

SVM

1. THE DUAL

• Primal

$$\min_{w, b} \frac{1}{2} \|w\|^2$$

data points
↓

$$\text{constraints} \rightarrow y_i (w^T x_i - b) \geq 1 \quad \forall i = 1, \dots, N$$

• Lagrangian

↪ α_i are Lagrangian multipliers
of each data point

$$\rightarrow \mathcal{L}(w, b, \alpha) = \frac{1}{2} w^T w - \sum_{i=1}^N \alpha_i ((w^T x_i - b) - 1)$$

• Dual

$$\max_{\alpha_i \geq 0} \left[\min_{w, b} \mathcal{L}(w, b, \alpha) \right]$$

↑
fixed

• (Dual)

$$\max_{\alpha_i \geq 0} \left[\min_{w, b} \mathcal{L}(w, b, \alpha) \right]$$

↳ it seems solving the min problem first to find w, b with a fixed α_i ; then find an α that maximize everything

fixed

? not sure

↳ then in stationarity constraints with partial derivatives of w not $b = 0$

$$1 \rightarrow \frac{\partial \mathcal{L}}{\partial w} = 0 \Rightarrow w = \sum_{i=1}^N \alpha_i y_i x_i$$

~> But what about α_i ?

The only vectors contributing to the margin are support vectors not they are the only vectors that can contribute to the definition of w .

This means that for NON-supporting vectors x_i : $\alpha_i = 0$. if not they would contribute to w and that's not true.

$$2 \rightarrow \frac{\partial \mathcal{L}}{\partial b} = 0 \Rightarrow 0 = \sum_{i=1}^N \alpha_i y_i$$

$$0 = \alpha^T y$$

vector dot-product formulation of the previous sum

• We want to define the dual problem only in terms of $\max_{\alpha_i \geq 0}$

so we substitute w and b from stationary constraints to get:

$$\max_{\alpha_i \geq 0} \left[\sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_i \sum_j \alpha_i \alpha_j y_i y_j x_i^T x_j \right]$$

After getting α , we can obtain w and b using stationary constraints

This is simple because NON-SVs have $\alpha_i = 0$ and it means we need to sum only SV because all the others are zero in the sum.

The simplified formula is (changing the sign to get a MIN problem)

$$\min_{\alpha_i \geq 0} \left[\frac{1}{2} \sum_{\substack{i,j \\ \uparrow \\ \text{SVs}}}^N \alpha_i \alpha_j y_i y_j x_i^T x_j - \sum_{\substack{i \\ \uparrow \\ \text{SVs}}}^N \alpha_i \right]$$

Why Dool?

1. useful for KERNELS (see later)

2. with high-dimensional data,

when $p \gg N$ (attributes \gg record)

