

STAT 535

11/16/23

# lecture 15

SVM

optimization  
C - SVM

Poll · future lectures

HW 5 due  
HW 6 <sup>TB</sup> out

Q3 ≥ TxG

# Lecture V: Support Vector Machines

Marina Meilă  
`mmp@stat.washington.edu`

Department of Statistics  
University of Washington

November, 2023

## Linear SVM's

The margin and the expected classification error

Maximum Margin Linear classifiers

Linear classifiers for non-linearly separable data

✓  
optimization

## Non linear SVM

The “kernel trick”

Kernels

Prediction with SVM

## Extensions

$L_1$  SVM

Multi-class and One class SVM

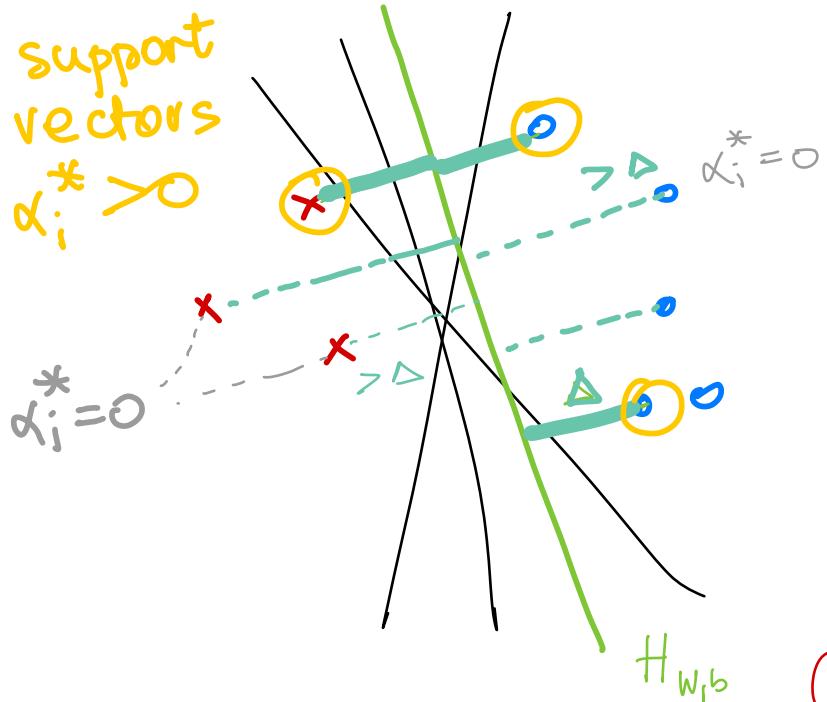
SV Regression

**Reading** HTF Ch.: Ch. 12.1–3, Murphy Ch.: Ch 14 (14.1,14.2–14.2.4 kernels, 14.4 and equations (14.28,14.29) kernel trick, 14.5.1.–3 Support Vector Machines), Bach Ch.: 7.1–7.4, 7.7

Additional Reading: C. Burges - “A tutorial on SVM for pattern recognition”

These notes: Appendices (convex optimization) are optional.

support vectors  
 $\alpha_i^* > 0$



Remark 1.

(P2) ...

$$\text{s.t. } y^i(w^T x^i + b) \geq 1$$

Remark 2

$$d(x, H_{w,b}) = \frac{|w^T x + b|}{\|w\|}$$

$$(P3) \quad \max_{w,b} \min_i |w^T x^i + b| \quad \text{s.t. } \|w\| = 1$$

Assume data  
linearly separable

margin  
 $\Delta_{w,b} = \min_i \text{dist}(x^i, H_{w,b})$

$$f(x) = w^T x + b \quad \text{linear classifier}$$

$$H_{w,b} = \{x \mid f(x) = 0\}$$

SVM pb want  $\arg \max_{w,b} \Delta_{w,b}$  s.t.

(P1)

$$y^i(w^T x^i + b) > 0$$

classify correctly

$$\text{s.t. } y^i(w^T x^i + b) \geq 1$$

$$\text{Remark 3 - } \min_i y^i(w^T x^i + b) = 1 \Rightarrow$$

$$(P4) \max_{w,b} \frac{1}{\|w\|} \text{ s.t. } y^i(w^T x^i + b) \geq 1$$

$$(P5) \min_{w,b} \|w\|^2 \text{ s.t. } \quad \text{CONVEX}$$

$$(P6) \min_{w,b} \|w\|^2 \text{ s.t. } \quad \text{QUADRATIC}$$

easier to analyse!

