

Recap

kernel perceptron $x \in \mathbb{R}^d$

$$h(x) = \text{sgn}\left(\sum_{i \in \mathcal{I}} \eta_i y_i \kappa(x_i, x)\right)$$

\downarrow
incorrectly classified indices

Polynomial kernel

$$\kappa(x, x') = (c + x \cdot x')^m$$

Last time: feature space associated
is finite-dimensional

$$\Phi(x) = \begin{bmatrix} x^{(1)}(x) \\ x^{(2)}(x) \\ \sqrt{2} x^{(1)}(x) x^{(2)}(x) \\ x^{(1)^2}(x) \\ x^{(2)^2}(x) \\ \sqrt{c} \end{bmatrix}$$

$$\kappa(x, x') = \Phi(x) \cdot \Phi(x')$$

\uparrow
Euclidean dot product

Dim of feature space =

of monomials of up to
degree m .

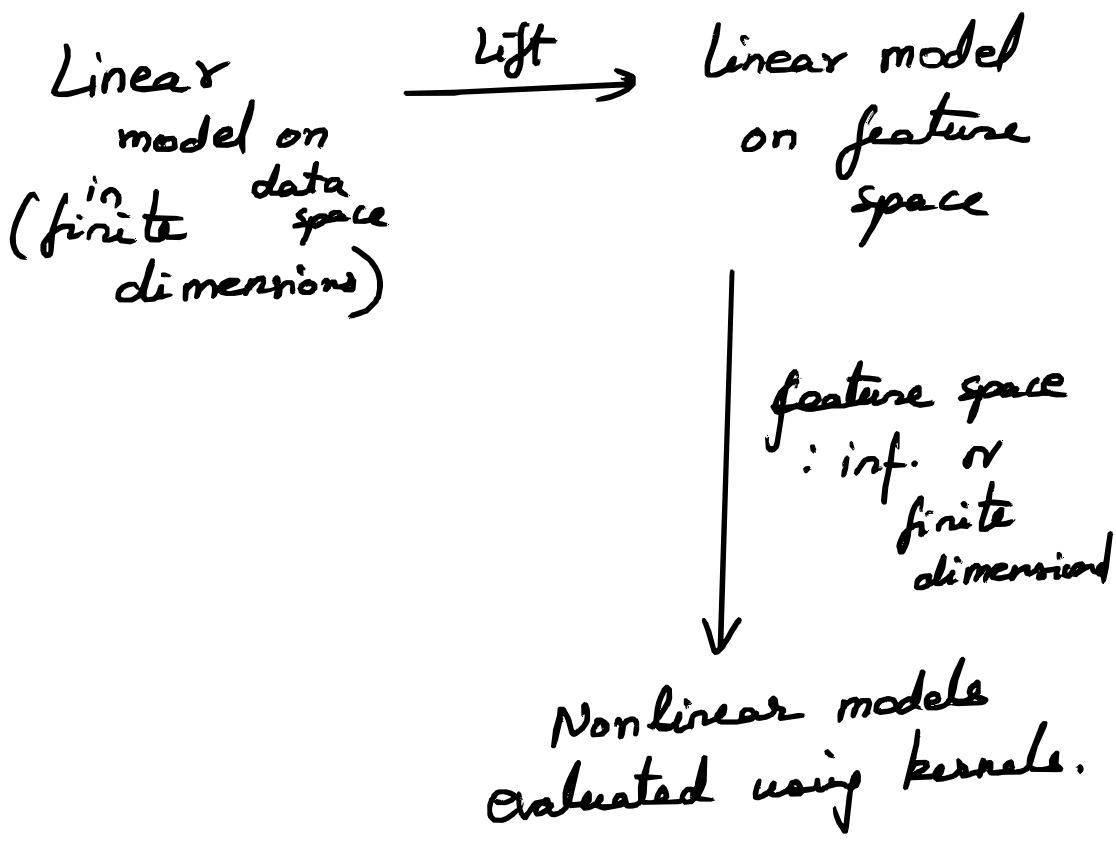
=

$$x^{(1)^{m_1}} x^{(2)^{m_2}} \dots x^{(d)^{m_d}}$$

$$\text{degree } m_1 + m_2 + \dots + m_d \leq m$$

$$\circ \circ \dots \circ \dots \circ ||||$$

$$\binom{m+d}{d}$$



Kernel: a scalar function on $\mathcal{X} \times \mathcal{X}$.

