

Lecture 12: 2/27/17

Kernels

$$D = \{(x^{(i)}, y^{(i)})\}_{i=1}^N \quad x^{(i)} \in \mathcal{X}, y^{(i)} \in \mathcal{Y}$$

Def: Function $k(x, z)$ is a kernel iff

$$\exists \phi: \mathcal{X} \rightarrow \mathbb{R}^D \text{ s.t. } k(x, z) = \phi(x)^T \phi(z)$$

Def: An $N \times N$ matrix K is a kernel matrix iff

① K is symmetric: $K_{ij} = K_{ji}$

② K is positive-semidefinite: $\forall a \in \mathbb{R}^N, a^T K a \geq 0$

(Mercer's Theorem)

Def: For dataset D , kernel k , the Gram matrix is

$$K_{ij} = k(x^{(i)}, x^{(j)})$$

↑ compute once at the start of training

Note: Can show k is a kernel fn. if for all datasets D the corresponding Gram matrix is symmetric and positive-semidef.

Closure Properties:

Let k_1 and k_2 be kernels.

We can define a new kernel as...

$$k(x, z) \triangleq k_1(x, z) + k_2(x, z)$$

$$k(x, z) \triangleq k_1(x, z) k_2(x, z)$$

$$k(x, z) \triangleq k_1(x, z) + c, \quad c > 0$$

$$k(x, z) \triangleq c k_1(x, z), \quad c > 0$$

$$k(x, z) \triangleq \exp(k_1(x, z))$$

Ex:

$$k(x, z) = x^T z$$

$$k(x, z) = x^T z + 1$$

$$k(x, z) = (x^T z + 1)(x^T z + 1)$$

$$k(x, z) = (x^T z + 1)^d$$

Background: Quadratic Programming

Unconstrained

$$\hat{\Theta} = \underset{\Theta \in \mathbb{R}^M}{\operatorname{argmin}} J(\Theta)$$

Constrained

$$\hat{\Theta} = \underset{\Theta \in \mathbb{R}^M}{\operatorname{argmin}} J(\Theta)$$

s.t. $g(\Theta) \leq \vec{b}$

Ex # 1:

Explicit Form

$$\hat{\vec{x}} = \underset{\vec{x}}{\operatorname{argmin}} 2x_1^2 + x_1x_2 + x_2^2 + x_1 + x_2$$

s.t.

$$x_1 \leq 1/4$$
$$-1/2 x_1 - x_2 \leq 1/2$$
$$x_2 \leq 3/4$$

Matrix Form

$$\underset{\vec{x}}{\operatorname{min}} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 2 & 1/2 \\ 1/2 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 1 \\ 1 \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$$

s.t.

$$\begin{bmatrix} 1 & 0 \\ -1/2 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 1/4 \\ 1/2 \\ 3/4 \end{bmatrix}$$

Component-wise
vector inequality

Quadratic Program (Standard Form)

Quadratic objective: $\underset{\vec{x}}{\operatorname{min}} \frac{1}{2} \vec{x}^T Q \vec{x} + \vec{c}^T \vec{x}$

Linear constraint: s.t. $A \vec{x} \leq \vec{b}$

$$Q \in \mathbb{R}^{M \times M} \quad \vec{c} \in \mathbb{R}^M$$

$$A \in \mathbb{R}^{N \times M} \quad \vec{b} \in \mathbb{R}^N$$

of constraints

of parameters

Solvers

- Interior Points
- Conjugate Gradient
- Ellipsoid Method
- Primal-Dual

Special Case:

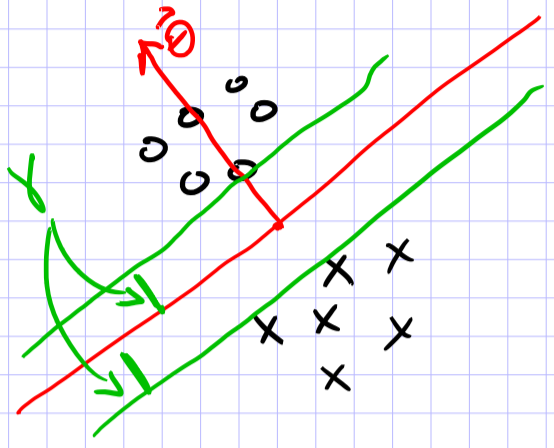
IF Q is positive-definite, the problem is convex.

$$a^T Q a > 0, \forall a \in \mathbb{R}^M$$

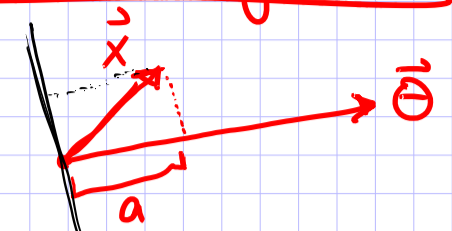
SVM

Recall: Data D is linearly separable iff

$$\exists \vec{\theta} \in \mathbb{R}^M \text{ s.t. } y^{(i)} \vec{\theta}^T \vec{x}^{(i)} \geq \gamma, \forall i$$
$$\gamma > 0$$
$$\|\vec{\theta}\|_2 = 1$$



