# Extra Topic (Reading Material)
# Lagrange Multipliers
# & Duality of SVMs

Rishabh Iyer

University of Texas at Dallas

Most Slides are from Nick Rouzzi! Thanks Nick!!
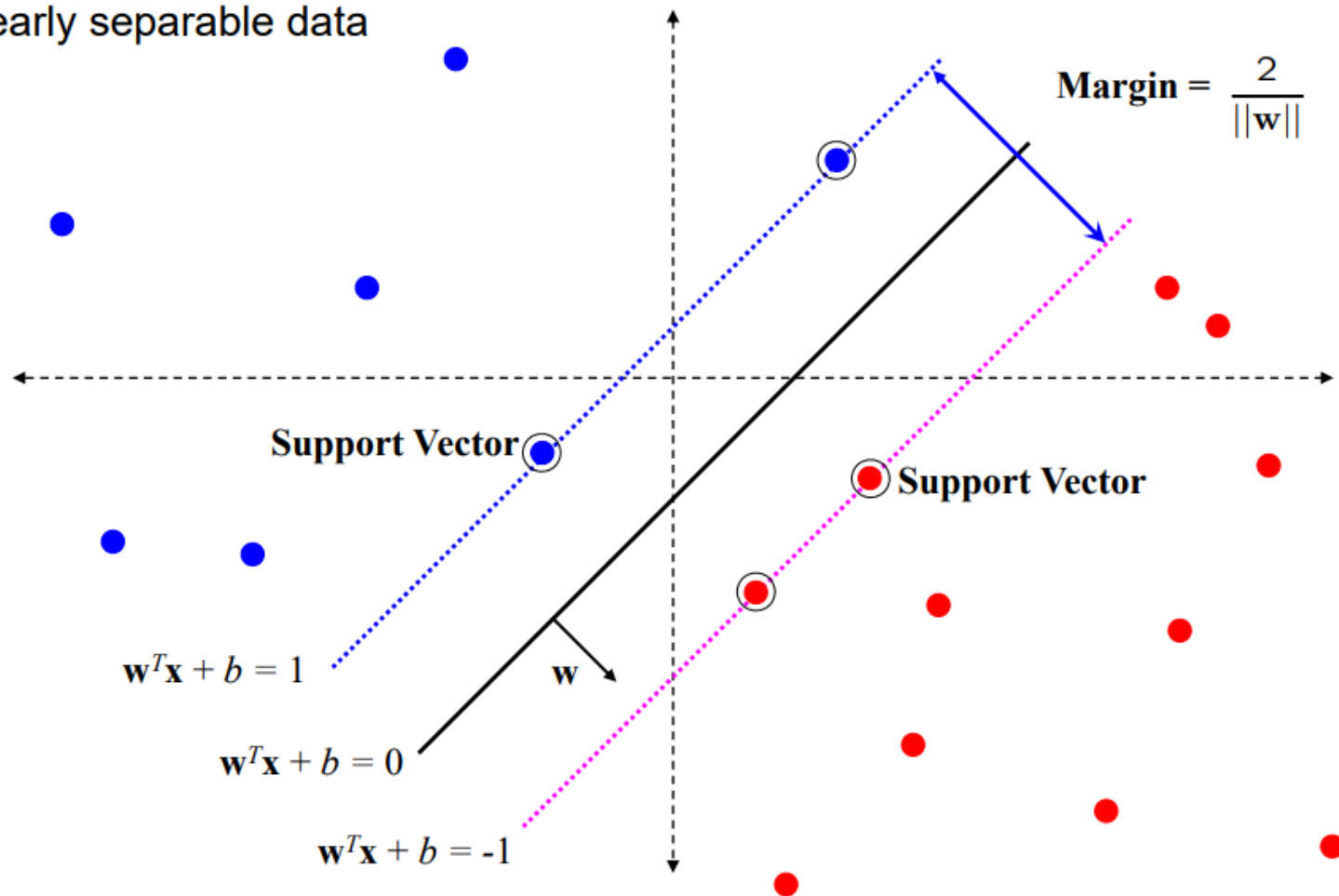
# The Strategy So Far…

- Choose hypothesis space

- Construct loss function (ideally convex)

- Minimize loss to "learn" correct parameters

# Recap: SVM

linearly separable data

$$\text{Margin} = \frac{2}{||\mathbf{w}||}$$

**Support Vector**

**Support Vector**

$\mathbf{w}^T\mathbf{x} + b = 1$

$\mathbf{w}$

$\mathbf{w}^T\mathbf{x} + b = 0$

$\mathbf{w}^T\mathbf{x} + b = -1$

# SVM Optimization Problem

- Recall: The SVM optimization problem:

$$\min_{w,b} \|w\|^2$$

 such that

$$y^{(i)}\left(w^T x^{(i)} + b\right) \geq 1, \text{for all } i$$

- This is a standard quadratic programming problem

  - Falls into the class of convex optimization problems

  - Can be solved with many specialized optimization tools (e.g., quadprog() in MATLAB)

