

BIOSTAT527

5/1/23

# Lecture 11

SVM

# Lecture V: Support Vector Machines and Kernel Machines

Marina Meilă  
`mmp@stat.washington.edu`

Department of Statistics  
University of Washington

April, 2023

## Linear SVM's

The margin and the expected classification error ✓

Maximum Margin Linear classifiers ✓

Linear classifiers for non-linearly separable data

optimization  
matters !!

Analysis of  
solution

## Non linear SVM

The "kernel trick"

Kernels

Prediction with SVM

## Extensions

$L_1$  SVM

Multi-class and One class SVM

SV Regression

Reading AoNPS Ch.: Ch. 12.1–3, HTF 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) 7.1–7.4, 7.7

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

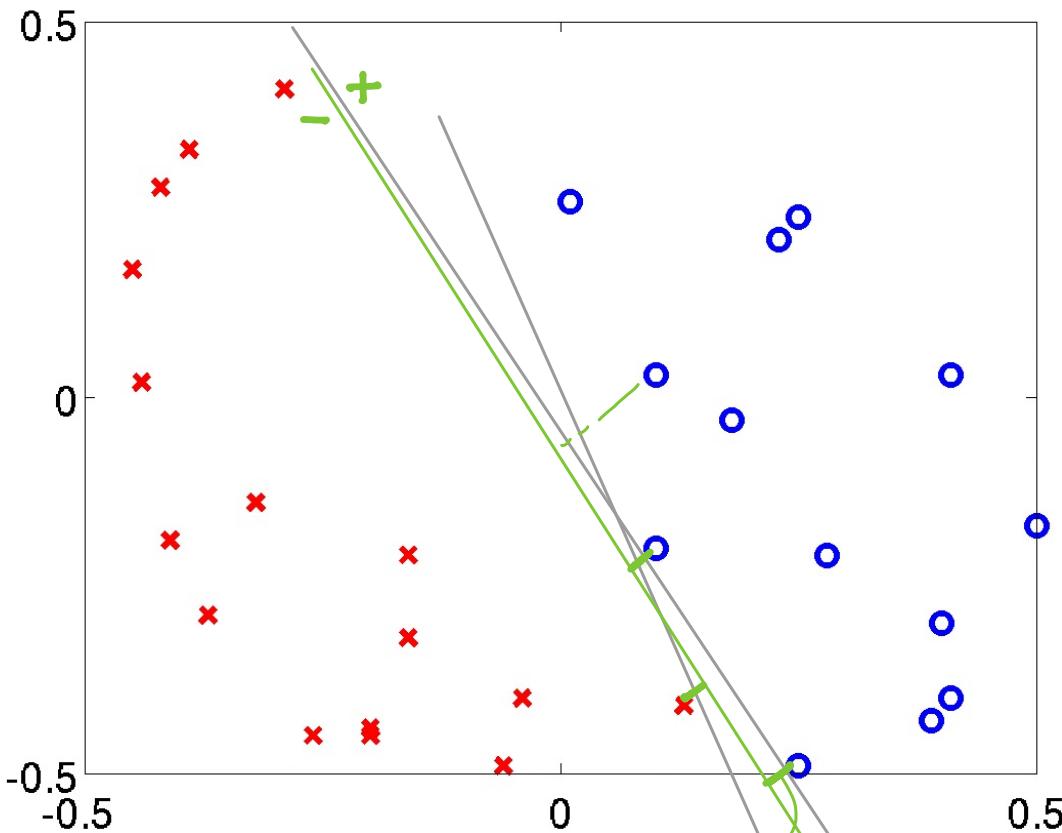
These notes: Appendices (convex optimization) are optional.

Data Linearly separable  $\Rightarrow \infty$  linear classifiers  
How to choose  $\hat{w}$ ?

$$f(x) = w^T x + b$$
$$x, w \in \mathbb{R}^d$$

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

Idea 1. Max Margin



Parametric  
LDA  
LR  
Perceptron

## Second idea

The **Second idea** is to formulate (10) as a **quadratic** optimization problem.