Positive definite kernels

A kernel $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ is called PD kernel if for any integer $p \in \mathbb{N}$ and $c \in \mathbb{R}^p$, $\sum_{i=1}^p c_i c_j \kappa(x_i, x_j) \geq 0$.

$\kappa(x_i, x_j) = \kappa(x_j, x_i)$ (symmetric)

Gram matrix : Given p data points $x_1, \dots, x_p \in \mathcal{X}$ a gram matrix, G , associated with a kernel $\kappa: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{C}$ is a $p \times p$ matrix with elements $G_{ij} = \kappa(x_i, x_j)$

$1 \leq i \leq p$
 $1 \leq j \leq p$.

PD kernel \iff Gram matrix
symmetric positive
semi-definite.

$$\sum_{i=1}^P \sum_{j=1}^P c_i c_j \kappa(x_i, x_j) = \mathbf{c}^T \mathbf{G} \mathbf{c} \geq 0$$

\iff \mathbf{G} is PSD

\iff all eigenvalues of \mathbf{G} are
real, non-negative

claim: if κ satisfies $\kappa(x, x') = \langle \Phi(x), \Phi(x') \rangle$

for some $\Phi \in \mathcal{H}$ (feature space),

then κ is a PD kernel

(Part of Mercer's
theorem)

Want: \mathbf{G} is PSD.

$$\begin{aligned} \sum_{i,j=1}^P c_i c_j \kappa(x_i, x_j) &= \sum_{i,j=1}^P c_i c_j \langle \Phi(x_i), \Phi(x_j) \rangle \\ &= \left\langle \sum_{i=1}^P c_i \Phi(x_i), \sum_{j=1}^P c_j \Phi(x_j) \right\rangle \\ &= \left\| \sum_{i=1}^P c_i \Phi(x_i) \right\|^2 \geq 0 \end{aligned}$$

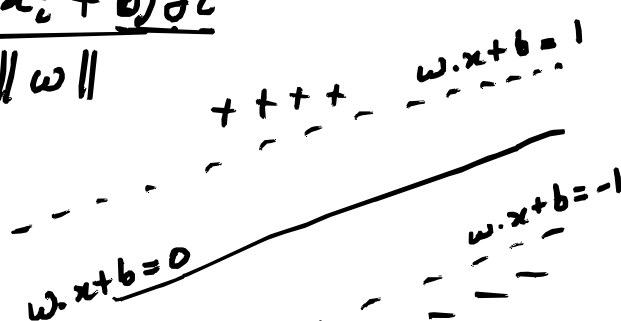
Support Vector Machine

Geometric Margin: \min distance b/w training points & separating hyperplane

Last time: Max margin \rightarrow separating hyperplane.

$$h(x) = \text{sgn}(\omega \cdot x + b)$$

$$\max_{\omega, b} \min_i \frac{(\omega \cdot x_i + b) y_i}{\|\omega\|}$$



Canonical form of hyperplane:

$$\min_i y_i(\omega \cdot x_i + b) = 1$$

$$\max_{\omega, b} \min_i \left| \frac{\omega}{\|\omega\|} \cdot x_i + \frac{b}{\|\omega\|} \right|$$

$$\max_{v \in S^{d-1}, b} \min_i |v \cdot x_i + b|$$

Derivation of max-margin classifier

(A) $\forall i \in [n], y_i(\omega \cdot x_i + b) \geq 1$

$$\max_{\omega, b} \min_i \frac{y_i(\omega \cdot x_i + b)}{\|\omega\|}$$

(B) $\max_{\omega, b} \frac{1}{\|\omega\|}$

s.t. $y_i(\omega \cdot x_i + b) \geq 1$
 $i \in [n]$

(C) $\min_{\omega, b} \frac{\|\omega\|^2}{2}$

s.t. $y_i(\omega \cdot x_i + b) \geq 1$

(SVM original formulation)

quadratic program convex problem
with affine constraints

Convex optimization

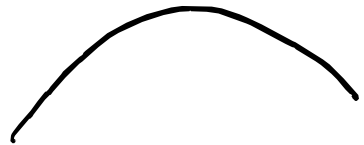
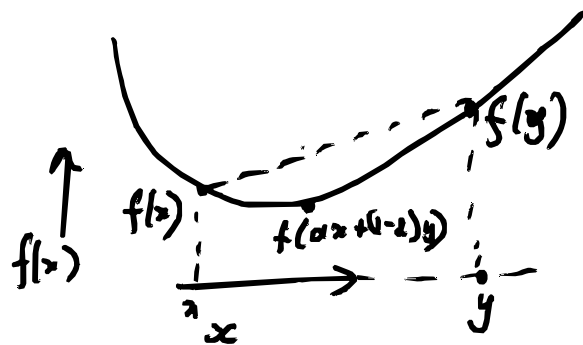
Convex function: line drawn between points on the image lies above the function