# Constrained Optimization

A mathematical detour, we'll come back to SVMs soon!

$$\min_{x \in \mathbb{R}^n} f_0(x) \quad \longleftarrow \quad \frac{1}{2}\|w\|^2$$

subject to:

$(w, b)$

$$\longrightarrow f_i(x) \leq 0, \qquad i = 1, \ldots, m \quad \longleftarrow \quad 1 - y^{(i)}(w^T x^{(i)} + b)$$
$$h_i(x) = 0, \qquad i = 1, \ldots, p \qquad\qquad\qquad \leq 0$$

# Constrained Optimization

$$\min_{x \in \mathbb{R}^n} f_0(x)$$

$f_0$ is not necessarily convex

subject to:

$$f_i(x) \leq 0, \qquad i = 1, \ldots, m$$
$$h_i(x) = 0, \qquad i = 1, \ldots, p$$

# General Optimization

$$\min_{x \in \mathbb{R}^n} f_0(x)$$

subject to:

Constraints do not need to be linear

$$f_i(x) \leq 0, \qquad i = 1, \ldots, m$$
$$h_i(x) = 0, \qquad i = 1, \ldots, p$$

# Example

$$\min_{x \in \mathbb{R}^2} x_1 \log x_1 + x_2 \log x_2$$

subject to:

$$x_1 + x_2 = 1$$
$$x_1 \geq 0$$
$$x_2 \geq 0$$

# Example

$$\overbrace{\min_{x \in \mathbb{R}^3} x_1 \log x_1 + x_2 \log x_2}^{f_0(x)}$$

subject to:

$$h_1(x) = 0 \longrightarrow 1 - x_1 - x_2 = 0 \quad h_i(x)$$

$$f_1(x) \leq 0 \longrightarrow \left. \begin{array}{l} -x_1 \leq 0 \\ -x_2 \leq 0 \end{array} \right\} f_i(x)$$

$$f_2(x) \leq 0 \longrightarrow$$

# Lagrangian

$$Obj$$
$$\downarrow$$

$$Ineq\ Const.$$
$$\downarrow$$

$$Eq\ Const.$$
$$\downarrow$$

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x)$$

$$\downarrow$$
$$f_i(x) \le 0$$

$$\llcorner h_i(x) = 0$$

- Incorporate constraints into a new objective function

- $\lambda \ge 0$ and $\nu$ are vectors of ***Lagrange multipliers***

- The Lagrange multipliers can be thought of as enforcing soft constraints

# Example

$$\min_{x \in \mathbb{R}^3} x_1 \log x_1 + x_2 \log x_2$$

subject to:

$$1 - x_1 - x_2 = 0 \quad \leftarrow \nu_1 \rightarrow Eq.$$
$$-x_1 \leq 0 \quad \leftarrow \lambda_1 \quad \Big\} \; ineq.$$
$$-x_2 \leq 0 \quad \leftarrow \lambda_2$$

$$L(x_1, x_2, \nu_1, \lambda_1, \lambda_2)$$
$$= x_1 \log x_1 + x_2 \log x_2 + \nu_1 \cdot (1 - x_1 - x_2) - \lambda_1 x_1 - \lambda_2 x_2$$

# Duality

- Construct a <span style="color:red">dual function</span> by minimizing the Lagrangian over the primal variables

$$g(\lambda, \nu) = \inf_x L(x, \lambda, \nu)$$

(handwritten: min)

- $g(\lambda, \nu) = -\infty$ whenever the Lagrangian is not bounded from below for a fixed $\lambda$ and $\nu$

# Example

$$\min_{x \in \mathbb{R}^2} x_1 \log x_1 + x_2 \log x_2$$

subject to:

$$1 - x_1 - x_2 = 0$$
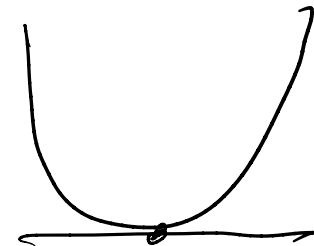$$-x_1 \leq 0$$
$$-x_2 \leq 0$$