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

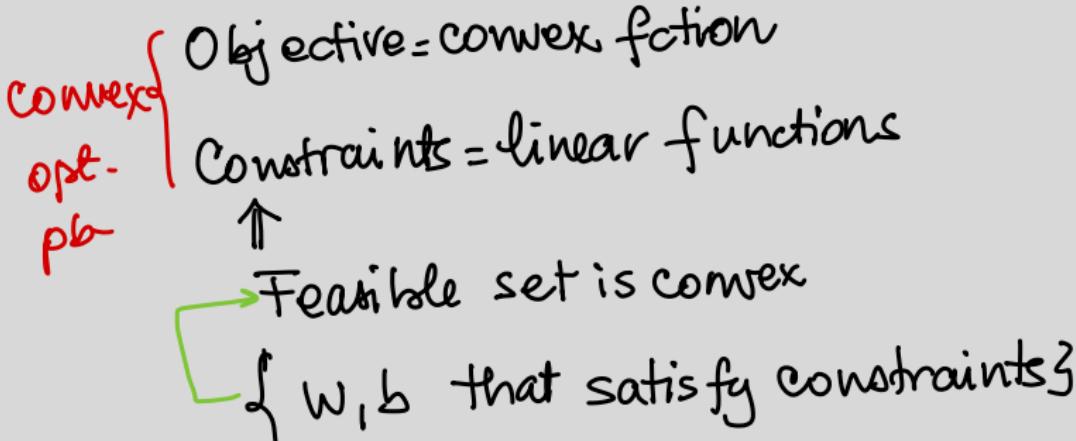
primal variables
objective
constraints
convex opt. pb.

(13)

This is the **Linear SVM (primal) optimization problem**

- ▶ This problem has a strongly convex **objective**  $\|w\|^2$ , and **constraints**  $y^i(w^T x^i + b)$  linear in  $(w, b)$ .
- ▶ Hence this is a convex problem, and can be studied with the tools of convex optimization.

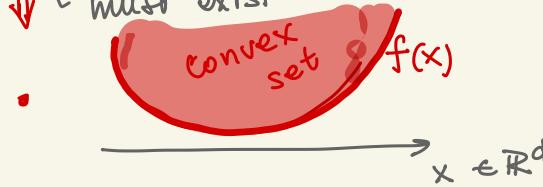
Optimization in the analysis of SVM - pb. solution



## Convex function

- $\nabla^2 f(x) \succeq 0$  for all  $x$

Hessian positive definite  
must exist



- Convex set  $A \Leftrightarrow$

$$x, x' \in A \Rightarrow \underbrace{tx + (1-t)x'}_{t \in [0,1]} \in A$$

line segment  $[x, x']$

SVM

$$(P) \min_{w,b} \frac{1}{2} \|w\|^2$$

s.t.  $\underbrace{-\left(y_i(w^T x_i + b) - 1\right)}_{\text{standard form}} \leq 0 \leftarrow \alpha_i$   
 $\alpha_i \geq 0$

$x, w \in \mathbb{R}^d$

$d+1$  variables

$w \in \mathbb{R}^d, b \in \mathbb{R}$

$n$  constraints

Lagrange function

$$L(\underbrace{w, b}_{\text{primal variables}}, \underbrace{\alpha_{1:n}}_{\text{dual variables}}) = \frac{1}{2} \|w\|^2 + \sum_{i=1}^n \alpha_i \underbrace{\left[ -y_i(w^T x_i + b) + 1 \right]}_{\text{constraint}}$$

KKT Conditions for optimum

$$\star \quad \frac{\partial L}{\partial w} = 0 = w - \sum_{i=1}^n \alpha_i y_i x_i \Rightarrow w = \sum_{i=1}^n \alpha_i y_i x_i$$

$$\circledcirc \quad \frac{\partial L}{\partial b} = 0$$

$$\alpha_i [\dots] = 0 \Rightarrow \begin{cases} \alpha_i = 0 \text{ and } [\dots] \neq 0 \leftarrow g_i < 0 & \text{slack} \\ \alpha_i > 0 \text{ and } [\dots] = 0 & \leftarrow \text{constraint is } \underline{\text{tight}} \end{cases}$$

$$\square \quad \frac{\partial L}{\partial \alpha_i} = 0 \text{ if } \alpha_i > 0 \quad \alpha_i = 0 = [\dots] \text{ w.p. 0 ignore!}$$

can't increase  $g_i$

$$w = \sum_{i=1}^n \alpha_i y_i x_i \implies w \text{ is linear combination of } x^{1:n}$$

↓  
d parameters? n params?  
NO      #params ≈  
        # { $\alpha_i > 0$ }

Ex: Logistic regression

when is  $\hat{p} = \sum \alpha_i x^i$ ?

