

Kernels

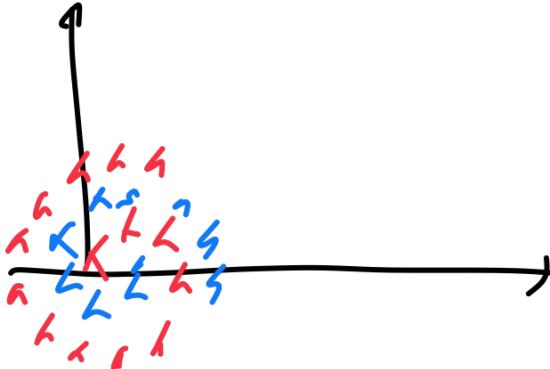
let us consider a long feature vector

$$\tilde{\phi}(x^{(i)}) = (1, x^{(i)}, (x^{(i)})^2, (x^{(i)})^3, \dots, (x^{(i)})^d)$$

We want to minimize the loss

$$l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \tilde{\phi}(x^{(i)}))^2$$

$$\beta^{(k+1)} = \beta^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \tilde{\phi}(x^{(i)})) \cdot \tilde{\phi}(x^{(i)})$$



$$x = (x_1, \dots, x_D)$$

$$D = 1000$$

$$\tilde{\phi} = 10^9$$

$$\tilde{\phi}(x) = \begin{bmatrix} x_1^1 \\ \vdots \\ x_D^1 \\ x_1^2 \\ \vdots \\ x_D^2 \\ x_1^3 \\ \vdots \\ x_1 x_2 x_3 \end{bmatrix}$$

$\rightarrow 1 + D + D^2 + D^3$

Idea: write $\beta^{(k)}$ as $\beta^{(k)} = \sum_{j=1}^N \phi(x^{(j)}) \lambda_j$

start $\beta^{(0)} = \vec{0}$

$$\beta^{(k+1)} = \beta^{(k)} + \gamma/N \sum_{i=1}^N (t^{(i)} - \beta^\top \tilde{\phi}(x^{(i)})) \cdot \tilde{\phi}^{(i)}$$

$$= \sum_{j=1}^N \underbrace{\phi(x^{(j)}) \lambda_j}_\downarrow + \frac{\gamma}{N} \sum_{j=1}^N (t^{(j)} - \sum_{j=1}^N \lambda_j \phi(x^{(j)})^\top \underbrace{\phi(x^{(j)})}_{\tilde{\phi}^{(j)}})$$

$$= \sum_{j=1}^N \underbrace{\phi(x^{(j)}) (\lambda_j)}_{\text{green bracket}} + \frac{\gamma}{N} \left(t^{(j)} - \sum_{i=1}^N \lambda_i \phi(x^{(i)})^\top \phi(x^{(j)}) \right)$$

$$\lambda_j^{(k+1)} = \lambda_j^{(k)} + \frac{\gamma}{N} \left(t^{(j)} - \sum_{i=1}^N \lambda_i \phi(x^{(i)})^\top \phi(x^{(j)}) \right)$$

Kernels

$$y(x) = \beta^T \phi(x)$$

$$\sum_{i=1}^{K'} \lambda_i^{(k)} \phi(x^{(i)})$$

Let $\phi(x^{(i)})$ to denote a feature vector (e.g. polynomial features)

$$\beta^{(k+1)} \leftarrow \beta^{(k)} + \gamma \sum_{i=1}^N (\ell^{(i)} - \beta^T \phi(x^{(i)})) \cdot \phi(x^{(i)})$$

$\phi(x^{(i)})$ depends on D (dimension) \rightarrow can't expect to store when D is large

How about storing β as $\sum_{i=1}^N \lambda_i^{(0)} \phi(x^{(i)})$

$$\sum_{i=1}^N \lambda_i^{(0)} \phi(x^{(i)})$$

Always possible to make such a decomposition for $\beta^{(0)} = 0$

$$\beta^{(k+1)} = \sum_{i=1}^N \lambda_i \phi(x^{(i)}) + \gamma \sum_{i=1}^N (t^{(i)} - \sum_{j=1}^N \lambda_j \phi(x^{(j)})^T \phi(x^{(i)})) \phi(x^{(i)})$$

$$= \sum_{i=1}^N \phi(x^{(i)}) \left(\lambda_i + \gamma \left(t^{(i)} - \sum_{j=1}^N \lambda_j \phi(x^{(j)})^T \phi(x^{(i)}) \right) \right)$$