# Example

$$\min_{x\in\mathbb{R}^3} x_1 \log x_1 + x_2 \log x_2$$

subject to:

$$1 - x_1 - x_2 = 0$$
$$-x_1 \leq 0$$
$$-x_2 \leq 0$$

$$L(x_1, x_2, \nu_1, \lambda_1, \lambda_2)$$
$$= x_1 \log x_1 + x_2 \log x_2 + \nu_1 \cdot (1 - x_1 - x_2) - \lambda_1 x_1 - \lambda_2 x_2$$

$$\frac{\partial L}{\partial x_1} = 1 + \log x_1 - \nu_1 - \lambda_1 = 0 \Rightarrow x_1 = \exp(\nu_1 + \lambda_1 - 1)$$

$$\frac{\partial L}{\partial x_2} = 1 + \log x_2 - \nu_1 - \lambda_2 = 0 \Rightarrow x_2 = \exp(\nu_1 + \lambda_2 - 1)$$

$$g(\lambda, \nu) = \text{Subst. } x_1 \, \& \, x_2 \text{ into } L(x_1, x_2, \lambda_1, \lambda_2, \nu_1)$$

# Example

$$\min_{x \in \mathbb{R}^3} x_1 \log x_1 + x_2 \log x_2$$

subject to:

$$1 - x_1 - x_2 = 0$$
$$-x_1 \leq 0$$
$$-x_2 \leq 0$$

# The Primal Problem

$$\min_{x \in \mathbb{R}^n} f_0(x) = \underset{\lambda \geq 0, \nu}{\max} \ L(x, \lambda, \nu)$$

subject to:

$$f_i(x) \leq 0, \qquad i = 1, \dots, m$$
$$h_i(x) = 0, \qquad i = 1, \dots, p$$

Equivalently,

$$\underset{x}{\inf} \ \underset{\lambda \geq 0, \nu}{\sup} \ L(x, \lambda, \nu)$$

$\min \quad \max$

**Why are these equivalent?**

16

# The Primal Problem

$$\min_{x \in \mathbb{R}^n} f_0(x)$$

subject to:

$$f_i(x) \leq 0, \qquad i = 1, \ldots, m$$
$$h_i(x) = 0, \qquad i = 1, \ldots, p$$

*Original Problem*

Equivalently,

*Min     max*

$$\inf_{x} \sup_{\lambda \geq 0, \nu} L(x, \lambda, \nu)$$

$$\sup_{\lambda \geq 0, \nu} \left[ f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x) \right] = \infty$$

whenever $x$ violates the constraints

# The Dual Problem

Equivalently,

$$\underset{\lambda \geq 0, \nu}{\text{sup}} \ g(\lambda, \nu)$$

(handwritten: max)

$$\underset{\lambda \geq 0, \nu}{\text{sup}} \ \underset{x}{\text{inf}} \ L(x, \lambda, \nu)$$

(handwritten: max, min)

- The dual problem is always concave, even if the primal problem is not convex

  - For each $x$, $L(x, \lambda, \nu)$ is a linear function in $\lambda$ and $\nu$

  - Minimum (or infimum) of linear functions is concave!

# Primal vs. Dual

$$\underbrace{\sup_{\lambda \geq 0, \nu} \inf_x L(x, \lambda, \nu)}_{\text{Dual.}} \overbrace{}^{g(\lambda, \nu)} \leq \underbrace{\inf_x \sup_{\lambda \geq 0, \nu} L(x, \lambda, \nu)}_{\text{Primal.}} \overbrace{}^{f_0(x)}$$

- Why?

  - $g(\lambda, \nu) \leq L(x, \lambda, \nu)$ for all $x$ $\qquad g(\lambda, \nu) = \min_x L(x, \lambda, \nu)$

  - $L(x', \lambda, \nu) \leq f_0(x')$ for any feasible $x'$, $\lambda \geq 0$

    - $x$ is feasible if it satisfies all of the constraints

  - Let $x^*$ be the optimal solution to the primal problem and $\lambda \geq 0$

$$g(\lambda, \nu) \leq L(x^*, \lambda, \nu) \leq f_0(x^*)$$

$$L(x, \lambda, \nu) = f_0(x) + \sum_{i=1}^{m} \lambda_i f_i(x) + \sum_{i=1}^{p} \nu_i h_i(x) = 0$$

$\leq 0$

# Duality

- Under certain conditions, the two optimization problems are equivalent

$$\sup_{\lambda \geq 0, \nu} \inf_x L(x, \lambda, \nu) = \inf_x \sup_{\lambda \geq 0, \nu} L(x, \lambda, \nu)$$

  - This is called <span style="color:red">strong duality</span>

- If the inequality is strict, then we say that there is a <span style="color:red">duality gap</span>

