

Data Mining Techniques

CS 6220 - Section 3 - Fall 2016

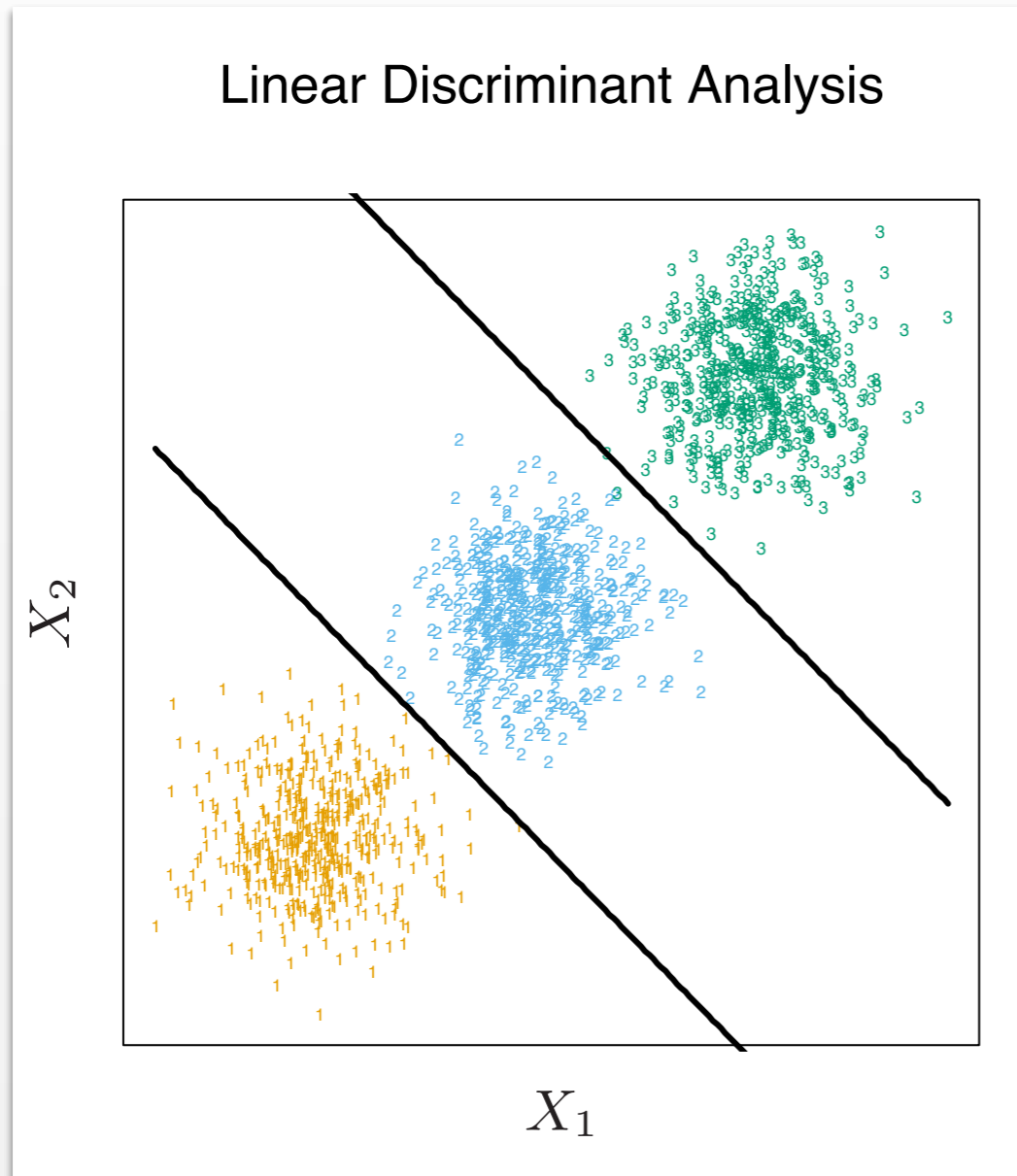
Lecture 5: Classification 2

Jan-Willem van de Meent
(*credit*: Zhao, CS 229, Bishop)



Generative Learning Algorithms

Linear Discriminant Analysis



Algorithm

- Mean for each class

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n:y_n=k} \mathbf{x}_n$$

- Covariance for each class

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n:y_n=k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Average covariance

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_k N_k \boldsymbol{\Sigma}_k$$

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Linear Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma})$$

Quadratic Discriminant Analysis

$$p(\mathbf{x} | y = k) = \mathcal{N}(\mathbf{x}; \boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Algorithm

- Mean for each class

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n:y_n=k} \mathbf{x}_n$$

- Covariance for each class

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n:y_n=k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

- Average covariance

$$\boldsymbol{\Sigma} = \frac{1}{N} \sum_k N_k \boldsymbol{\Sigma}_k$$

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Predict using posterior

$$y^* = \operatorname{argmax}_k p(y = k | \mathbf{x}^*)$$

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Predict using posterior

$$y^* = \operatorname{argmax}_k p(y = k | \mathbf{x}^*)$$

Generative Model

$$y_n \sim \text{Discrete}(\boldsymbol{\pi})$$
$$x_n | y_n = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Bayes Rule

$$p(\mathbf{x} | y) = \frac{p(y | \mathbf{x})p(\mathbf{x})}{p(y)}$$

$$p(y | \mathbf{x}) = \frac{p(\mathbf{x} | y)p(y)}{p(\mathbf{x})}$$



Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Predict using posterior

$$y^* = \operatorname{argmax}_k p(y = k | \mathbf{x}^*)$$

Generative Model

$$y_n \sim \text{Discrete}(\boldsymbol{\pi})$$
$$x_n | y_n = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Generative Learning

- Treat features as “*observations*”
- Treat class labels as “*latent variables*”
- Calculate ML estimates of parameters
- Predict according to MAP value

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Predict using posterior

$$y^* = \operatorname{argmax}_k p(y = k | \mathbf{x}^*)$$

Generative Model

$$y_n \sim \text{Discrete}(\boldsymbol{\pi})$$
$$x_n | y_n = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Maximum Likelihood Estimates