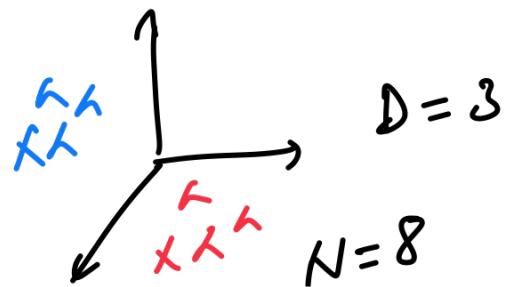
$$\beta^{(k)} = \sum_{i=1}^N \lambda_i^{(k)} \phi(x^{(i)}) \rightarrow \lambda^{(k+1)}$$

$$\beta^{(k+1)} = \sum_{i=1}^N \lambda_i^{(k+1)} \phi(x^{(i)})$$

$$\lambda_i^{(k+1)} = \lambda_i^{(k)} + \gamma \left(t^{(i)} - \sum_{j=1}^N \lambda_j^{(k)} \underbrace{\phi(x^{(j)})^T \phi(x^{(i)})}_{(*)} \right)$$

General idea: iterate over the $\lambda_i^{(k)}$ instead $\beta_j^{(k)}$

→ Storage goes from $O(D)$ to N



$$K(x^{(i)}, x^{(j)}) = k(i, j) = \phi(x^{(i)})^T \phi(x^{(j)})$$

→ K is of size N^2

$$\begin{aligned} \text{Why is it interesting to have } & \langle \phi(x^{(i)}), \phi(x^{(j)}) \rangle \\ & = \phi(x^{(i)})^T \phi(x^{(j)}) \end{aligned}$$

Consider polynomial features of max degree 3

$$\begin{aligned} \phi(x) &= [1, x_1, x_1^2, x_2, x_2^2, \dots, x_1^3, x_2^3, \dots, x_D^3] \\ \phi(z) &= [1, z_1, z_1^2, z_2, z_2^2, \dots, z_1^3, z_2^3, \dots, z_D^3] \end{aligned}$$

$$\begin{aligned} \phi(x)^T \phi(z) &= 1 + \sum x_i z_i + \sum x_i x_j z_j z_j + \sum x_i x_j x_k z_i z_j z_k \\ &= 1 + (x^T z) + (x^T z)^2 + (x^T z)^3 \end{aligned}$$

$$\gamma_i^{(k+1)} \leftarrow \gamma_i^{(k)} + \gamma (t^{(i)} - \sum_{j=1}^N \gamma_j^{(k)} \phi(x^{(i)})^T \phi(x^{(j)}))$$

$\leftarrow \gamma_i^{(k)} + \gamma (t^{(i)} - \sum_{j=1}^N \gamma_j^{(k)} K(i, j))$

If we want to focus on K instead of ϕ

What are the "good" K matrices that we could use
to learn the properties of our data?

$K \rightarrow$ Square Symmetric

→ positive semi-definite (psd)

To see why K has to be psd

take $z \in \mathbb{R}^N$

$$\text{if } \exists \phi(x^{c_i}), \phi(x^{g_j}) \text{ s.t. } K(i, j) \\ = \phi(x^{c_i})^\top \phi(x^{g_j})$$

then

$$\underbrace{z^\top K z}_{=} = \sum_{i,j} z_i k_{ij} z_j$$

$$= \sum_{i,j} z_i \phi(x^{c_i})^\top \phi(x^{g_j}) z_j$$

$$= \sum_{i,j,k} z_i \phi_k(x^{c_i}) \phi_k(x^{g_j}) z_j$$

$$= \sum_k \left(\sum_i \underbrace{z_i \phi_k(x^{c_i})}_{} \right) \left(\sum_j \underbrace{z_j \phi_k(x^{g_j})}_{} \right)$$

$$= \sum_k \left(\sum_i z_i \phi_k(x^{(i)}) \right)^2$$

$\Rightarrow z^T K z \geq 0 \Leftrightarrow K$ positive semi-definite matrix

\Rightarrow Theorem (Mercer) For K to be a valid kernel (Mercer kernel) it is necessary and sufficient to have K symmetric and positive semi-definite

$$k(i,j) = \frac{1}{\sigma} \exp \left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma^2} \right)$$