!!!

$$(P) \min_{w,b} \frac{1}{2} \|w\|^2 \text{ s.t. } y^i(w^T x^i + b) \geq 1 \quad \text{for all } i=1:w$$

Primal  
[opt. pb.]

Convex opt. analysis

$$f(x) = \underline{w^T x + b}$$

$$\hat{y} = \operatorname{sgn} f(x)$$

# Optimization with Lagrange multipliers

<sup>2</sup> The **Lagrangean** of (13) is

$$L(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_i \alpha_i [y^i (w^T x^i + b) - 1]. \quad (16)$$

## [KKT conditions]

At the optimum of (13)

$$w = \sum_i \alpha_i y^i x^i \quad \text{with } \alpha_i \geq 0 \quad (17)$$

and  $b = y^i - w^T x^i$  for any  $i$  with  $\alpha_i > 0$ .

- ▶ **Support vector** is a data point  $x^i$  such that  $\alpha_i > 0$ .
- ▶ According to (17), the final decision boundary is determined by the support vectors (i.e. does not depend explicitly on any data point that is not a support vector).

---

<sup>2</sup>The derivations of these results are in the Appendix

(P)

Primal  
[opt. pb.]

$$\min_{\underline{w}, \underline{b}} \frac{1}{2} \|\underline{w}\|^2 \quad \text{s.t.}$$

constraints

$$y^i (\underline{w}^T \underline{x}^i + b) \geq 1 \quad \text{for all } i=1:w \leftarrow \alpha_i$$

Convex opt. analysis  
↓

Lagrangian function

$$L: L = \frac{1}{2} \|\underline{w}\|^2 + \sum_{i=1}^n \alpha_i [1 - y^i (\underline{w}^T \underline{x}^i + b)]$$

Optimum:  $(\underline{w}^*, \underline{b}^*, \alpha^*)$

$\min_{\underline{w}}$        $\max_{\alpha}$

KKT conditions

at opt

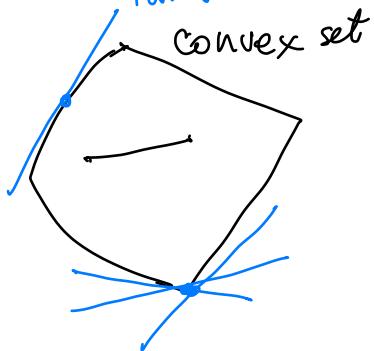
$$\begin{cases} \frac{\partial L}{\partial \underline{w}} = 0 \\ \frac{\partial L}{\partial b} = 0 \end{cases}$$

dual variable

$$1 - y^i (\underline{w}^T \underline{x}^i + b) \leq 0$$

$$\sum_{i=1}^n \alpha_i [1 - y^i (\underline{w}^T \underline{x}^i + b)] = L(\underline{w}, \underline{b}, \alpha)$$

tangent



at opt

$$\frac{\partial L}{\partial w} = 0$$

$$w + \sum_i (-y^i \alpha_i x^i)$$

$$\Rightarrow w^* = \sum_{i=1}^n \alpha_i^* y^i x^i$$

$$\alpha_i^* \geq 0 \quad (\text{no proof})$$

$$\frac{\partial L}{\partial b} = - \left[ \sum \alpha_i^* y^i \right] \quad \text{constraint on } \alpha_i^* \text{'s}$$

$$\begin{aligned} L(\alpha) &= \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i,j} \alpha_i \alpha_j y^i y^j (x^i)^T x^j + 0 \\ &= 1^T \alpha - \frac{1}{2} \alpha^T \bar{G} \alpha = g(\alpha) \end{aligned}$$

dual pb  
(D)

$$\max_{\alpha \in \mathbb{R}^n} g(\alpha) \quad \text{s.t.} \quad \alpha_i \geq 0$$

$$\sum_{i=1}^n \alpha_i y^i = 0 \quad \underbrace{\text{linear}}$$

solution  $\alpha^*$

$$L = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i [1 - y^i (w^T x^i + b)]$$

