

(SVC + Kernels) = SVM Non-linear data

Thursday 2 May 2024 3:24 PM

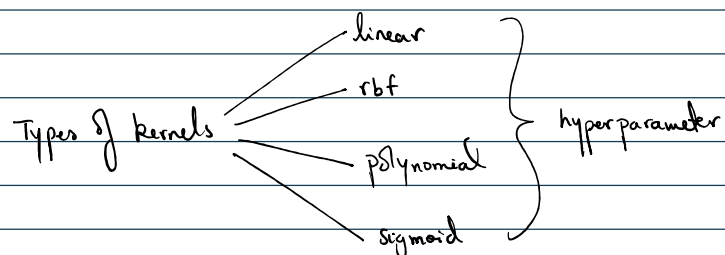
Kernel Trick → Why is it called a trick & not a transformation?

1) input $\xrightarrow{\text{kernel}}$ higher dimension

2) kernels are designed in such a way that the data is linearly separable in higher dimensions.

3) SVM

4) project down to lower dimension



You
in svm why do we call it kernel trick and not kernel transformation? explain in simple terms

ChatGPT
Sure, let's break it down in simple terms.

Imagine you have a bunch of data points that you want to classify into two groups, like whether an email is spam or not based on certain features. Now, in traditional Support Vector Machines (SVM), you try to find a straight line (or a hyperplane in higher dimensions) that separates these two groups as best as possible.

But what if your data isn't nicely separated by a straight line? What if it's all twisted and tangled up? This is where the "kernel trick" comes in.

Instead of trying to untangle your data by transforming it into a higher-dimensional space directly (which could be very computationally expensive), the kernel trick allows SVM to find that separating boundary as if it were working in a higher-dimensional space, without actually having to compute all those extra dimensions.

So, in a sense, it's like performing a transformation on your data without actually doing all the heavy lifting of computing the transformation explicitly. That's why it's called a "trick" rather than a "transformation". It makes SVM flexible enough to handle complex data without bogging down in calculations.

↓

Mathematics of SVM

1) Solving constraint optimization problems

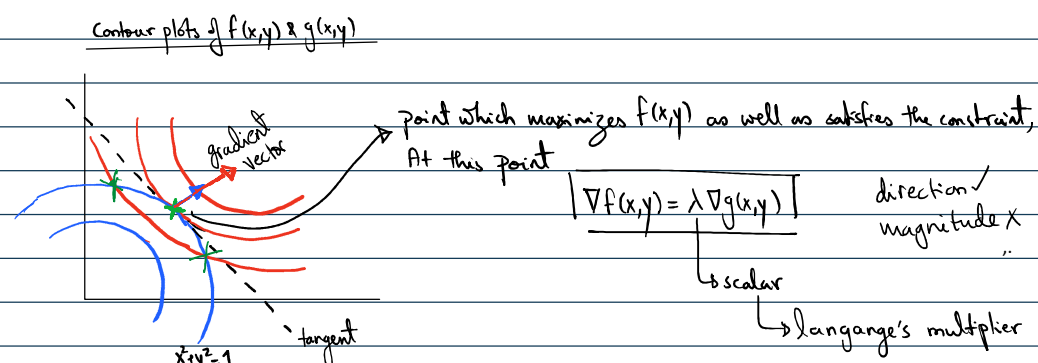
e.g. $\underset{x,y}{\text{argmax}} x^2y$ s.t. $x^2+y^2=1$

→ The gradient of a function at a point is a vector that points in the direction of the steepest ascent/maximum change of the function at that point. The magnitude of the gradient vector is equal to the rate of increase of the function in that direction.

→ Contour lines of a function are curves that connect points where the function has the same value.

- * The gradient at a point is \perp to the contour line passing through that point.
- * The gradient points in the direction where the function increases most rapidly.
- * The magnitude of the gradient indicates how steeply the function is increasing.