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n: y_n = k} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n: y_n = k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

$$\boldsymbol{\pi}_k = ?$$

Linear Discriminant Analysis

Predict using likelihood

$$y^* = \operatorname{argmax}_k p(\mathbf{x}^* | y = k)$$

Predict using posterior

$$y^* = \operatorname{argmax}_k p(y = k | \mathbf{x}^*)$$

Generative Model

$$y_n \sim \text{Discrete}(\boldsymbol{\pi})$$
$$\mathbf{x}_n | y_n = k \sim \mathcal{N}(\boldsymbol{\mu}_k, \boldsymbol{\Sigma}_k)$$

Maximum Likelihood Estimates

$$\boldsymbol{\mu}_k = \frac{1}{N_k} \sum_{n: y_n = k} \mathbf{x}_n$$

$$\boldsymbol{\Sigma}_k = \frac{1}{N_k} \sum_{n: y_n = k} (\mathbf{x}_n - \boldsymbol{\mu}_k)^\top (\mathbf{x}_n - \boldsymbol{\mu}_k)$$

$$\pi_k = \frac{N_k}{N}$$

Naive Bayes

Example: Spam Filtering

Features: Words in E-mail

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Labels: Spam or not Spam

$$y_n \in \{0, 1\}$$

Naive Bayes

Features: Words in E-mail

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Generative Model

$$y_n \sim \text{Bernoulli}(\mu)$$
$$x_{nd} \mid y_n = k \sim \text{Bernoulli}(\phi_{kd})$$

Conditional Independence

$$p(\mathbf{x}_n \mid y_n) = \prod_{d=1}^D p(x_{nd} \mid y_n)$$

Naive Bayes

Features: Words in E-mail

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Generative Model

$$y_n \sim \text{Bernoulli}(\mu)$$
$$x_{nd} \mid y_n = k \sim \text{Bernoulli}(\phi_{kd})$$

Maximum Likelihood

$$\mu = \frac{1}{N} \sum_{n=1}^N I[y_n = 1]$$
$$\phi_{kd} = \frac{1}{N_k} \sum_{n:y_n=k} I[x_{nd} = 1]$$

Online Estimation and Smoothing

Features: Words in E-mail

$$\mathbf{x} = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{array}{l} \text{a} \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{array}$$

Suppose word d not in training set

$$\phi_{0d} = \phi_{1d} = 0$$

$$p(\mathbf{x} | y = 0) = p(\mathbf{x} | y = 1) = 0$$

Bayes Rule

$$\begin{aligned} p(y | \mathbf{x}) &= \frac{p(\mathbf{x} | y)p(y)}{p(\mathbf{x})} \\ &= \frac{0}{0} \end{aligned}$$

Online Estimation and Smoothing

Generative model with prior

$$\mu \sim \text{Beta}(1, 1)$$

$$\phi_{kd} \sim \text{Beta}(1, 1)$$

$$y_n \sim \text{Bernoulli}(\mu)$$

$$x_{nd} \mid y_n = k \sim \text{Bernoulli}(\phi_{kd})$$

Posterior Mean

$$\mu^*, \phi^* = \mathbb{E}_{p(\mu, \phi \mid \mathbf{x}_{1:N}, \mathbf{y}_{1:N})}[\mu, \phi]$$

Conjugacy

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1 - \mu)^{N-m}$$

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1 - \mu)^{b-1}$$

$$p(\mu|m) = \frac{p(m, \mu)}{p(m)}$$

$$\propto \text{Bin}(m|N, \mu) \text{Beta}(\mu|a, b)$$

$$\propto \mu^{m+(a-1)} (1 - \mu)^{(N-m)+(b-1)}$$

Online Estimation and Smoothing

Generative model with prior

$$\mu \sim \text{Beta}(1, 1)$$

$$\phi_{kd} \sim \text{Beta}(1, 1)$$

$$y_n \sim \text{Bernoulli}(\mu)$$

$$x_{nd} \mid y_n = k \sim \text{Bernoulli}(\phi_{kd})$$

Posterior Mean

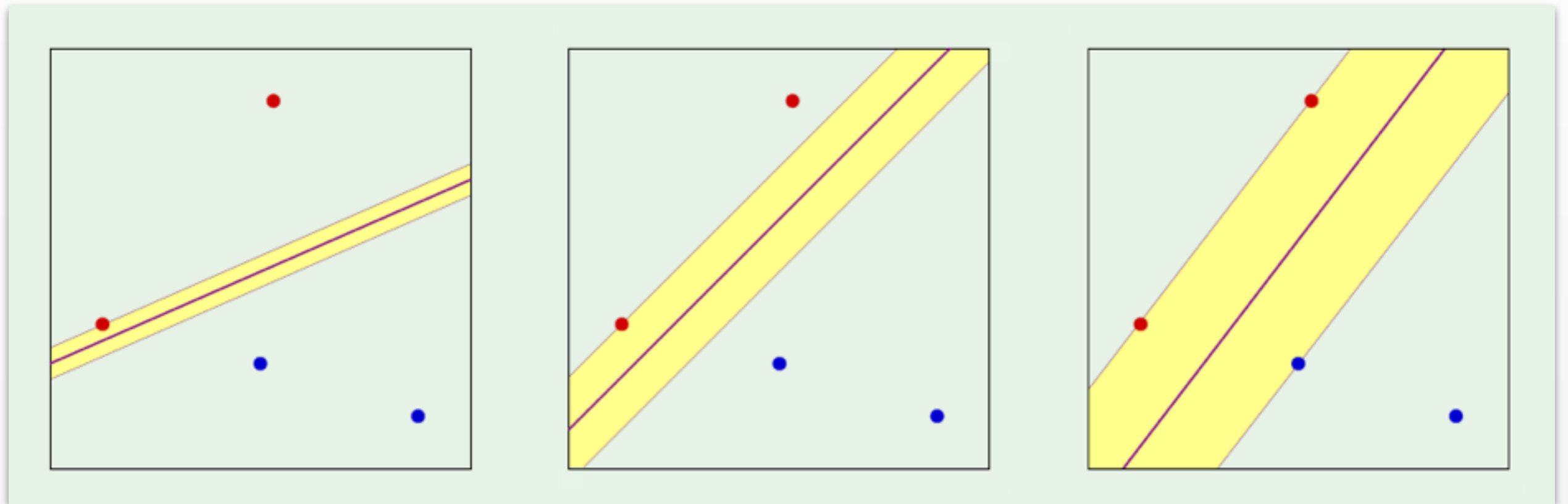
$$\mu^*, \phi^* = \mathbb{E}_{p(\mu, \phi \mid \mathbf{x}_{1:N}, \mathbf{y}_{1:N})}[\mu, \phi]$$

