
- **Local Decisions:** Unlike some other kernels, the RBF kernel makes "local" decisions. That is, the effect of each data point is limited to a certain region around that point. This can make the model more robust to outliers and create complex decision boundaries.
- **Flexibility:** The RBF kernel has a parameter γ (related to the standard deviation of the Gaussian distribution) that determines the complexity of the decision boundary. By tuning this parameter, we can adjust the trade-off between bias and variance, allowing for a flexible range of decision boundaries.
- **Universal Approximation Property:** The RBF kernel has a property known as the "universal approximation" property, meaning it can approximate any continuous function to a certain degree of accuracy given enough data points. This makes it highly versatile and capable of modelling a wide variety of relationships in data.
- **General-Purpose:** The RBF kernel does not make any strong assumptions about the data and can therefore be a good choice in many different situations, making it a versatile, general-purpose kernel.

distance

Local Decision Boundary

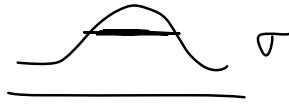
19 July 2023

18:28



Effect of Gamma

19 July 2023 11:11



The parameter γ in the Radial Basis Function (RBF) kernel of a Support Vector Machine (SVM) is a hyperparameter that determines the spread of the kernel and therefore the decision region.

The effect of γ can be summarized as follows:

- If γ is too large, the exponential will decay very quickly, which means that each data point will only have an influence in its immediate vicinity. The result is a more complex decision boundary, which might overfit the training data.
- If γ is too small, the exponential will decay slowly, which means that each data point will have a wide range of influence. The decision boundary will therefore be smoother and more simplistic, which might underfit the training data.

In a sense, γ in the RBF kernel plays a role similar to that of the inverse of the regularization parameter: it controls the trade-off between bias (underfitting) and variance (overfitting). High γ values can lead to high variance (overfitting) due to more flexibility in shaping the decision boundary, while low γ values can lead to high bias (underfitting) due to a more rigid, simplistic decision boundary.

Tuning the γ parameter using cross-validation or a similar technique is typically a crucial step when training SVMs with an RBF kernel.

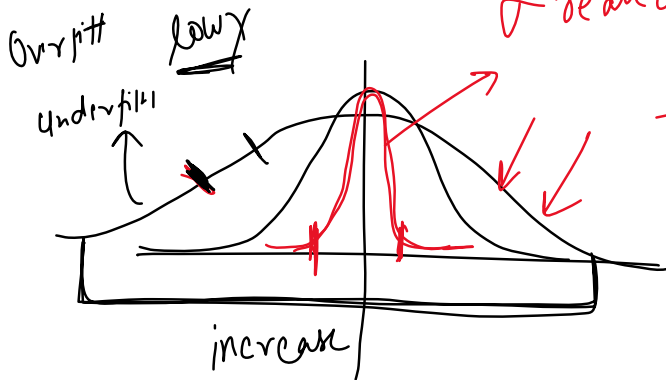
$$e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