let $f(x,y)=x^2y$ & $g(x,y)=x^2+y^2$



$$\nabla f(x,y) = \begin{bmatrix} \frac{\partial f}{\partial x} \\ \frac{\partial f}{\partial y} \end{bmatrix} = \begin{bmatrix} 2xy \\ x^2 \end{bmatrix} \quad \nabla g(x,y) = \begin{bmatrix} \frac{\partial g}{\partial x} \\ \frac{\partial g}{\partial y} \end{bmatrix} = \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$\begin{bmatrix} 2xy \\ x^2 \end{bmatrix} = \lambda \begin{bmatrix} 2x \\ 2y \end{bmatrix}$$

$$2xy = \lambda 2x \quad \text{--- ①}$$

$$x^2 = \lambda 2y \quad \text{--- ②}$$

$$x^2 + y^2 = 1 \quad \text{--- ③}$$

$$y = \pm \frac{1}{\sqrt{3}} \quad x = \pm \sqrt{\frac{2}{3}}$$

4 set of points
which satisfy
the constraint

→ Out of these 4 points $\left(\sqrt{\frac{2}{3}}, \frac{1}{\sqrt{3}}\right)$ & $\left(-\sqrt{\frac{2}{3}}, \frac{1}{\sqrt{3}}\right)$
result in max value for $f(x,y)$.

2) Lagrange multiplier

$\underset{x,y}{\operatorname{argmax}} x^2 y \quad \text{s.t.} \quad x^2 + y^2 = 1$
↓ constraint optimization problem
↓ optimization problem

$$L(x,y,\lambda) = \underset{x,y}{\operatorname{argmax}} f(x,y) - \lambda(g(x,y) - 1)$$

$$\begin{array}{ccc} \frac{\partial L}{\partial x} = 0 & \frac{\partial L}{\partial y} = 0 & \frac{\partial L}{\partial \lambda} = 0 \\ \downarrow & \downarrow & \downarrow \\ 2xy = \lambda 2x & x^2 = \lambda 2y & x^2 + y^2 = 1 \end{array}$$

3) SVM in n-Dimensions

equation of hyperplane: $w^T x + b = 0$

$$\underset{w,b}{\operatorname{argmin}} \left(\frac{\|w\|^2}{2} + \frac{1}{n} \sum_{i=1}^n \xi_i \right) \quad \text{s.t.} \quad \left\{ y_i(w^T x + b) \geq 1 - \xi_i \right\}$$

↑ soft margin formulation in n-D

4) Constrained optimization problems with inequality

$$\text{e.g.} \quad \min_x f(x) = x^2 \quad \text{s.t.} \quad x - 1 \leq 0$$

→ KKT conditions

→ They generalize the method of Lagrange multipliers to handle inequality constraints.

→ They generalize the method of Lagrange multipliers to handle inequality constraints.

→ KKT conditions play a key role in deriving the dual problem from the primal problem.

$$\begin{array}{c} \text{Primal} \\ \text{e.g. } \min_x f(x) = x^2 \text{ s.t. } x-1 \leq 0 \end{array} \rightarrow \begin{array}{c} \text{Lagrangian} \\ L(x, \lambda) = x^2 - \lambda(x-1) \end{array}$$

Conditions:

$$1) \frac{\partial L}{\partial x} = 0 \quad \frac{\partial L}{\partial \lambda} = 0 \quad 2) x-1 \leq 0$$

$$3) \lambda \geq 0 \quad 4) \lambda(x-1) = 0$$

x & λ which satisfy these 4 conditions

will be the solution to our problem

5) Duality

→ optimization theory

→ The primal problem is the original optimization problem that you are trying to solve. It involves finding the minimum or maximum of a particular objective function, subject to certain constraints.

→ The dual problem is a related optimization problem that is derived from the primal problem. It provides a lower or upper bound on the solution to the primal problem.
 → easier to solve.

→ if strong duality holds, then solving the dual problem can directly give the solution to the primal problem.

6) SVM dual problem derivation (Hard margin SVM)

Primal form

training data norm

$$\arg \min_{w, b} \frac{\|w\|^2}{2} \text{ s.t. } y_i(w \cdot x_i + b) \geq 1 \quad \forall i$$

↳ for every row 1 constraint

Lagrangian

$$L(w, b, \alpha) = \frac{\|w\|^2}{2} - \sum_{i=1}^n \alpha_i [y_i(w \cdot x_i + b) - 1]$$

$$L(w, b, \alpha) = \frac{\|w\|^2}{2} - \sum_{i=1}^n \alpha_i y_i w \cdot x_i - \sum_{i=1}^n \alpha_i y_i b + \sum_{i=1}^n \alpha_i$$

$$\frac{\partial L}{\partial w} = \frac{2w}{2} - \sum_{i=1}^n \alpha_i y_i x_i = 0$$

$$\frac{\partial L}{\partial b} = - \sum_{i=1}^n \alpha_i y_i = 0$$

$$\textcircled{1} \quad w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\textcircled{2} \quad \sum_{i=1}^n \alpha_i y_i = 0$$

