

03/30/15

①

## SVM DUALITY + KERNELS

## ① Hard-margin SVM dual

General Lagrangian Duality | SVM (Special case) <sup>OP</sup>

$$\min_{\vec{w}} f(\vec{w})$$

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

$$\text{s.t. } g(\vec{w}) \leq 0$$

:  $\alpha$  $\alpha_i$ 

$$y_i (\vec{w}^T \vec{x}_i + b) \geq 1 \quad \forall i$$

$$h(\vec{w}) = 0$$

:  $\beta$ 

N/A

Lagrangian Multiplier / Dual Variables

$$\rightarrow L(\vec{w}, \alpha, \beta) = f(\vec{w}) + \alpha g(\vec{w}) + \beta h(\vec{w}) = \frac{1}{2} \vec{w}^T \vec{w} + \sum_{i=1}^N \alpha_i [1 - y_i (\vec{w}^T \vec{x}_i + b)]$$

\*KKT Condition #1

$$\frac{\partial L}{\partial \vec{w}} (\vec{w}^*, \alpha^*, \beta^*) = 0$$

$$\frac{\partial L}{\partial \vec{w}} = \vec{w}^T + \sum_{i=1}^N \alpha_i [-y_i \vec{x}_i^T] = 0$$

$$\Rightarrow \boxed{\vec{w}^* = \sum \alpha_i^* y_i \vec{x}_i}$$

$\vec{w}$  is a linear combination of data!

$$\text{Also } \frac{\partial L}{\partial b} = \sum_{i=1}^N \alpha_i [-y_i] = 0 \Rightarrow \boxed{\sum_{i=1}^N \alpha_i^* y_i = 0}$$

constraint on  $\vec{a}$

→ Dual Function

$$L_D(\alpha, \beta) = \min_{\vec{w}} f(\vec{w}) + \alpha g(\vec{w}) + \beta h(\vec{w})$$

$$L_D(\alpha_1, \dots, \alpha_N) = \frac{1}{2} \vec{w}^{*T} \vec{w}^* + \sum_{i=1}^N \alpha_i [1 - y_i (\vec{w}^{*T} \vec{x}_i + b^*)]$$

$$L_D(\vec{\alpha}) = \frac{1}{2} \left[ \sum_{i=1}^N \alpha_i y_i x_i \right]^T \left[ \sum_{j=1}^N \alpha_j y_j x_j \right]$$

$$+ \sum_{i=1}^N \alpha_i \left[ 1 - y_i \left( \left( \sum_{j=1}^N \alpha_j y_j x_j \right)^T x_i + b^* \right) \right]$$

$$= \frac{1}{2} \left[ \sum_i \alpha_i y_i x_i \right]^T \left[ \sum_j \alpha_j y_j x_j \right] + \sum_i \alpha_i - \left[ \sum_i \alpha_i y_i x_i \right]^T \left[ \sum_j \alpha_j y_j x_j \right]$$

$$- \underbrace{\sum_i \alpha_i y_i b^*}_{=0}$$

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

$$= \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N \alpha_i \alpha_j y_i y_j \vec{x}_i^T \vec{x}_j$$

So SVM Dual Program

$$\max_{\alpha_1, \dots, \alpha_N}$$

$$\vec{\alpha}^T \mathbf{1} - \frac{1}{2} \vec{\alpha}^T H \vec{\alpha}$$

$$\text{s.t.}$$

$$\sum_{i=1}^N \alpha_i y_i = 0$$

$$\alpha_i \geq 0$$

QP

→ N vars

→ N+1 constraints

easy positivity only non-trivial constraint

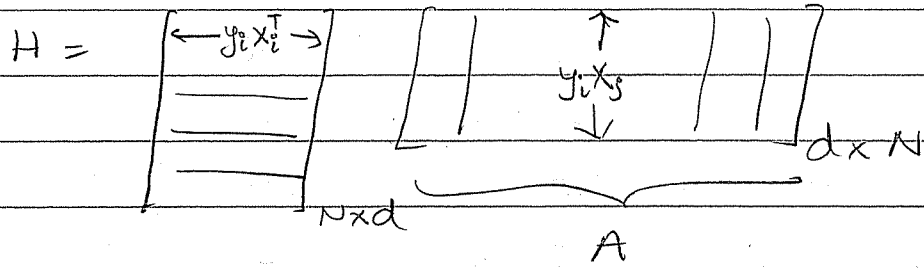
where  $H_{ij} = y_i y_j \vec{x}_i^T \vec{x}_j$

Is this a convex QP?