λ_i we get after
our
iterations

How do we recover $y(x)$ from λ_i^* ?

$$\begin{aligned}
 y(x) &= \beta^T \phi(x) = \left(\sum_{i=1}^N \lambda_i^* \phi(x^{(i)}) \right)^T \phi(x) \\
 &= \sum_{i=1}^N \lambda_i^* \phi(x^{(i)})^T \phi(x) \\
 &= \underbrace{\lambda^T}_{\lambda} K(x^{(i)}, x)
 \end{aligned}$$

as an example if $k(i,j)$ can be generate from a function $k(x,y)$ of the original data,

We can get the expression of our model $y(x)$ directly from that function $k(x, y)$

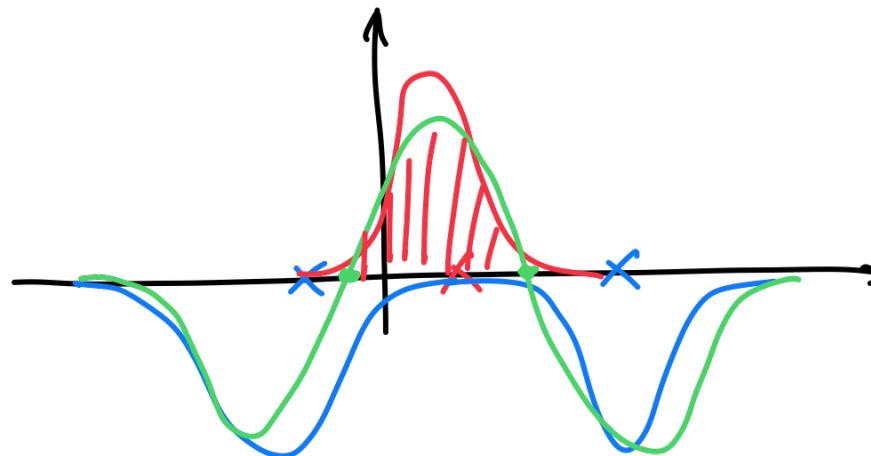
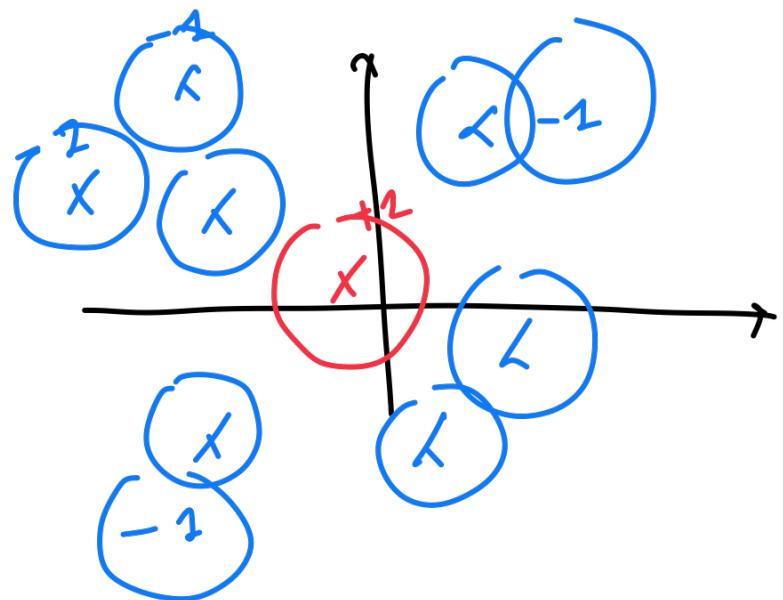
e.g take $k(\underline{x}, \underline{y}) = \frac{1}{\sigma} \exp\left(-\frac{\|\underline{x} - \underline{y}\|^2}{\sigma^2}\right)$