$$\mu^* = \frac{N_1 + 1}{N + 2}$$

$$\phi_{kd}^* = \frac{N_{kd} + 1}{N_k + 2}$$

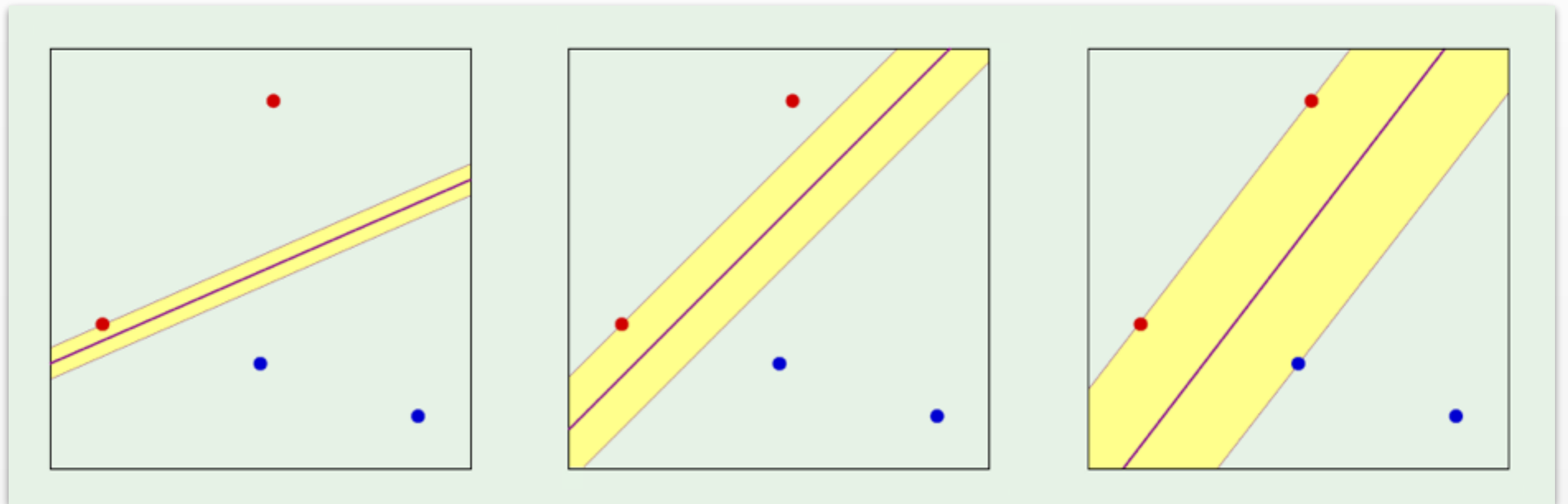
Support Vector Machines

Intuition



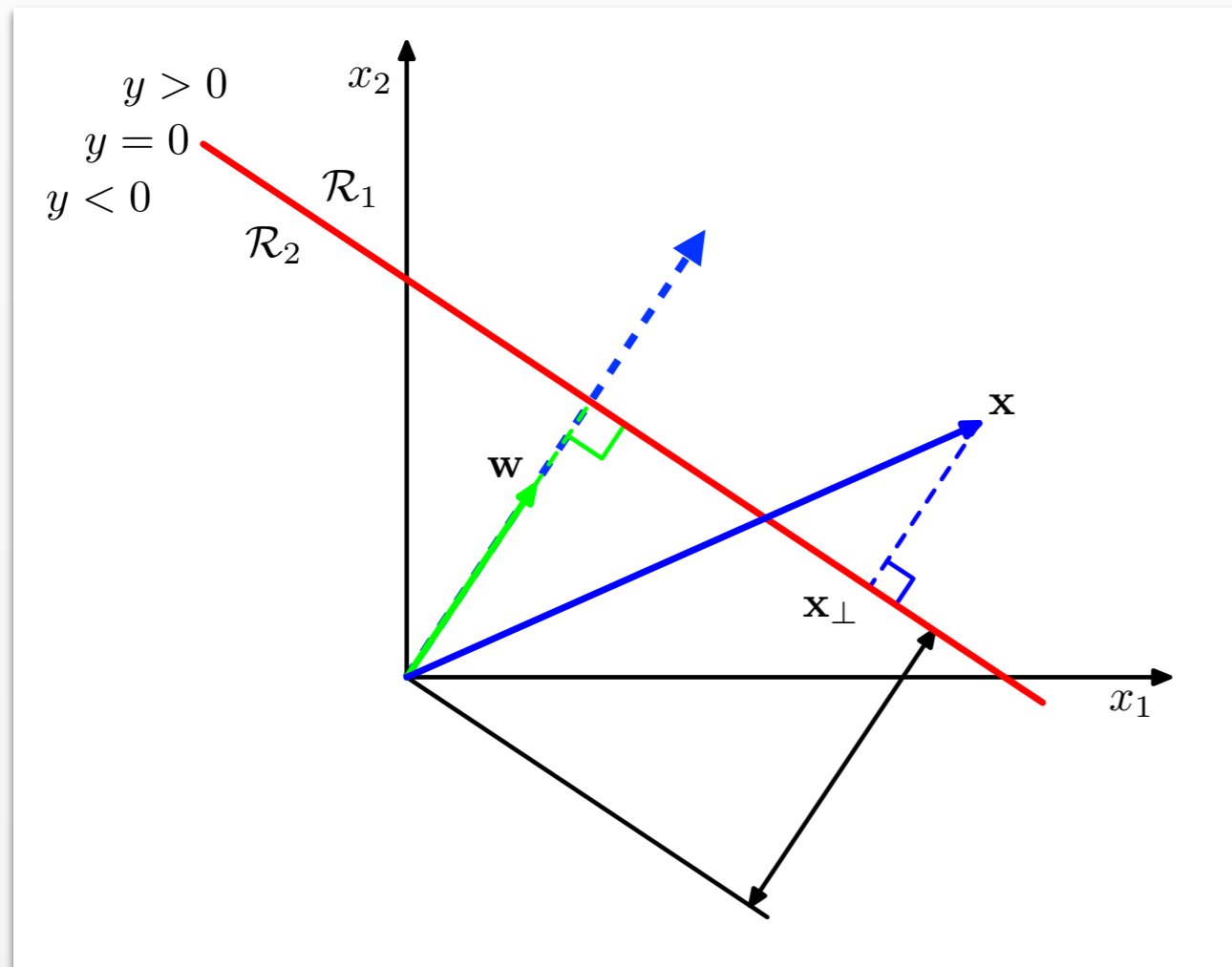
Which of these linear classifiers is the best?