$w^* = \text{linear combination of } x^i$

$$\|w\|^2 = w^T w$$

$$-w^T \sum \alpha_i y^i x^i = -w^T \bar{w}$$

$$\alpha = \begin{bmatrix} \alpha_1 \\ \vdots \\ \alpha_n \end{bmatrix}$$

concave, quadratic

$$G = [(x^i)^T x^j]_{i,j=1:n}$$

Gram matrix

$$1 = \begin{bmatrix} 1 \\ \vdots \\ 1 \end{bmatrix} \in \mathbb{R}^n$$

$$\bar{G} = [y^i y^j (x^i)^T (x^j)]_{i,j}$$

$$= \begin{bmatrix} y^1 & \dots \\ \vdots & \ddots \end{bmatrix} G \begin{bmatrix} y^1 & \dots \\ \vdots & \ddots \end{bmatrix} \leq 0$$

(P)

Primal  
[opt. pb.]

(D)

objective

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad \text{s.t.}$$

$$y^i (w^T x^i + b) \geq 1 \quad = 1$$

for all  $i=1:w \leftarrow x_i$

$$\max_{\alpha \in \mathbb{R}^n} g(\alpha) \quad \text{s.t.} \quad \alpha_i \geq 0$$

$$\sum_{i=1}^n \alpha_i y^i = 0 \quad \parallel$$

linear

Theory

$\alpha_i^* > 0$  tight  
 $\alpha_i^* = 0$  slack constraints

support vectors

$(y^i, x^i)$

=  $x^i$ 's on the margin

$$w^* = \sum_{\substack{i=1 \\ \alpha_i^* > 0}}^n \alpha_i^* y^i x^i$$

STATISTICS, finally

- depends on subset of  $(x^i, y^i)$
- directly  $\rightarrow$  complexity
- depends on data
- # params = # S.V.

$b^*$

for S.V.:  $y^i (w^T x^i + b) = 1$   $\Rightarrow$  low variance

$\hookrightarrow$  solve for  $\underline{b}^*$   $w^T x^i + b = y^i$

## SVM Training

1.  $\{(\mathbf{x}_i, y_i)\} \rightarrow \text{compute } \mathbf{G}, \bar{\mathbf{G}}$
2. Solve (D) dual pb  $\Rightarrow \alpha_i^*, i=1:n$
3.  $\mathbf{w} = \sum_{i=1}^n \alpha_i^* y_i \mathbf{x}_i$   
 $\alpha_i^* > 0$
4. for i support vector : solve for b  
 $\alpha_i^* > 0$

## Prediction

new  $\mathbf{x}$  :  $f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b = \sum_{i, \alpha_i^* > 0} \alpha_i^* y_i (\mathbf{x}_i)^T \mathbf{x} + b$

$\hat{y} = \text{sgn } f(\mathbf{x})$

- So far: linear  $f$ , data linearly separable ✓
- Next: —||—, data NOT —||— —||—
- —||— : Non Linear  $f$

Computation efficient  
(D)  $\alpha \in \mathbb{R}^n$  for  $n \ll d$

$n$  variables

(n+1) constraints

(P)  $d+1$  variables

$n$  constraints

$\mathbf{w} \in \mathbb{R}^d$

$b \in \mathbb{R}$

efficient for  $d$  small  
 $n$  large

## Dual SVM optimization problem

- ▶ Any convex optimization problem has a **dual** problem. In SVM, it is both illuminating and practical to solve the dual problem.
- ▶ The dual to problem (13) is

$$\max_{\alpha_{1:n}} \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y^i y^j x^{iT} x_j \text{ s.t. } \alpha_i \geq 0 \text{ for all } i \text{ and } \sum_i \alpha_i y^i = 0. \quad (18)$$

- ▶ This is a **quadratic** problem with  $n$  variables on a convex domain.
- ▶ Dual problem in matrix form