for example with images,

in the primal we have to use the whole dataset

so $N \cdot p$ operations, while in the dual

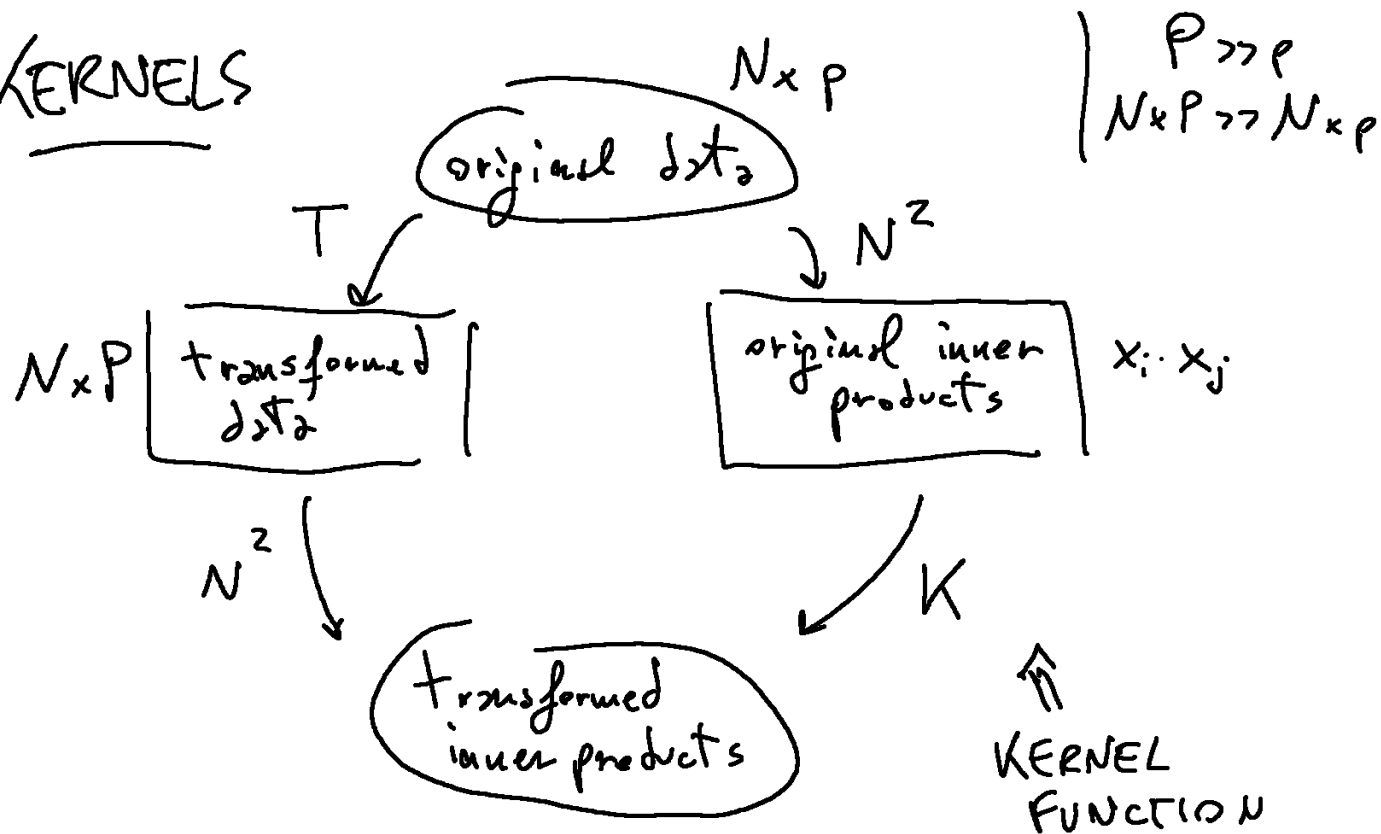
we need to compute only $\langle x_i^T, x \rangle$ which is only N^2 .

Then if $p \gg N$ and the dual is better

$$Np \gg N^2$$

② I didn't get why we used Np in the primal

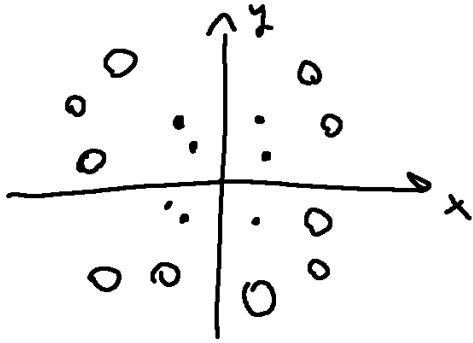
KERNELS



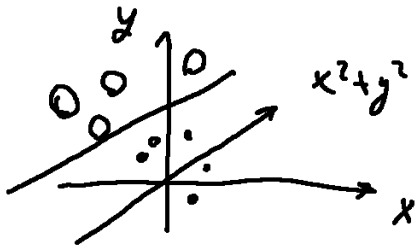
- What is the benefit of kernel functions?
We never have to send data to higher dimensional space.
- In the context of SVM, a kernel function is any function that only uses the inner products of the original data but it's able to use those inner products and send them into transformed inner products without visiting the higher dimensional space.

- Why do we need to transform data into higher dimensional data?

↳ Because if we had points like this,



we cannot use an hyperplane to separate the dataset and therefore we need to engineer another dimension.



In this case, it will be x^2+y^2 , (which in two-dimensions is a circle).

But this means transforming

$$x_i \mapsto T(x_i)$$

- Why do we only need inner products?

↳ Because we are solving the Dual problem of SVM

PRIMAL

$$\begin{cases} \min_{w, b} \frac{1}{2} \|w\|^2 \\ y_i (w^T x_i - b) \geq 1 \end{cases}$$

↔ DUAL

$$\begin{cases} \min_{\alpha_i \geq 0} \sum_i \sum_j \alpha_i \alpha_j y_i y_j \overbrace{x_i^T x_j} \\ \alpha^T y = 0 \leftarrow ? \end{cases}$$

$$x_i \cdot x_j \mapsto T(x_i \cdot x_j)$$

or with Kernel Functions

$$\mapsto K(x_i \cdot x_j)$$