Max Margin Classifiers



Idea: Maximize the *margin* between two *separable* classes

Max Margin Classifiers



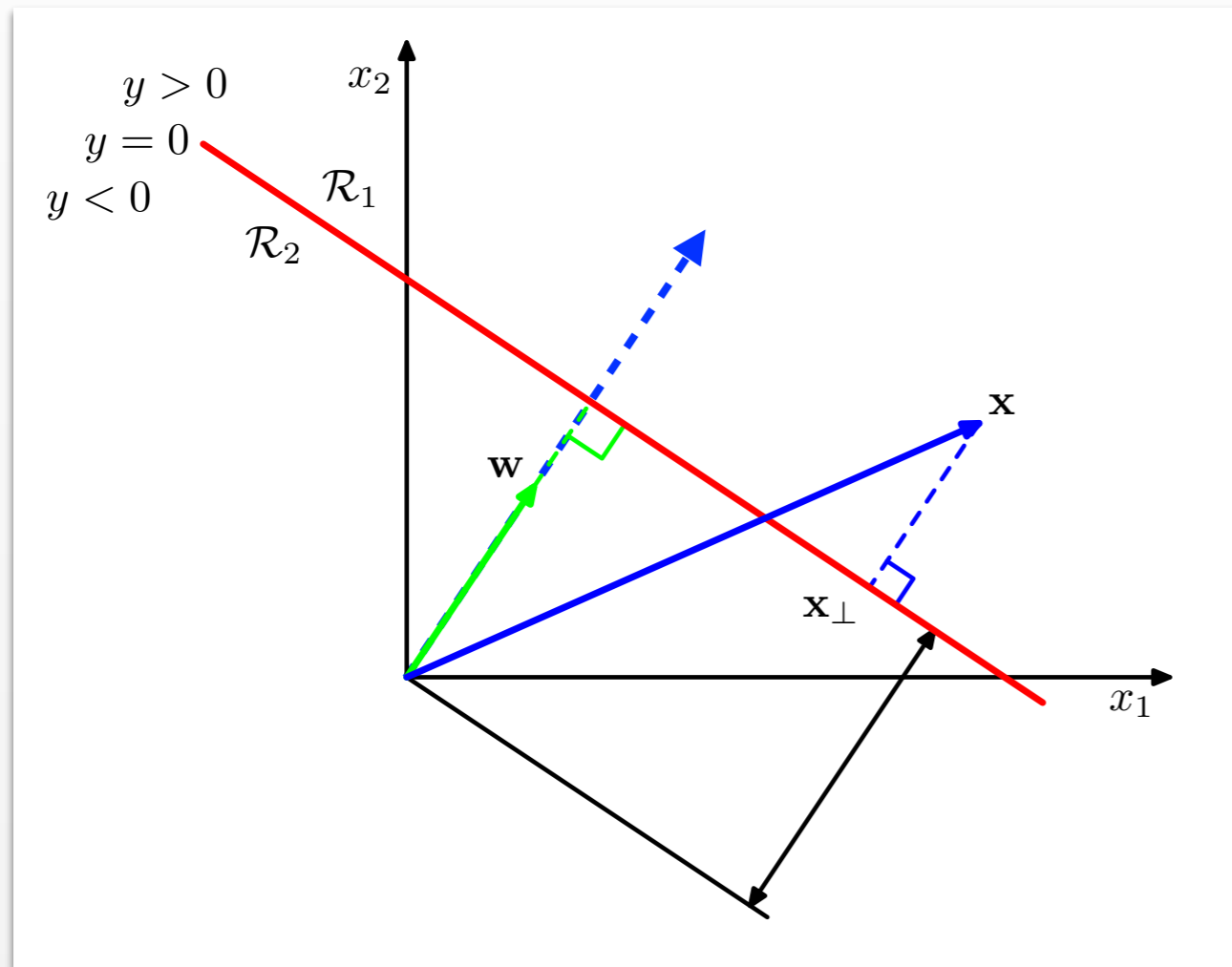
$$h(\mathbf{x}; \mathbf{w}, b) = \text{Sign}(\mathbf{w}^\top \mathbf{x} + b)$$