$$\boxed{\frac{1}{2\sigma^2} = \gamma}$$

$$\sigma^2 = 1$$

$$\sigma = 10 \quad \sigma^2 = 100$$

$$\sigma = 0.1 \quad \sigma^2 = 0.01$$



Bias variance tradeoff

σ reduce $\rightarrow 0.1$

$\sigma \uparrow$ range \uparrow accuracy \uparrow
locality \uparrow

$\gamma \rightarrow$ hyperparameter

$\gamma \downarrow \rightarrow$ local \uparrow

$\gamma \uparrow \rightarrow$ local \downarrow

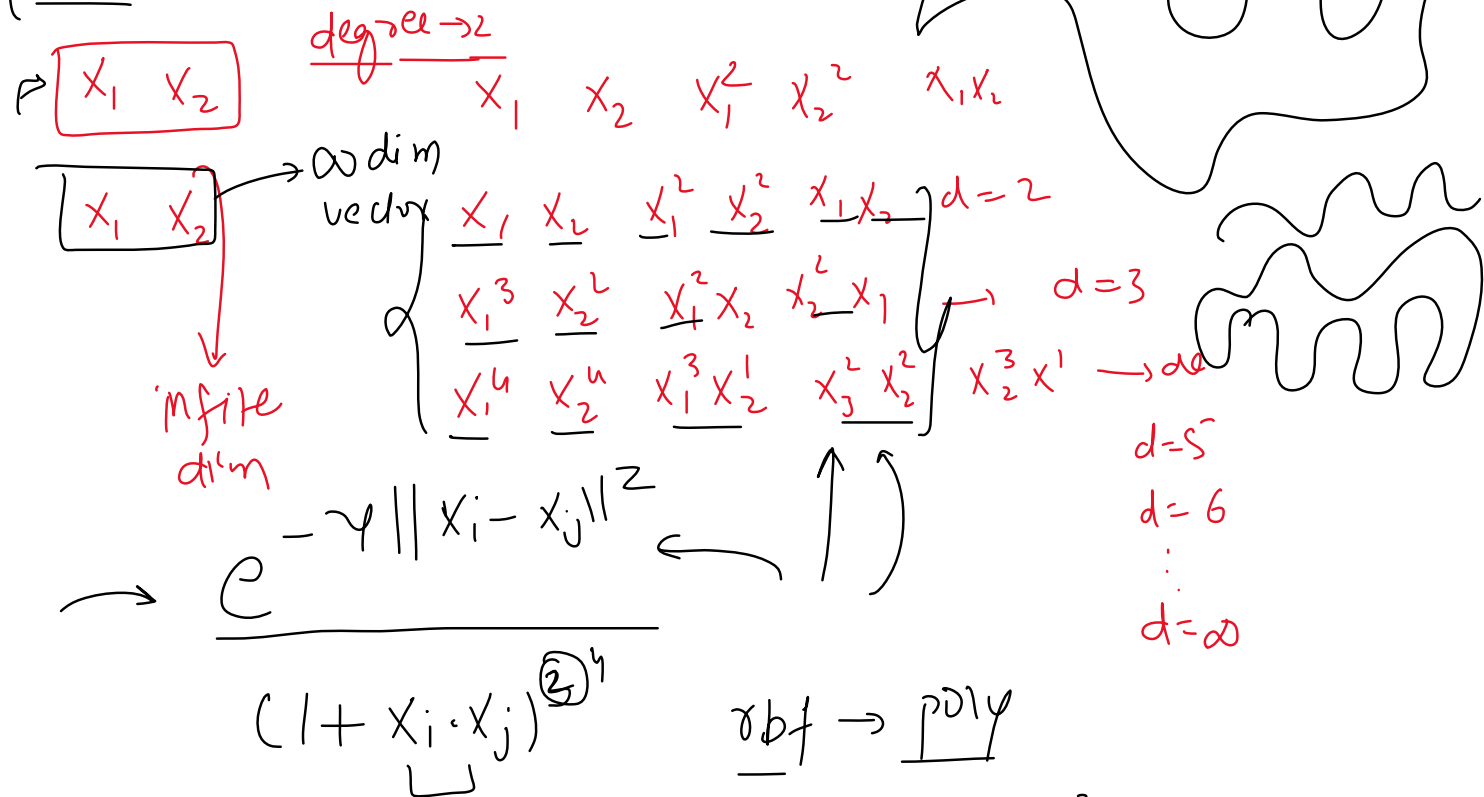
$$\sigma \propto \frac{1}{\gamma}$$

$\sigma \downarrow \gamma \uparrow \rightarrow$ overfitting
 $\sigma \uparrow \gamma \downarrow \rightarrow$ underfitting

Relationship Between RBF and Polynomial Kernel

19 July 2023 14:36

Infinite Dimensional Mapping: The RBF kernel implicitly maps input data to an infinite-dimensional feature space, which allows for even greater flexibility in forming decision boundaries



$$K(x_i, x_j) = e^{-\frac{\|x_i - x_j\|^2}{2\sigma^2}}$$

$$\sigma = 1 \quad \sigma^2 = 1$$

$$= e^{-\frac{\|x_i - x_j\|^2}{2}}$$

$$= e^{-\frac{(x_i - x_j)^T (x_i - x_j)}{2}} = e^{-\frac{(x_i^T - x_j^T)(x_i - x_j)}{2}}$$

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

$$x_j = [x_{j1} \ x_{j2}]$$

$$x_i^T x_j =$$

$$x_j^T x_i =$$

$$= e^{-\frac{x_i^T x_i + x_j^T x_j - x_i^T x_j - x_j^T x_i}{2}}$$

$$= e^{-\frac{x_i^T x_i + x_j^T x_j - 2x_i^T x_j}{2}}$$

$$= e^{-\frac{1}{2} [x_i^T x_i + x_j^T x_j]} e^{x_i^T x_j}$$

$$= \frac{1}{2} [x_i^T x_i + x_j^T x_j - 1]$$

$$x_i^T \cdot x_j$$

$$\begin{aligned}
&= C e^{1 + x_i^T x_j - 1} \\
&= C e^{1 + x_i^T x_j} e^{-1} \\
&= C' e^{1 + x_i^T x_j} \rightarrow = C' \sum_{k=0}^{\infty} \frac{(1 + x_i^T x_j)^k}{k!}
\end{aligned}$$

$C' \sum_{k=0}^{\infty} \frac{K_{\text{poly}}(x_i, x_j)^k}{k!}$

$$e^x = \sum_{k=0}^{\infty} \frac{x^k}{k!}$$

$$C' \left[1 + \frac{(1 + x_i^T x_j)}{1!} + \frac{(1 + x_i^T x_j)^2}{2!} + \frac{(1 + x_i^T x_j)^3}{3!} + \dots \right]$$

$$\textcircled{\partial b f} = \sum_{i=2}^{\infty} \frac{K_{\text{poly}}(x_i, x_j)}{i!} + \frac{(1 + x_i^T x_j)^{\infty}}{\infty!}$$

1. **String kernels:** These are used for classifying text or sequences, where the input data is not numerical. String kernels measure the similarity between two strings. For example, a simple string kernel might count the number of common substrings between two strings.
2. **Chi-square kernel:** This kernel is often used in computer vision problems, especially for histogram comparison. It's defined as $K(x, y) = \exp(-\gamma \chi^2(x, y))$, where $\chi^2(x, y)$ is the chi-square distance between the histograms x and y .
3. **Intersection kernel:** This is another kernel commonly used in computer vision, which computes the intersection between two histograms (or generally non-negative feature vectors).
4. **Hellinger's kernel:** Hellinger's kernel, or Bhattacharyya kernel, is used for comparing probability distributions and is popular in image recognition tasks.
5. **Radial basis function network (RBFN) kernels:** These are similar to the standard RBF kernel, but the centers and widths of the RBFs are learned from the data, rather than being fixed a priori.
6. **Spectral kernels:** These kernels use spectral analysis techniques to compare data points. They can be particularly useful for dealing with cyclic or periodic data.

Doubts

20 July 2023 15:42

1. What is complementary Slackness
2. Why did min change to max?
3. Observations
 - a. Effect on high dimensional data
4. Prediction is done?