- ▶ Denote  $\alpha = [\alpha_i]_{i=1:n}$ ,  $y = [y^i]_{i=1:n}$ ,  $G_{ij} = x^{iT} x_j$ ,  $\bar{G}_{ij} = y^i y^j x^{iT} x_j$ ,  $G = [G_{ij}] \in \mathbb{R}^{n \times n}$ ,  $\bar{G} = [\bar{G}_{ij}] \in \mathbb{R}^{n \times n}$ .

$$\max_{\alpha \in \mathbb{R}^n} 1^T \alpha - \frac{1}{2} \alpha^T \bar{G} \alpha \quad \text{s.t. } \alpha \succeq 0 \text{ and } y^T \alpha = 0. \quad (19)$$

- ▶  $g(\alpha) = 1^T \alpha - \frac{1}{2} \alpha^T \bar{G} \alpha$  is the **dual objective function**
- ▶  $G$  is called the **Gram matrix** of the data. Note that  $\bar{G} = \text{diag}\{y^{1:n}\}^T G \text{diag}\{y^{1:n}\}$ .
- ▶ At the dual optimum
  - ▶  $\alpha_i > 0$  for constraints that are satisfied with equality, i.e. **tight**
  - ▶  $\alpha_i = 0$  for the **slack** constraints

# Non-linearly separable problems and their duals

The **C-SVM**

$$\begin{array}{l} C \uparrow \text{less slack} \\ \Rightarrow \text{bias } \downarrow \\ \text{var } \rightarrow \end{array} \quad \begin{array}{l} \text{minimize}_{w, b, \xi} \\ \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \end{array}$$

*$C > 0$  regularization  
as little slack as possible*

$$\begin{array}{l} \text{s.t.} \\ y^i(w^T x^i + b) \geq 1 - \xi_i \\ \xi_i \geq 0 \end{array} \quad \begin{array}{l} \text{slack variables} \end{array} \quad (20)$$

In the above,  $\xi_i$  are the **slack variables**. Dual<sup>3</sup>:

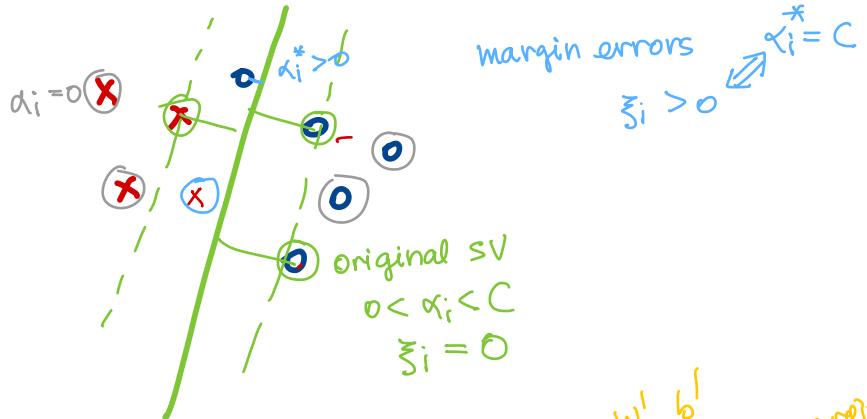
$$\begin{array}{ll} \text{maximize}_{\alpha} & \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y^i y_j x^i x^j \\ \text{s.t.} & \boxed{C \geq \alpha_i \geq 0 \text{ for all } i} \\ & \sum_i \alpha_i y^i = 0 \end{array} \quad (21)$$