To learn $y(x)$ we first need the N^2

$$K(i, j) = \frac{1}{\sigma} \exp\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma^2}\right) \text{ for every } x^{(i)}, x^{(j)}$$

then learn the $\lambda_i^{(k)}$

Finally you get your model as $y(x) = \sum_{i=1}^N \lambda_i^* \exp\left(-\frac{\|x^{(i)} - x\|^2}{\sigma^2}\right)$



Kernels

General idea → switch from an approach based on feature vectors (here on the dimension D)

to an approach depending on the number of training points N

$$\rightarrow l(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \phi(x^{(i)}))^2$$

$$\beta^{(k+1)} \leftarrow \beta^{(k)} + \gamma \frac{1}{N} \sum_{i=1}^N (t^{(i)} - \beta^T \phi(x^{(i)})) \phi(x^{(i)})$$

general idea $\beta^{(k)} = \sum \lambda_i^{(k)} \phi(x^{(i)})$

$$\begin{aligned} \beta^{(k+1)} &\leftarrow \sum_{i=1}^N \lambda_i^{(k)} \phi(x^{(i)}) + \frac{1}{N} \sum_{i=1}^N \left(t^{(i)} - \sum_{j=1}^N \lambda_j^{(k)} \phi(x^{(j)})^\top \phi(x^{(i)}) \right) \phi(x^{(i)}) \\ &\leftarrow \sum_{i=1}^N \left(\lambda_i^{(k)} + \frac{1}{N} \sum_{j=1}^N \left(t^{(j)} - \sum_{j=1}^N \lambda_j^{(k)} \phi(x^{(j)})^\top \phi(x^{(i)}) \right) \right) \phi(x^{(i)}) \end{aligned}$$

$$\lambda_i^{(k+1)} \leftarrow \lambda_i^{(k)} + \frac{1}{N} \underbrace{\left(t^{(i)} - \sum_{j=1}^N \lambda_j^{(k)} \phi(x^{(j)})^\top \phi(x^{(i)}) \right)}_{K(i,j) \lambda_j}$$

We then introduce $K(x^i, x^j)$ to store the scalar products between the feature vectors.

We can then study the class of K matrices for which there exist a decomposition in feature vectors

→ Mercer's Kernel : Positive semi definite K
(Symmetric)

e.g $K(x, y) = \exp\left(-\frac{\|x-y\|^2}{\sigma}\right)$

Claim $K(x, y) = \exp\left(-\frac{\|x - y\|^2}{\sigma^2}\right)$ is a valid kernel
 (meaning
 $\exists \phi(x^{(i)}) \phi(x^{(j)})$)

such that

$$K(x^{(i)}, x^{(j)}) = \phi(x^{(i)})^\top \phi(x^{(j)})$$

Proof

$$\begin{aligned} \exp\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma^2}\right) &= e^{-\frac{1}{\sigma^2} \|x^{(i)}\|^2 - \frac{1}{\sigma^2} \|x^{(j)}\|^2} \\ &= e^{-\frac{1}{\sigma^2} \|x^{(i)}\|^2} e^{-\frac{1}{\sigma^2} \|x^{(j)}\|^2} \end{aligned}$$

$$e^{\frac{2}{\sigma^2} x^{(i)^\top} x^{(j)}}$$

Taylor

$$\begin{aligned} &\left(1 + \frac{2}{\sigma^2} \sum_{k=1}^D x_k^{(i)} x_k^{(j)} + \left(\frac{2}{\sigma^2}\right)^2 \frac{\left(\sum_k x_k^{(i)} x_k^{(j)}\right)^2}{2!} \right. \\ &\quad \left. + \frac{1}{3!} \left(\frac{2}{\sigma^2}\right)^3 \left(\sum_k x_k^{(i)} x_k^{(j)}\right)^3 + \dots\right) \end{aligned}$$

$$\exp\left(-\frac{\|x^{(i)} - x^{(j)}\|^2}{\sigma}\right) = e^{-\frac{1}{\sigma} \|x^{(i)}\|^2} e^{-\frac{1}{\sigma} \|x^{(j)}\|^2}$$