Vector Projection



The projection of \vec{x} onto $\vec{\theta}$ has length

$$a = \frac{\vec{\theta}^T \vec{x}}{\|\vec{\theta}\|_2}$$

Below we "invent" SVMs through a sequence of optimization problems

Key idea: Find $\vec{\theta}$ with maximum margin

- Intuitively reasonable
- Theoretically well-grounded

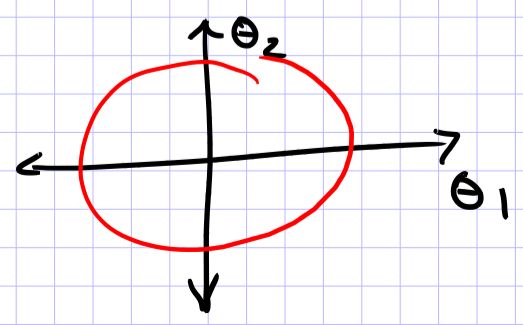
SVM (Linearly Separable Case)

Opt #1

Goal: Maximize the margin s.t. all points are correctly classified

$$\max_{\gamma, \vec{\theta}} \gamma$$
$$\text{s.t. } y^{(i)} \vec{\theta}^T \vec{x}^{(i)} \geq \gamma \quad \forall i$$
$$\|\vec{\theta}\|_2 = 1$$

Problem: $\|\vec{\theta}\|_2 = 1$ is nonlinear and nonconvex



Opt #2

Define $\vec{w} = \frac{\vec{\Theta}}{\gamma}$

If $\|\vec{\Theta}\|_2 = 1$, then

$$(\|w\|_2)^2 = \frac{1}{\gamma^2} (\|\vec{\Theta}\|_2)^2 = \frac{1}{\gamma^2}$$

\Rightarrow minimizing $(\|w\|_2)^2$ is equivalent to maximizing $\sum_{n=1}^M (w_n)^2$

~~$\max_{\gamma, \vec{\Theta}, \vec{w}} \gamma$~~ $\rightarrow \min_{\vec{w}} (\|w\|_2)^2$

~~s.t. $\gamma^{(i)} \vec{\Theta}^T \vec{x}^{(i)} \geq \gamma \quad \forall i$~~

~~$\|\vec{\Theta}\|_2 = 1$ irrelevant~~

$\min_{\vec{w}} (\|w\|_2)^2$

s.t. $y^{(i)} \vec{w}^T x^{(i)} \geq 1, \forall i$

Linear Constraints

Quadratic Objective

★ Standard Form for a Quadratic Program

★ The problem is convex

Problem: not linearly separable?

Opt #3

Goal: Maximize margin and Minimize # of errors

$$\min_{\vec{w}} (\|w\|_2)^2 + C (\# \text{ of errors on train})$$

Problem: NP-Hard to solve this!
Sign() is a step function

Opt # 4

Goal: Maximize margin, Minimize Hinge Loss

Def: the Hinge Loss is $\max(0, 1 - y^{(i)} \vec{w}^T x^{(i)})$

$$\min_{\vec{w}} (\|\vec{w}\|_2)^2 + C \left(\sum_{i=1}^N \max(0, 1 - y^{(i)} \vec{w}^T x^{(i)}) \right)$$

↕ or Equivalently

$$\min_{\vec{w}, \vec{e}} (\|\vec{w}\|_2)^2 + C \left(\sum_{i=1}^N e_i \right)$$

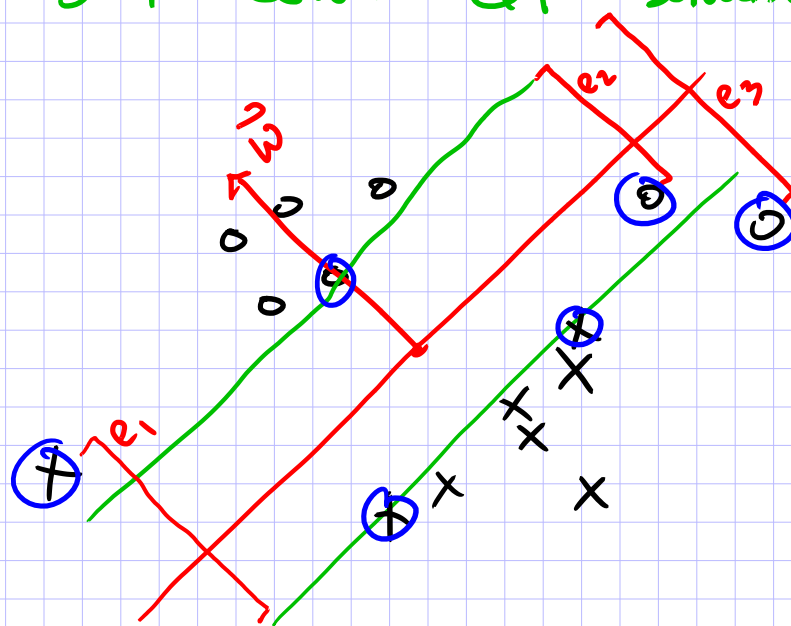
$$\text{s.t. } y^{(i)} \vec{w}^T x^{(i)} \geq 1 - e_i, \forall i$$

$$e_i \geq 0$$

called "slack variables"

* C is user-chosen hyperparameter which trades off between max margin and min error goals.

* Still convex QP solvable by black-box solvers



blue circled points are "support vectors"