# 7- PCA and SVMs

measure: distance is called measure when:   ✿ a measure is called metric, if and only
① symmetry = $d(\theta_1, \theta_2) = d(\theta_2, \theta_1)$    if it holds the triangular inequality
② self similarity = $d(\theta_1, \theta_1) = 0$       $\hookrightarrow d(\theta_1, \theta_3) \leq d(\theta_1, \theta_2) + d(\theta_2, \theta_3)$
③ positivity = $d(\theta_1, \theta_2) = 0 \Longleftrightarrow \theta_1 = \theta_2$

minkowski metric = $d(x,y) = \left( \sum_{i=1}^{p} |x_i - y_i|^q \right)^{\frac{1}{q}}$   $\left( \ell_1 = \text{manhattan} \quad \ell_2 = \text{euclidean} \quad \ell_\infty = \text{tchebychev} \atop \qquad\qquad\qquad\qquad\qquad = \max(|x_i - y_i|) \right)$
( $L_q$ norm)
$\hookrightarrow$ as the number of dimension increases, minkowski distance becomes meaningless

reducing dimentionality =
$\hookrightarrow$ feature selection based on domain knowledge (eliminatin redundant and irrelevant features)
$\hookrightarrow$ gain-based feature selection (trying different subsets and taking the best one)
$\hookrightarrow$ distribution based feature selection (principle component analysis)

$cov(X,Y) = E\left[ (X-\mu_x)(Y-\mu_y) \right] = E(XY) - E(x).E(y)$
covariance matrix = $\Sigma_{ij} = E\left[ (x_i - \mu_i)(x_j - \mu_j) \right]$

mahalanobis distance = $\sqrt{(\vec{x}-\vec{y})^T \Sigma^{-1} (\vec{x}-\vec{y})}$ (distance between a point and a distribution)
$\hookrightarrow$ distance of the test point from the center of mass divided by the width of the
ellipsoid in the direction of the test point
$\hookrightarrow$ unlike euclidean, it also seeks to measure the correlation between variables
✿ if all variables are independent, mahalanobis dist degrades into normalized euclidean dist.

principle component analysis (PCA) = technique that transforms high-dimensional data
into lower-dim while preserving as much info using dependencies between variables.
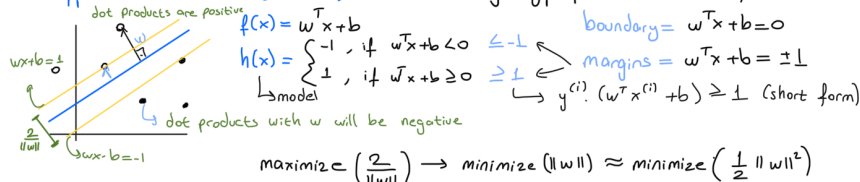  w = direction of first principle component (unit vector)      ✿ independent dimensions have
variance = of the data along the direction = $w^T \Sigma w$ ⟵ covariance matrix    largest variance
$\hookrightarrow \Sigma w = \lambda w \longrightarrow \lambda = w^{-1}\Sigma w$  ✿ highest eigenvalue is the first principle component
• each succeding component is orthogonal to the previous ones

support vector machines (SVMs) = the boundary hyperplane can have infinite dimensions



dot products are positive  $f(x) = w^T x + b$       boundary = $w^T x + b = 0$
wx+b=$\frac{1}{0}$   $h(x) = \begin{cases} -1, & \text{if } w^T x + b < 0 \quad \leq -1 \\ 1, & \text{if } w^T x + b \geq 0 \quad \geq 1 \end{cases}$  margins = $w^T x + b = \pm 1$
         $\hookrightarrow$ model          $\hookrightarrow y^{(i)} \cdot (w^T x^{(i)} + b) \geq 1$ (short form)
$\frac{2}{\|w\|}$  wx-b=-1
$\hookrightarrow$ dot products with w will be negative
         maximize $\left( \frac{2}{\|w\|} \right) \longrightarrow$ minimize $(\|w\|) \approx$ minimize $\left( \frac{1}{2} \|w\|^2 \right)$

✿ $\min_{w,b} \max_{\alpha \geq 0} \left\{ \frac{1}{2}\|w^2\| - \sum_{i=1}^{n} \alpha_i \left[ y^{(i)} (w.x^{(i)} + b) - 1 \right] \right\}$   $\alpha_i = $ lagrange multipliers $\genfrac{}{}{0pt}{}{\text{(nonzero only)}}{\text{for support vectors}}$
  $\hookrightarrow w = \sum_{i=1}^{n} \alpha_i \, y^{(i)} . x^{(i)}$   $\hookrightarrow b = \frac{1}{|SV|} \sum_{i \in sv} w.x^{(i)} - y^{(i)}$   $SV = $ support vectors set $(\alpha_i > 0)$

soft margins = allows misclassification