$\phi(x^{(i)})^T$
 $\phi(x^{(j)})$

* $(1, \sqrt{\frac{2}{\sigma}} x^{(i)}, \sqrt{\left(\frac{2}{\sigma}\right)^2 \frac{1}{2!}} (x_k^{(i)} x_{k+1}^{(i)})_{k, l}, \dots)$

$\sqrt{\left(\frac{2}{\sigma}\right)^3 \frac{1}{3!}} (x_k^{(i)} x_{k+1}^{(i)} x_{k+2}^{(i)})_{k, l, m, \dots}$

$$\begin{aligned} & \left(\sum_{k=1}^D x_k^{(i)} x_k^{(j)} \right)^2 \\ &= \sum_{k=1}^D \sum_{l=1}^D x_k^{(i)} x_l^{(i)} x_k^{(j)} x_l^{(j)} = \end{aligned}$$

$(1, \sqrt{\frac{2}{\sigma}} x^{(j)}, \sqrt{\left(\frac{2}{\sigma}\right)^3 \frac{1}{2!}} (x_k^{(j)} x_{k+1}^{(j)})_{k, l}, \dots)$

Formal statement of Mercer:

X measure space, we will $k: L^2(X \times X) \rightarrow \mathbb{R}$ is a valid kernel iff there exist some feature map $\varphi: X \rightarrow H$ H separable Hilbert space

$$k(x, x') = \langle \varphi(x), \varphi(x') \rangle_H$$

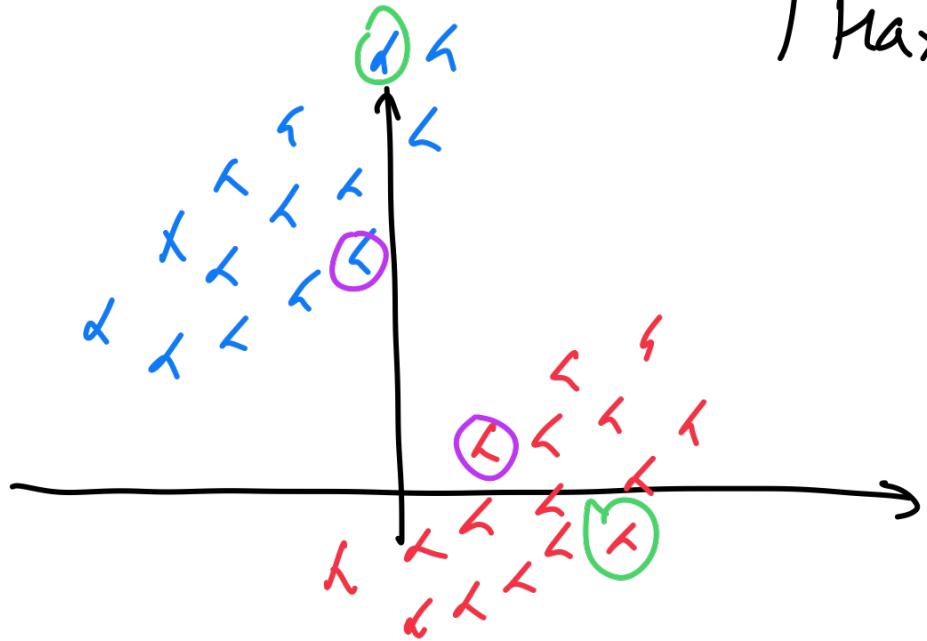
Given a kernel k we say that k satisfies Mercer's condition if