- if  $\alpha_i = 0$  (slack)  $\rightarrow w$  does not depend on  $(x^i, y^i)$

- $\alpha_i > 0 \Rightarrow x^i$  is support vector
- [• if few  $\alpha_i > 0$  then solution is sparse]

**¶**  $\frac{\partial L}{\partial b} = 0 = -\sum_{i=1}^n \alpha_i y_i \leftarrow \text{constraint on } \alpha_{1:n} !!$

for sv    $\alpha_i > 0$

$$\begin{aligned} y^i (w^T x^i + b) &= 1 \\ \Rightarrow [b = y^i - w^T x^i \in \mathbb{R}] \\ &\text{same for all } i \text{ S.V.} \end{aligned}$$

Solving SVM problem

• plug in  $w$ , constraint  $\sum y^i \alpha_i = 0$  in  $L$

Dual SVM problem

$$(D) \quad \max_{\alpha_{1:n}} \underbrace{1^T \alpha - \frac{1}{2} \alpha^T \tilde{G} \alpha}_{n \text{ variables} \uparrow} \quad \text{s.t.} \quad \begin{aligned} \alpha_i &\geq 0 \\ \sum_{i=1}^n \alpha_i y^i &= 0 \end{aligned}$$

quadratic objective  
+ linear constraints

Dual objective

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

$$-\left(\sum \alpha_i y^i x^i\right)^T w = -\|w\|^2$$

$W^T W$

$$\|w\|^2 = \|\sum \alpha_i y^i x^i\|^2 = \alpha^T \tilde{G} \alpha$$

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

$$G = \left[ (x^i)^T x^j \right]_{i,j=1:n} \quad \text{Gram matrix}$$

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

$$\tilde{G} = \left[ y^i y^j (x^i)^T x^j \right]_{i,j=1:n}$$

# 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

## 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

$$\begin{aligned} & \text{minimize}_{w,b,\xi} \quad \frac{1}{2} \|w\|^2 + C \sum_i \xi_i \\ \text{s.t.} \quad & y^i(w^T x^i + b) \geq 1 - \xi_i \\ & \xi_i \geq 0 \end{aligned} \tag{20}$$

regularization constant

slack variables

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

$$\begin{aligned} & 1^T \alpha - \frac{1}{2} \alpha^T G \alpha \quad \xrightarrow{\text{maximize } \alpha} \quad \sum_i \alpha_i - \frac{1}{2} \sum_i \alpha_i \alpha_j y^i y_j x^{iT} x_j \quad \checkmark \\ \text{s.t.} \quad & C \geq \alpha_i \geq 0 \text{ for all } i \\ & \sum_i \alpha_i y^i = 0 \quad \checkmark \end{aligned} \tag{21}$$