Epi. of function is convex set

$$\alpha f(x) + (1-\alpha) f(y) \geq$$

$$f(\alpha x + (1-\alpha)y)$$

$$\alpha \in (0,1)$$



Max concave or min convex
→ convex optimization

KKT condition: necessary & sufficient conditions for solution of convex program

Constrained optimization problem:

$$\text{Primal problem} \left\{ \begin{array}{l} p^* = \min_{\omega} f(\omega) \\ \text{s.t.} \quad g_i(\omega) \leq 0 \\ i = 1, \dots, p \\ \omega \in \mathbb{R}^d. \end{array} \right.$$

Convex problem:

f, g_i convex & differentiable

Lagrangian

$$\mathcal{L}(\omega, \lambda) = f(\omega) + \sum_{i=1}^p \lambda_i g_i(\omega)$$

λ : dual variables $\lambda \in \mathbb{R}^p$

Dual problem:

$$d^* = \max_{\lambda} \min_{\omega} \mathcal{L}(\omega, \lambda)$$

$$\lambda_i \geq 0$$

$$d^* \leq p^*$$

$$d^* = p^* \quad (\text{strong duality})$$

All local minima are global minima.

KKT : . f is convex, diff
 . g_i are convex, diff

. Slater's condition: $\exists \omega_0 \in \text{dom}(\omega)$
 s.t. $g_i(\omega_0) < 0$ for every i
 or
 $g_i(\omega_0) \leq 0$ for every i
 & g_i is affine.

Then, $\exists \lambda, \lambda_i \geq 0$, s.t.
 at a minimum ω , the following hold:
 $\lambda = [\lambda_1, \dots, \lambda_p]^T \in \mathbb{R}^p$

$$i) \quad \nabla_{\omega} \mathcal{L}(\omega, \lambda) = 0$$

$$(ii) \quad \nabla_{\lambda} \mathcal{L}(\omega, \lambda) \leq 0$$

$$g_i(\omega) \leq 0 \quad \forall i$$

(iii) Complementarity conditions

$$\lambda_i = 0 \quad \text{or} \quad g_i(\omega) = 0, \quad \forall i$$

$$\lambda_i g_i(\omega) = 0 \quad \forall i.$$

SVM

$$\min_{\omega, b} \frac{\|\omega\|^2}{2}$$

$$\text{s.t.} \quad y_i (\omega \cdot x_i + b) \geq 1 \quad \forall i \in [n]$$

$$\mathcal{L}(\omega, b, \lambda) = \frac{\|\omega\|^2}{2} - \sum_{i=1}^n \lambda_i (y_i (\omega \cdot x_i + b) - 1)$$

$$f(\omega) = \frac{\|\omega\|^2}{2}$$

$$g_i(\omega) = -y_i (\omega \cdot x_i + b) + 1$$

KKT conditions

$$\nabla_{\omega} \mathcal{L}(\omega, b, \lambda) = \omega - \sum_{i=1}^n \lambda_i y_i x_i = 0$$

$$\Rightarrow \omega = \sum_{i=1}^n \lambda_i y_i x_i$$

$$\nabla_b \mathcal{L}(\omega, \lambda) = 0$$

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

Physical meaning of term "support vector"
 x_i where $\lambda_i \neq 0$.

Think of each data point x_i as exerting a force $\lambda_i y_i \omega$ on separating hyperplane.

Then, force balance on hyperplane gives:

$$\sum_i F_i = \sum_{i=1}^m \lambda_i y_i \omega = 0 \quad (\nabla_b \mathcal{L} = 0)$$

Torque balance gives:

$$\sum_{i=1}^m x_i \times F_i = \sum_{i=1}^m x_i \times \lambda_i y_i \omega$$

$$\text{Cross product} = \sum_{i=1}^m \lambda_i y_i x_i \times \omega$$

$$\left(\because \sum_{i=1}^m \lambda_i y_i x_i = \omega \text{ from } \nabla_{\omega} \mathcal{L} = 0 \right) = \omega \times \omega = 0$$