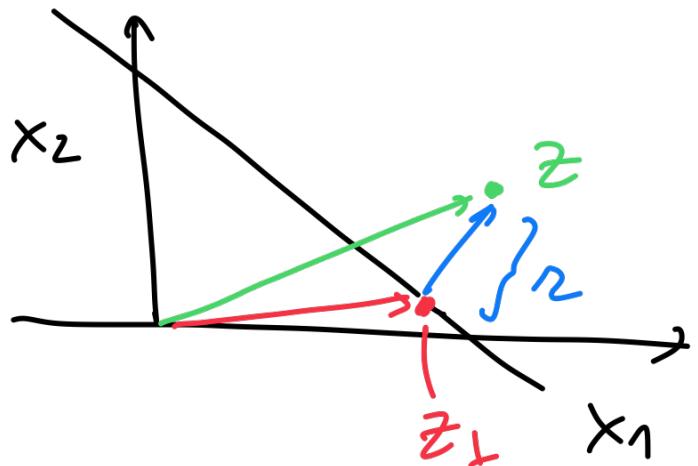
$$\int_{X \times X} k(x, x') f(x) f(x') dx dx' \geq 0 \text{ for all } f \in L^2(X)$$

Mercer theorem

k valid kernel iff k is Mercer

Support Vector Machines / Kernel Vector Machines
/ Sparse Vector Machines
/ Maximum Margin classifier

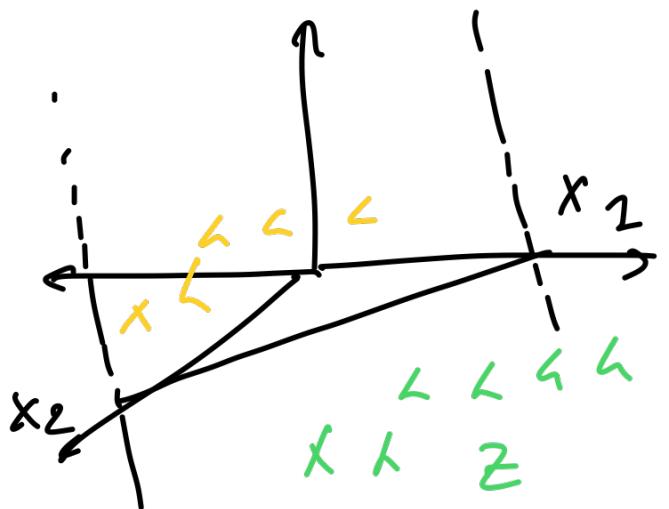




$$z = z_{\perp} + \frac{(\beta_1, \beta_2)}{\|(\beta_1, \beta_2)\|} \cdot r ?$$

Plane $\beta_1 x_1 + \beta_2 x_2 + \beta_0 = 0$

Normal vector $= (\beta_1, \beta_2)$



$$x_3 = \beta_1 x_1 + \beta_2 x_2 + \beta_0$$

$$\ell(\beta) = \frac{1}{N} \sum_{i=1}^N (t^{(i)} - (\beta_0 + \beta_1 x_1 + \beta_2 x_2))^2$$

$$x_2 = -\frac{\beta_1}{\beta_2} x_1 - \frac{\beta_0}{\beta_2}$$

$$z = z_{\perp} + \frac{(\beta_1, \beta_2)}{\|(\beta_1, \beta_2)\|} \cdot r = (z_{\perp})_2, (z_{\perp})_2 + \frac{(\beta_1, \beta_2)}{\|(\beta_1, \beta_2)\|} r$$

$\|z_1, z_2\|$

$$\beta_0 + \beta_1 (z_{\perp})_2 + \beta_2 (z_{\perp})_2 = 0$$

$$\underbrace{\beta_0 + \beta_1 z_1 + \beta_2 z_2}_{y(z)} = \underbrace{\beta_0 + \beta_1 (z_{\perp})_2 + \beta_2 (z_{\perp})_2}_{= 0} + \frac{\beta_1 \beta_1}{\|(\beta_1, \beta_2)\|} r + \frac{\beta_2 \beta_2}{\|(\beta_1, \beta_2)\|} r$$

$$= r \frac{\|(\beta_1, \beta_2)\|^2}{\|(\beta_1, \beta_2)\|} = r (\|\beta_1, \beta_2\|)$$