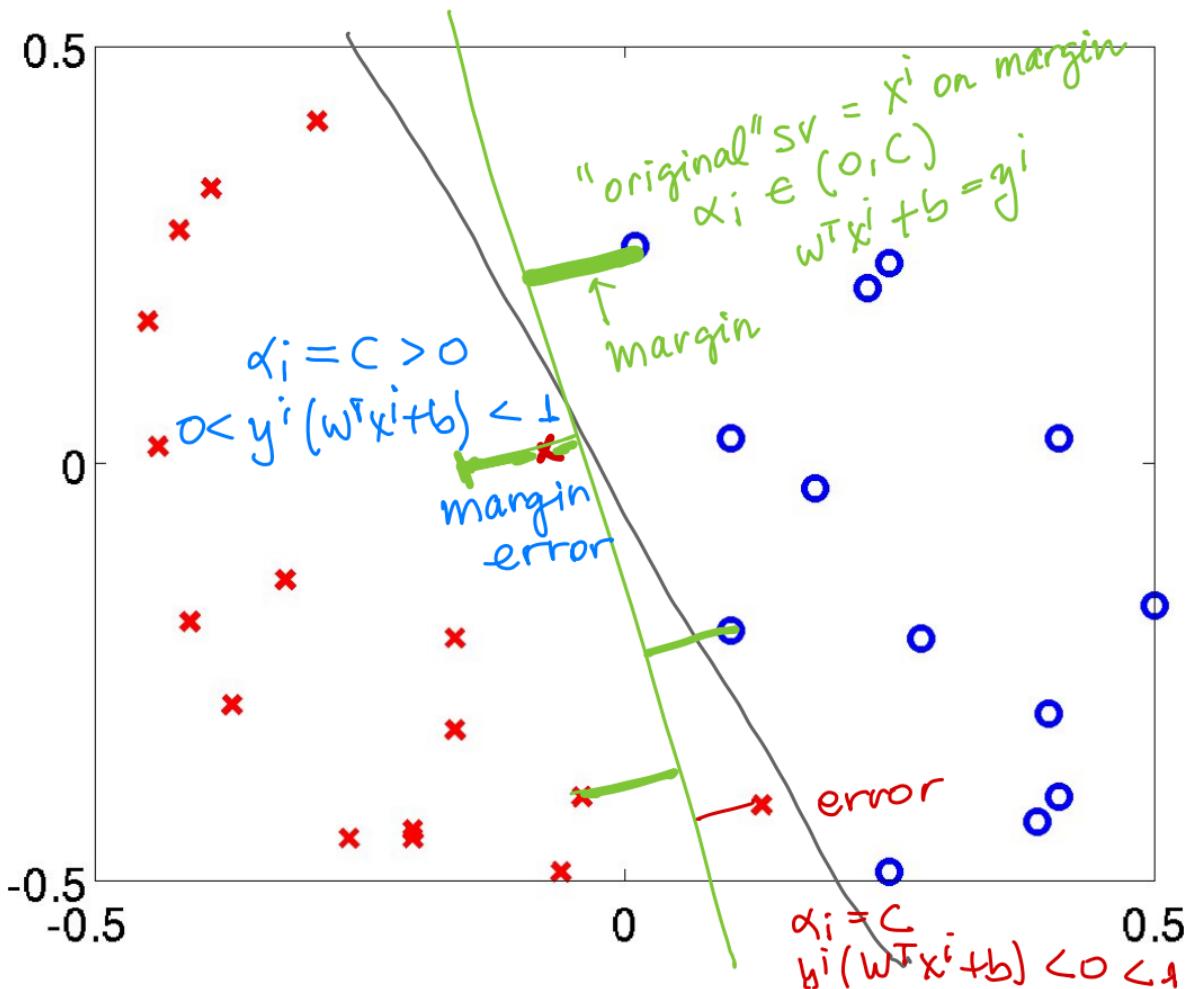
new constraints

⇒ 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)  $\xi_i = 0$
- ▶  $\alpha_i = C$  data point  $x^i$  cannot be classified with margin 1 (margin error)  
 $\Leftrightarrow y^i(w^T x^i + b) < 1$   $\xi_i > 0$
- ▶  $\alpha_i = 0$   $x^i$  not s.v.  $y^i(w^T x^i + b) < 1 \Rightarrow < 0$  error  
 $\Rightarrow > 0$  proper margin error

---

<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  $\alpha_i \geq 0, \xi_i \geq 0, \mu_i \geq 0$



## The $\nu$ -SVM

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

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

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

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

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

Dual<sup>4</sup>:

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

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

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

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

**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  $\alpha_i \geq 0, \delta \geq 0, \mu_i \geq 0$

## A simple error bound

$$L_{01}(f_n) \leq E \left[ \frac{\#\text{support vectors of } f_{n+1}}{n+1} \right] \quad (30)$$

where  $f_n$  denotes the SVM trained on a sample of size  $n$ .

**Exercise** Use the Homework 6 to prove this result.