↔ is H PSD? Yes!

(2)

Why?



$$H = A^T A$$

$$\Rightarrow \alpha^T H \alpha = \alpha^T A^T A \alpha = \|A \alpha\|_2^2 \geq 0$$

$$\Rightarrow H \succeq 0$$

$\Rightarrow$  Dual QP objective is concave

$\Rightarrow$  QP is a convex OPT problem

$\Rightarrow$  can solve easily

But Primal QP was also convex.

So why do we care about the dual?

Answers:

① Exposes structure about the problem!

$$\text{KKT\#1} \Rightarrow \bar{b} = \sum x_i y_i \bar{x}_i$$

$$\text{KKT\#2 (Complementary Slackness)} \quad \alpha_i^* g_i(w^*) = 0$$

$$\Rightarrow \alpha_i^* [1 - y_i (w^{*T} x_i + b^*)] = 0$$

$$\text{Thus, } \alpha_i^* > 0 \Rightarrow y_i (w^{*T} x_i + b^*) = 1$$

AHA! So  $\alpha_i^* > 0 \Rightarrow i$  lies on the Margin!   
 Not just  $\geq$

$$\text{So } \vec{w} = \sum_{i=1}^N \alpha_i y_i \vec{x}_i$$

$$= \sum_{\substack{i \in \text{support} \\ \text{vectors}}} \alpha_i y_i \vec{x}_i$$

Mathematically ~~precise~~ precise way of saying  
"other data points don't matter!"

Isn't that beautiful?

Answer (2)

Dual OP has only 1 non-trivial constraint.

Unconstrained or "nearly unconstrained" QPs are sometimes easier to solve.

e.g. Platt's SMO (Sequential Minimal OP) algorithm solves SVM dual OP easily.

Answer (3) The "Kernel Trick"

Observation:

In the dual OP, data  $\{\vec{x}_i\}$  only appears in pairs as  $\vec{x}_i^T \vec{x}_j$

Idea: So if we increase dimensionality by feature augmentation

$$\vec{x}_i \rightarrow \Phi(\vec{x}_i)$$

$$\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^D \quad D \gg d$$

all we need to do is replace  $\vec{x}_i^T \vec{x}_j$  by  $\Phi(\vec{x}_i)^T \Phi(\vec{x}_j)$

What is a Kernel?

Intuitively, a kernel  $K: \mathbb{R}^d \times \mathbb{R}^d \rightarrow \mathbb{R}$  is a function  $K(\vec{x}_i, \vec{x}_j)$  that efficiently computes  $\phi(\vec{x}_i)^T \phi(\vec{x}_j)$  for some  $\phi$

Example: Polynomial Kernel of degree 2

$$\vec{u} \in \mathbb{R}^d \quad \vec{v} \in \mathbb{R}^d$$

$$K(\vec{u}, \vec{v}) = (\vec{u}^T \vec{v})^2$$

$$= \left( \sum_{i=1}^d u_i v_i \right)^2$$

$$= (u_1 v_1 + u_2 v_2 + \dots + u_d v_d)^2$$

$$= (u_1 v_1)^2 + \dots + (u_d v_d)^2 + (u_1 v_1)(u_2 v_2) + \dots +$$

$$= \sum_{i,j} u_i v_i u_j v_j$$

$$= \sum_i \sum_j (u_i u_j) (v_i v_j)$$

$$= \phi(\vec{u})^T \phi(\vec{v})$$

where  $\phi(\vec{u}) =$

$$\begin{pmatrix} u_1 u_1 \\ \vdots \\ u_1 u_d \\ u_2 u_1 \\ \vdots \\ u_2 u_d \\ \vdots \\ u_d u_1 \\ u_d u_d \end{pmatrix}$$

But  $\phi(u)^T \phi(x)$  takes  $O(d^2)$  time

$k(\vec{u}, \vec{v}) = (\vec{u}^T \vec{v})^2$  takes  $O(d)$  time