  - Size of gap measured by the difference between the two sides of the inequality

# Slater's Condition

For any optimization problem of the form

$$\min_{x \in \mathbb{R}^n} \underline{f_0(x)} \quad \leftarrow \text{Convex}$$

subject to:

$$\text{Convex.} \downarrow$$

$$\underline{f_i(x)} \leq 0, \qquad i = 1, \ldots, m$$

$$Ax = b \quad \leftarrow \text{Linear Equality.}$$

where $f_0, \ldots, f_m$ are <span style="color:red">convex functions</span>, strong duality holds if there exists an $x$ such that

$$f_i(x) < 0, \qquad i = 1, \ldots, m$$
$$Ax = b$$

# Dual SVM

$$\min_{w} \frac{1}{2} \|w\|^2 \leftarrow \text{Objective}$$

such that

$$y_i\left(w^T x^{(i)} + b\right) \geq 1, \text{ for all } i \leftarrow \text{InEq Constraints}$$

- Note that Slater's condition holds as long as the data is linearly separable

# Dual SVM

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_i \lambda_i (1 - y_i(w^T x^{(i)} + b))$$

Convex in $w$, so take derivatives to form the dual

$$\frac{\partial L}{\partial w_k} = w_k + \sum_i -\lambda_i y_i x_k^{(i)} = 0$$

$$\frac{\partial L}{\partial b} = \sum_i -\lambda_i y_i = 0$$

# Dual SVM

$$L(w, b, \lambda) = \frac{1}{2} w^T w + \sum_i \lambda_i (1 - y_i (w^T x^{(i)} + b))$$

Convex in $w$, so take derivatives to form the dual

$$w = \sum_i \lambda_i y_i x^{(i)}$$

$$\sum_i \lambda_i y_i = 0$$

# Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x^{(i)T} x^{(j)} + \sum_i \lambda_i$$

such that

$$k\left(x^{(i)}, x^{(j)}\right)$$

$$\sum_i \lambda_i y_i = 0$$

- By strong duality, solving this problem is equivalent to solving the primal problem

  - Given the optimal $\lambda$, we can easily construct $w$ ($b$ can be found by complementary slackness…)

$$x_{test}, \quad y_{pred} = \text{sign}(w^T x_{test} + b)$$

# Complementary Slackness

- Suppose that there is zero duality gap

- Let $x^*$ be an optimum of the primal and $(\lambda^*, \nu^*)$ be an optimum of the dual

$$f_0(x^*) = g(\lambda^*, \nu^*) \quad \leftarrow \text{Primal} = \text{Dual.}$$

$$= \inf_x \left[ f_0(x) + \sum_{i=1}^{m} \lambda_i^* f_i(x) + \sum_{i=1}^{p} \nu_i^* h_i(x) \right]$$

$$\leq f_0(x^*) + \sum_{i=1}^{m} \lambda_i^* f_i(x^*) + \sum_{i=1}^{p} \nu_i^* h_i(x^*) \quad = 0$$

$$= f_0(x^*) + \sum_{i=1}^{m} \lambda_i^* f_i(x^*)$$

$$\lambda_i^* \geq 0, \quad f_i(x^*) \leq 0$$

$$\leq f_0(x^*)$$

# Complementary Slackness

- This means that

$$\sum_{i=1}^{m} \lambda_i^* f_i(x^*) = 0$$

(handwritten: $\leq 0$ bracketed over $\lambda_i^* f_i(x^*)$)

- As $\lambda \geq 0$ and $f_i(x_i^*) \leq 0$, this can only happen if $\lambda_i^* f_i(x^*) = 0$ for all $i$

- Put another way,

  - If $f_i(x^*) < 0$ (i.e., the constraint is not tight), then $\lambda_i^* = 0$

  - If $\lambda_i^* > 0$, then $f_i(x^*) = 0$  (handwritten: $\leftarrow$ i is a support vector)

  - ONLY applies when there is no duality gap

# Dual SVM (Obtaining b)

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x^{(i)^T} x^{(j)} + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

- By complementary slackness, $\lambda_i^* > 0$ means that $x^{(i)}$ is a support vector (can then solve for $b$ using $w$)
- In particular,

$$b = y_i - w.x_i$$

for any $i$ where $\underline{\lambda_i > 0}$  [ Support Vectors ]

# Dual SVM

$$\max_{\lambda \geq 0} -\frac{1}{2} \sum_{i=1}^{M} \sum_{j=1}^{M} \lambda_i \lambda_j y_i y_j x^{(i)T} x^{(j)} + \sum_i \lambda_i$$

such that

$$\sum_i \lambda_i y_i = 0$$

$O(M^2)$

- Takes $O(m^2)$ time just to evaluate the objective function

  - Active area of research to try to speed this up