For any point z , the distance π of z to the plane is given by

$$\pi = \frac{y(z)}{\|(\beta_1, \beta_2)\|}$$

$$\text{dist} = \pi t(z) = \frac{y(z) \cdot t(z)}{\|(\beta_1, \beta_2)\|}$$

From π we can get $| \pi |$ by

Multiplying by $t^{(i)}$

(target of $z^{(i)}$)

provided $t^{(i)} = +1$ if

$z^{(i)}$ is

above

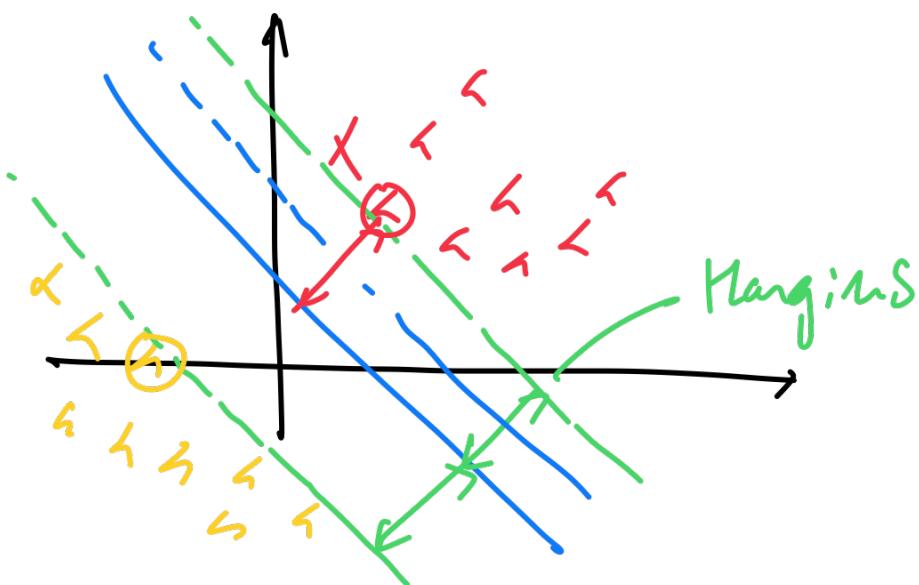
-1 if $z^{(i)}$ is

below

We can then look for the plane that maximizes its distance with respect to the closest $z^{(i)}$

$$\max_{(\beta_0, \beta_1, \beta_2)} t(z^{(i)}) = \frac{t(z^{(i)})y(z^{(i)})}{\|(\beta_1, \beta_2)\|} = \frac{t(z)(\beta_0 + \beta_1 z_1 + \beta_2 z_2)}{\|(\beta_1, \beta_2)\|}$$

± 1



ratio is invariant by rescaling of β

$$\beta \leftarrow \alpha \beta$$

$$\frac{t(z^{(i)}) (\cancel{\beta_0} + \cancel{\beta_1 z_1^{(i)}} + \cancel{\beta_2 z_2^{(i)}})}{\cancel{\|\beta_1, \beta_2\|}} \rightarrow$$

→ let us choose β such that for the closest $z^{(i)}$ we have

$$t(z^{(i)}) (\beta_0 + \beta_1 z_1^{(i)} + \beta_2 z_2^{(i)}) = 1$$

⇒ For the closest point the distance is $\frac{1}{\|\beta_1, \beta_2\|}$

⇒ For all the other points, necessarily we have

$$\text{distance} \geq \frac{1}{\|\beta_1, \beta_2\|}$$

For all the other points $t(z^{(j)}) (\beta_0 + \beta_1 z_1^{(j)} + \beta_2 z_2^{(j)}) \geq 1$

the problem then turns into

$$\max_{(\beta_0, \beta_1, \beta_2)} \frac{1}{\|(\beta_1, \beta_2)\|}$$

under the constraints

$$t^{(i)} (\beta_0 + \beta_1 z_1^{(i)} + \beta_2 z_2^{(i)}) \geq 1$$

for every i

We can further simplify the formulation into

$$\begin{array}{ll}\text{min}_{\beta} & \|\beta_1, \beta_2\|^2 \\ \text{s.t.} & t^{(i)}(\beta_0 + \beta_1 z_1^{(i)} + \beta_2 z_2^{(i)}) \geq 1\end{array}$$