$$y_n \in \{-1, 1\}$$

$$\mathbf{x} = (x_1, \dots, x_D)$$

$$\mathbf{w} = (w_1, \dots, w_D)$$

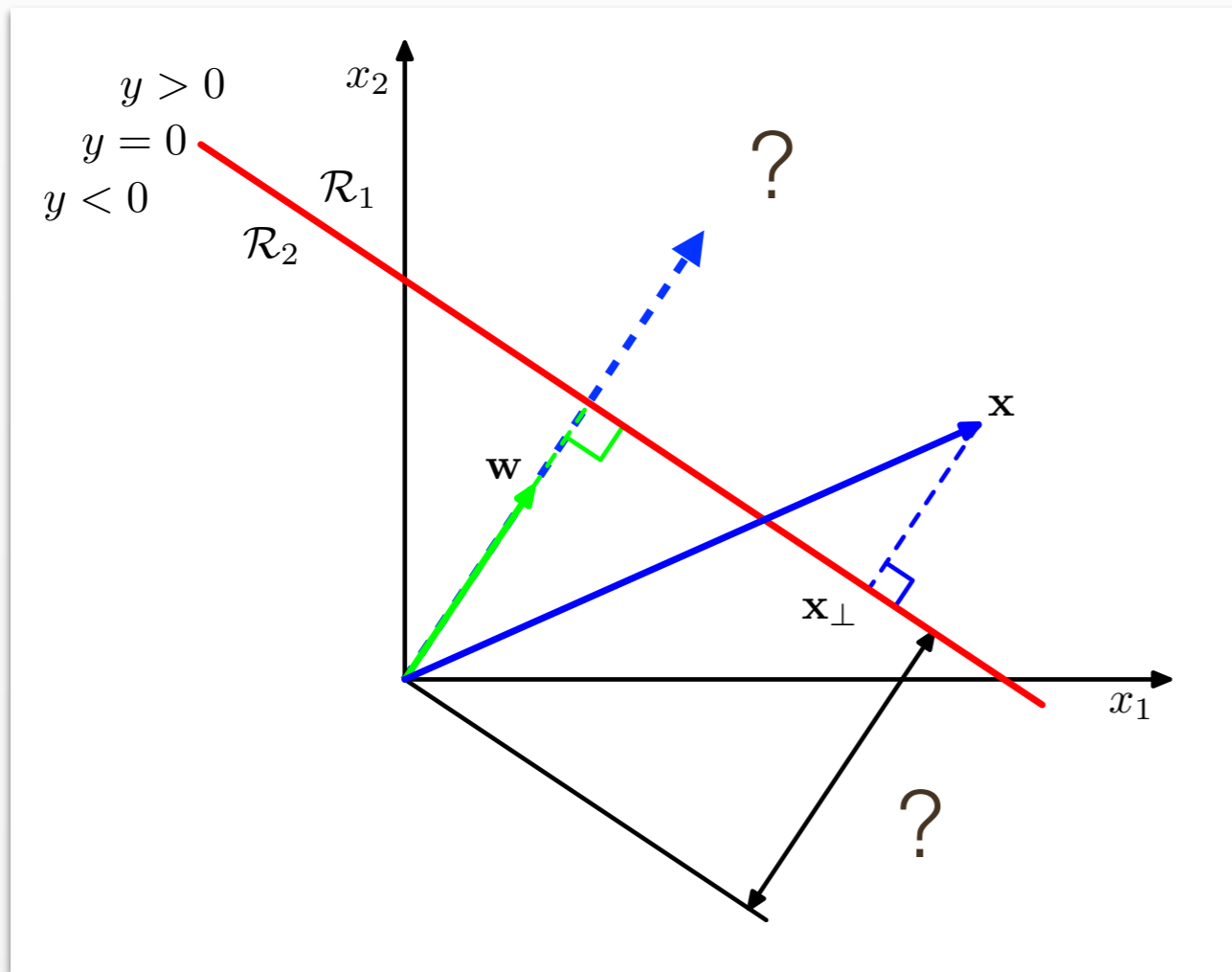
Max Margin Classifiers



$$\mathbf{w}^\top \mathbf{x} + b =$$

$$\|\mathbf{w}\| \left(\frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|} \right)$$