⇒ two types of SV

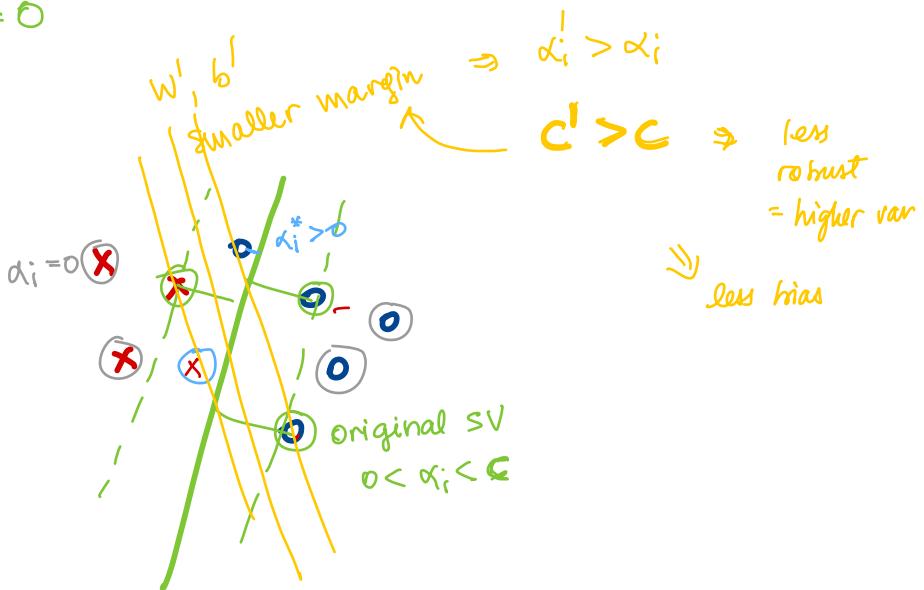
- $\alpha_i < C$  data point  $x^i$  is “on the margin”  $\Leftrightarrow y^i(w^T x^i + b) = 1$  (original SV)
- $\alpha_i = C$  data point  $x^i$  cannot be classified with margin 1 (**margin error**)  
 $\Leftrightarrow y^i(w^T x^i + b) < 1$

---

<sup>3</sup>Lagrangian  $L(w, b, \xi, \alpha, \mu) = \frac{1}{2} \|w\|^2 + C \sum_i \xi_i - \sum_i \alpha_i [y^i(w^T x^i + b) - 1 + \xi_i] - \sum_i \mu_i \xi_i$  with  $\alpha_i \geq 0, \xi_i \geq 0, \mu_i \geq 0$



In general:  
 $\alpha_i^* \uparrow \Leftrightarrow$  point  $i$  harder to classify



Lectures 1.5-2  
 —————  
 $\leq 1$       Boosting SVM  
 [kernel machines]  $\rightarrow$  R Fourier F. \*  $\xrightarrow{*}$  Double Descent  
 $\leq 1$       Neuro-tangent \* + ...

1.5 — Clust K-means + EM  
 1.5      Spectral clustering  
 $\leq 1$       Model selection AIC, BIC, struct Risk Min  $\xleftarrow{\text{VC dim}}$

4.5 lectures

Kernel tricks

RKHS  $\mathcal{H}$

## The $\nu$ -SVM

$$\underset{w, b, \xi, \rho}{\text{minimize}} \quad \frac{1}{2} \|w\|^2 - \nu\rho + \frac{1}{n} \sum_i \xi_i \quad (22)$$

$$\text{s.t.} \quad y^i(w^T x^i + b) \geq \rho - \xi_i \quad (23)$$

$$\xi_i \geq 0 \quad (24)$$

$$\rho \geq 0 \quad (25)$$

where  $\nu \in [0, 1]$  is a parameter.

Dual<sup>4</sup>:

$$\underset{\alpha}{\text{maximize}} \quad -\frac{1}{2} \sum_i \alpha_i \alpha_j y^i y^j x^{iT} x^j \quad (26)$$

$$\text{s.t.} \quad \frac{1}{n} \geq \alpha_i \geq 0 \text{ for all } i \quad (27)$$

$$\sum_i \alpha_i y^i = 0 \quad (28)$$

$$\sum_i \alpha_i \geq \nu \quad (29)$$

**Properties** If  $\rho > 0$  then:

- ▶  $\nu$  is an upper bound on #margin errors/n (if  $\sum_i \alpha_i = \nu$ )
- ▶  $\nu$  is a lower bound on #(original support vectors + margin errors)/n
- ▶  $\nu$ -SVM leads to the same  $w, b$  as C-SVM with  $C = 1/\nu$

<sup>4</sup>Lagrangian  $L(w, b, \xi, \rho, \alpha, \mu, \delta) = \frac{1}{2} \|w\|^2 - \nu\rho + \frac{1}{n} \sum_i \xi_i - \sum_i \alpha_i [y^i(w^T x^i + b) - \rho + \xi_i] - \sum_i \mu_i \xi_i - \delta\rho$   
with  $\alpha_i \geq 0, \delta \geq 0, \mu_i \geq 0$