sub ① & ② into Lagrangian form

$$L(w, b, \alpha) = \frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i x_i \right) \left(\sum_{j=1}^n \alpha_j y_j x_j \right) - \sum_{i=1}^n \alpha_i y_i x_i \left(\sum_{j=1}^n \alpha_j y_j x_j \right) + \sum_{i=1}^n \alpha_i$$

$$= -\frac{1}{2} \left(\sum_{i=1}^n \alpha_i y_i x_i \right) \left(\sum_{j=1}^n \alpha_j y_j x_j \right) + \sum_{i=1}^n \alpha_i$$

primal dual
min \rightarrow max
max \rightarrow min

$$\underset{\alpha_i}{\operatorname{argmax}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (x_i \cdot x_j) \quad \text{subject to } \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$$

\hookrightarrow dual form
easier to solve kernel friendly

1) $\alpha > 0$ only for support vectors; $\alpha = 0$ for non-support vectors

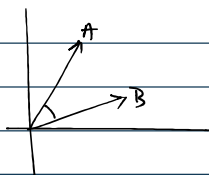
in dual form only support vectors are considered

7) Similarity

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

if $\|A\| = 1$ & $\|B\| = 1$

$$= A \cdot B \quad \hookrightarrow \text{dot product of } A \text{ \& } B$$



$$\underset{\alpha_i}{\operatorname{argmax}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j \underbrace{(x_i \cdot x_j)}_{\text{similarity}} \quad \text{subject to } \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$$

\hookrightarrow maximize the similarity of support vectors based on their signs.

8) Kernel SVM

$$\underset{\alpha_i}{\operatorname{argmax}} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j K(x_i \cdot x_j) \quad \text{subject to } \alpha_i \geq 0, \sum_{i=1}^n \alpha_i y_i = 0$$

$K(x_i, x_j) \rightarrow$ kernel \rightarrow similarity b/w x_i & x_j

if kernel function is $x_i \cdot x_j \rightarrow$ linear SVM

\rightarrow if in the similarity part of the dual form we use some other form of mathematical function to calculate the similarity we call it kernel SVM.

9) Polynomial Kernel

$K(x_i, x_j) = (r + x_i \cdot x_j)^d$ \rightarrow expansion of this will contain polynomial terms which mimics transforming our data into higher dimension without actually applying feature transformation. This way, we save memory space.

\rightarrow While the squared terms cater only to circular kind of shapes. But with the help of the interaction term we will be able to map a wider variety of shapes.

10) RBF kernel

$$K(x_i, x_j) = e^{-\gamma \|x_i - x_j\|^2} \quad \text{where } \gamma = 1/2\sigma^2$$

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

\uparrow hyperparameter

$\gamma \uparrow \sigma \downarrow$ similarity region \downarrow (low bias high variance) overfitting

$\gamma \downarrow \sigma \uparrow$ similarity region \uparrow (high bias low variance) underfitting

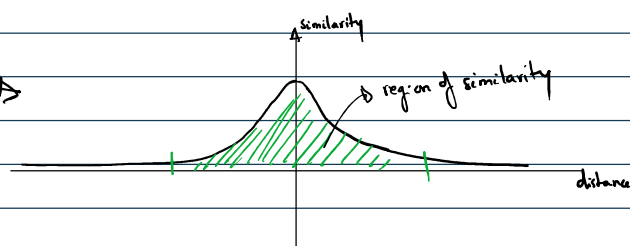
} bias variance trade off

1) Non-linear transformations

2) Local decisions

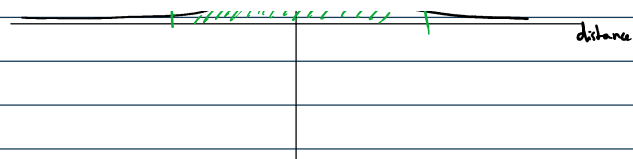
3) Flexibility

4) Universal Approximation Property



4) Universal Approximation Property

5) general purpose



* features must be normalized (b/w -1 & +1) ??

→ Infinite dimensional mapping: The RBF kernel implicitly maps input data to a so dimensional feature space, which allows for even greater flexibility in forming decision boundaries.

→ The reason lies in e^x being expressed as infinite series.

→ Similarity with KNN

$$K \text{ similar to } \sigma \propto \frac{1}{\lambda}$$

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

→ Sum of infinite polynomial kernels.

11) Custom kernels are possible