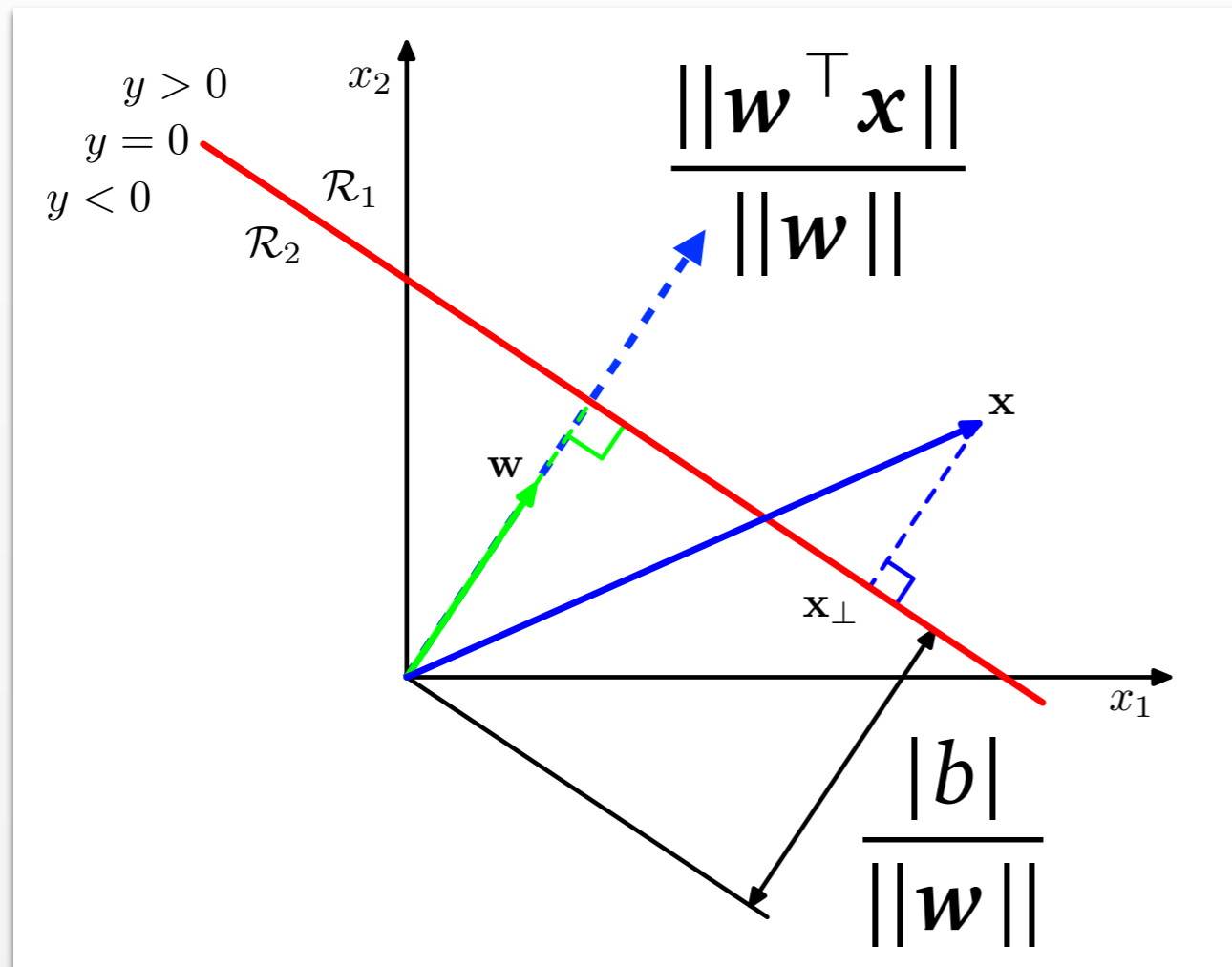
Max Margin Classifiers



$$w^{\top} x + b =$$
$$\|w\| \left(\frac{w^{\top} x}{\|w\|} + \frac{b}{\|w\|} \right)$$

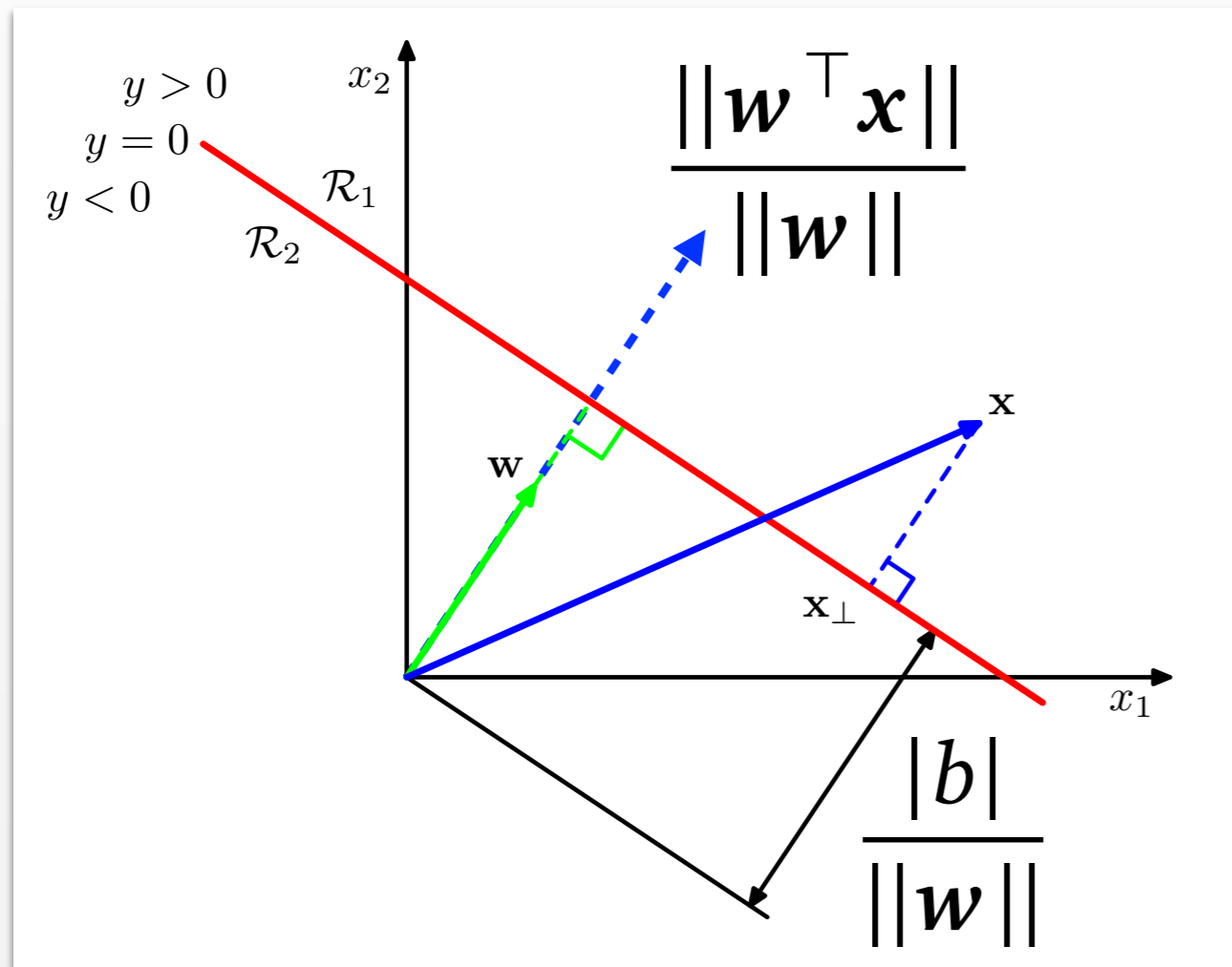
What are the lengths of these vectors?

Max Margin Classifiers



$$\mathbf{w}^\top \mathbf{x} + b = \|\mathbf{w}\| \left(\frac{\mathbf{w}^\top \mathbf{x}}{\|\mathbf{w}\|} + \frac{b}{\|\mathbf{w}\|} \right)$$

Max Margin Classifiers



$$w^{\top} x + b =$$
$$\|w\| \left(\frac{w^{\top} x}{\|w\|} + \frac{b}{\|w\|} \right)$$

Distance from plane: $\frac{1}{\|w\|} (w^{\top} x + b)$

Equivalent Optimization Problems

$$\max_{w, b} \hat{\gamma}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

Distance from plane: $\frac{1}{\|\mathbf{w}\|} (\mathbf{w}^\top \mathbf{x} + b)$

Equivalent Optimization Problems

$$\max_{w, b} \hat{\gamma}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

$$\max_{w, b, \gamma} \frac{\gamma}{\|\mathbf{w}\|}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \gamma \quad n = 1, \dots, N$$

Distance from plane: $\frac{1}{\|\mathbf{w}\|} (\mathbf{w}^\top \mathbf{x} + b)$

Equivalent Optimization Problems

$$\max_{w, b} \hat{\gamma}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

$$\max_{w, b, \gamma} \frac{\gamma}{\|\mathbf{w}\|}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \gamma \quad n = 1, \dots, N$$

$$\max_{w, b} \frac{1}{\|\mathbf{w}\|}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

Distance from plane: $\frac{1}{\|\mathbf{w}\|} (\mathbf{w}^\top \mathbf{x} + b)$

Equivalent Optimization Problems

$$\max_{w,b} \hat{\gamma} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

$$\max_{w,b,\gamma} \frac{\gamma}{\|\mathbf{w}\|} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \gamma \quad n = 1, \dots, N$$

$$\max_{w,b} \frac{1}{\|\mathbf{w}\|} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

$$\min_{w,b} \frac{1}{2} \|\mathbf{w}\|^2 \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

Distance from plane: $\frac{1}{\|\mathbf{w}\|} (\mathbf{w}^\top \mathbf{x} + b)$

Equivalent Optimization Problems

$$\max_{w, b} \hat{\gamma} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

$$\max_{w, b, \gamma} \frac{\gamma}{\|\mathbf{w}\|} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \gamma \quad n = 1, \dots, N$$

$$\max_{w, b} \frac{1}{\|\mathbf{w}\|} \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

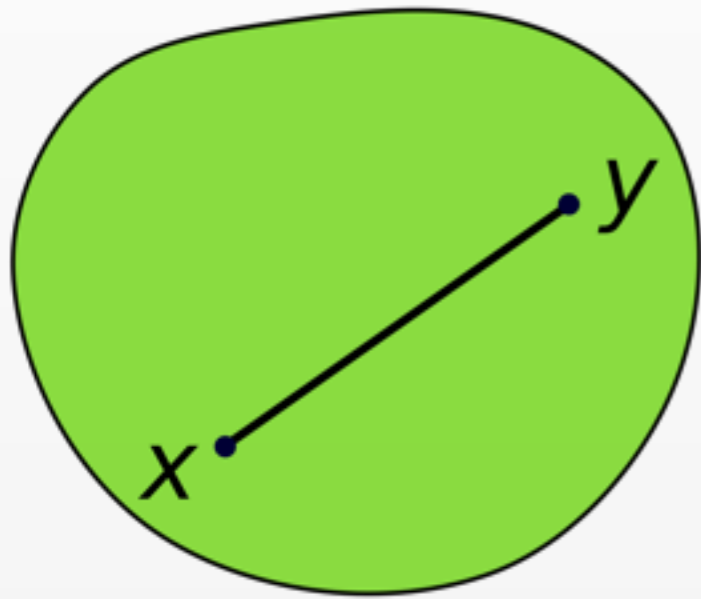
$$\min_{w, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

Distance from plane: $\frac{1}{\|\mathbf{w}\|} (\mathbf{w}^\top \mathbf{x} + b)$

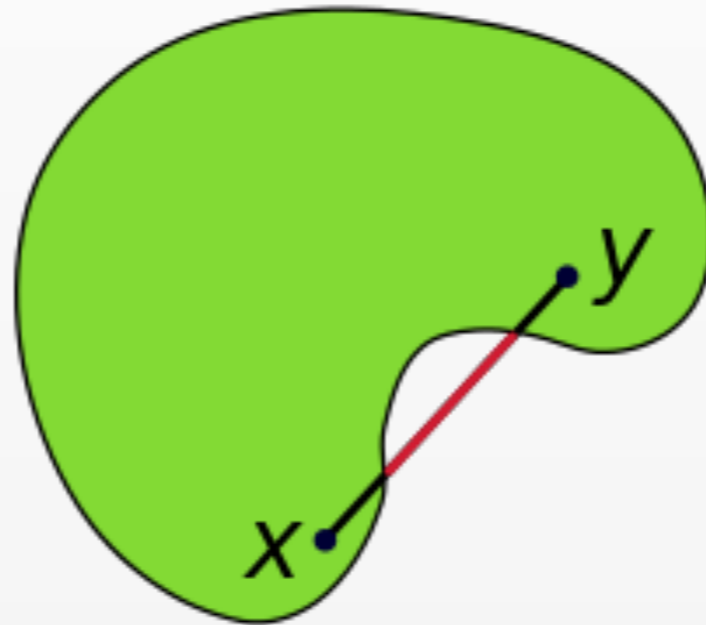
Intermezzo

Convex Optimization

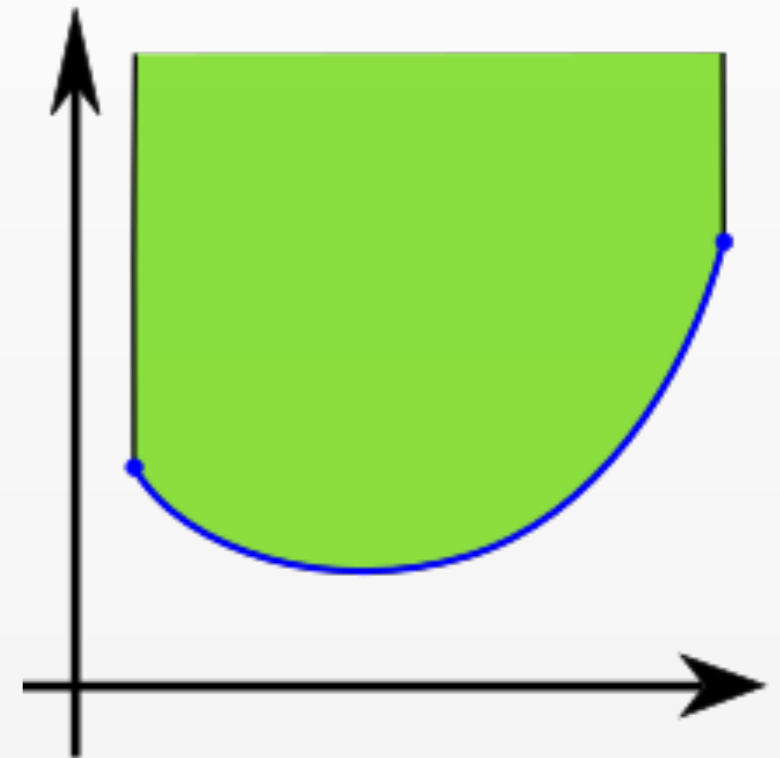
Convex Sets and Functions



Convex
Set



Non-convex
Set



Convex
Function

Lagrange Duality

Constrained Optimization Problem

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

Lagrangian

$$\mathcal{L}(w, \beta) = f(w) + \sum_{i=1}^l \beta_i h_i(w)$$

Optimum

$$\frac{\partial \mathcal{L}}{\partial w_i} = 0; \quad \frac{\partial \mathcal{L}}{\partial \beta_i} = 0,$$

Lagrange Duality

Primal Optimization Problem

$$\begin{aligned} \min_w \quad & f(w) \\ \text{s.t.} \quad & g_i(w) \leq 0, \quad i = 1, \dots, k \\ & h_i(w) = 0, \quad i = 1, \dots, l. \end{aligned}$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Primal Optimization Problem

$$\theta_{\mathcal{P}}(w) = \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta).$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Primal Optimization Problem

$$\theta_{\mathcal{P}}(w) = \begin{cases} f(w) & \text{if } w \text{ satisfies primal constraints} \\ \infty & \text{otherwise.} \end{cases}$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Primal Optimization Problem

$$\min_w \theta_{\mathcal{P}}(w) = \min_w \max_{\alpha, \beta: \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta),$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Dual Optimization Problem

$$\theta_{\mathcal{D}}(\alpha, \beta) = \min_w \mathcal{L}(w, \alpha, \beta).$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Dual Optimization Problem

$$\max_{\alpha, \beta: \alpha_i \geq 0} \theta_{\mathcal{D}}(\alpha, \beta) = \max_{\alpha, \beta: \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta).$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Relationship between Primal and Dual

$$d^* = \max_{\alpha, \beta : \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \quad ? \quad \min_w \max_{\alpha, \beta : \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*.$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Relationship between Primal and Dual

$$d^* = \max_{\alpha, \beta : \alpha_i \geq 0} \min_w \mathcal{L}(w, \alpha, \beta) \leq \min_w \max_{\alpha, \beta : \alpha_i \geq 0} \mathcal{L}(w, \alpha, \beta) = p^*.$$

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Duality gap p^-d^* is zero when*

f convex, g_i convex, h_i affine

Generalized Lagrangian

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

Lagrange Duality

Karush-Kuhn-Tucker (KKT) conditions at optimum

$$\frac{\partial}{\partial w_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, n$$

$$\frac{\partial}{\partial \beta_i} \mathcal{L}(w^*, \alpha^*, \beta^*) = 0, \quad i = 1, \dots, l$$

$$\alpha_i^* g_i(w^*) = 0, \quad i = 1, \dots, k$$

$$g_i(w^*) \leq 0, \quad i = 1, \dots, k$$

$$\alpha^* \geq 0, \quad i = 1, \dots, k$$

Dual complementarity

$g_i(w) = 0$ when $\alpha_i > 0$, $\alpha_i = 0$ when $g_i(w) < 0$

(back to SVMs)

Equivalent Optimization Problems

$$\max_{w, b} \hat{\gamma}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \hat{\gamma} \quad n = 1, \dots, N$$

$$\|\mathbf{w}\| = 1$$

$$\max_{w, b, \gamma} \frac{\gamma}{\|\mathbf{w}\|}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq \gamma \quad n = 1, \dots, N$$

$$\max_{w, b} \frac{1}{\|\mathbf{w}\|}$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

$$\min_{w, b} \frac{1}{2} \|\mathbf{w}\|^2$$

$$y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

SVMs as Convex Optimization

$$\min_{\mathbf{w}, b} \frac{1}{2} \|\mathbf{w}\|^2 \quad y_n(\mathbf{w}^\top \mathbf{x}_n + b) \geq 1 \quad n = 1, \dots, N$$

SVMs as Convex Optimization

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad y_n(w^\top x_n + b) \geq 1 \quad n = 1, \dots, N$$

Write as Convex Optimization Problem

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

SVMs as Convex Optimization

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad y_n(w^\top x_n + b) \geq 1 \quad n = 1, \dots, N$$

Write as Convex Optimization Problem

$$\mathcal{L}(w, \alpha, \beta) = f(w) + \sum_{i=1}^k \alpha_i g_i(w) + \sum_{i=1}^l \beta_i h_i(w).$$

$$f(w) = \frac{1}{2} \|w\|^2$$

$$g_i(w) = -y_i(w^\top x_i + b) + 1 \leq 0$$

SVMs as Convex Optimization

$$\min_{w, b} \frac{1}{2} \|w\|^2 \quad y_n(w^\top x_n + b) \geq 1 \quad n = 1, \dots, N$$

Generalized Lagrangian

$$\mathcal{L}(w, b, \alpha) = \frac{1}{2} \|w\|^2 - \sum_{i=1}^m \alpha_i [y^{(i)}(w^\top x^{(i)} + b) - 1].$$

(note: no equality constraints)

Dual Form

Generalized Lagrangian

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

Dual problem

$$\theta_D(\alpha) = \min_{w, b} \mathcal{L}(w, b, \alpha)$$

Solve for w

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0$$

Solve for b

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i y^{(i)} = 0.$$

Dual Form

Generalized Lagrangian

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

Dual problem

$$\theta_D(\alpha) = ?$$

Solve for w

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0$$

Solve for b

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i y^{(i)} = 0.$$

Dual Form

Generalized Lagrangian

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

Dual problem

$$\theta_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)} - b \sum_{i=1}^m \alpha_i y^{(i)}.$$

Solve for w

$$\nabla_w \mathcal{L}(w, b, \alpha) = w - \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} = 0$$

Solve for b

$$\frac{\partial}{\partial b} \mathcal{L}(w, b, \alpha) = \sum_{i=1}^m \alpha_i y^{(i)} = 0.$$

Dual Form

Generalized Lagrangian

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

Dual problem

$$\theta_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}$$

Dual Form

Generalized Lagrangian

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

Dual problem

$$\theta_D(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j (x^{(i)})^T x^{(j)}$$

Compute w

$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}.$$

Compute b

$$b^* = -\frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}.$$

Support Vectors

Compute w

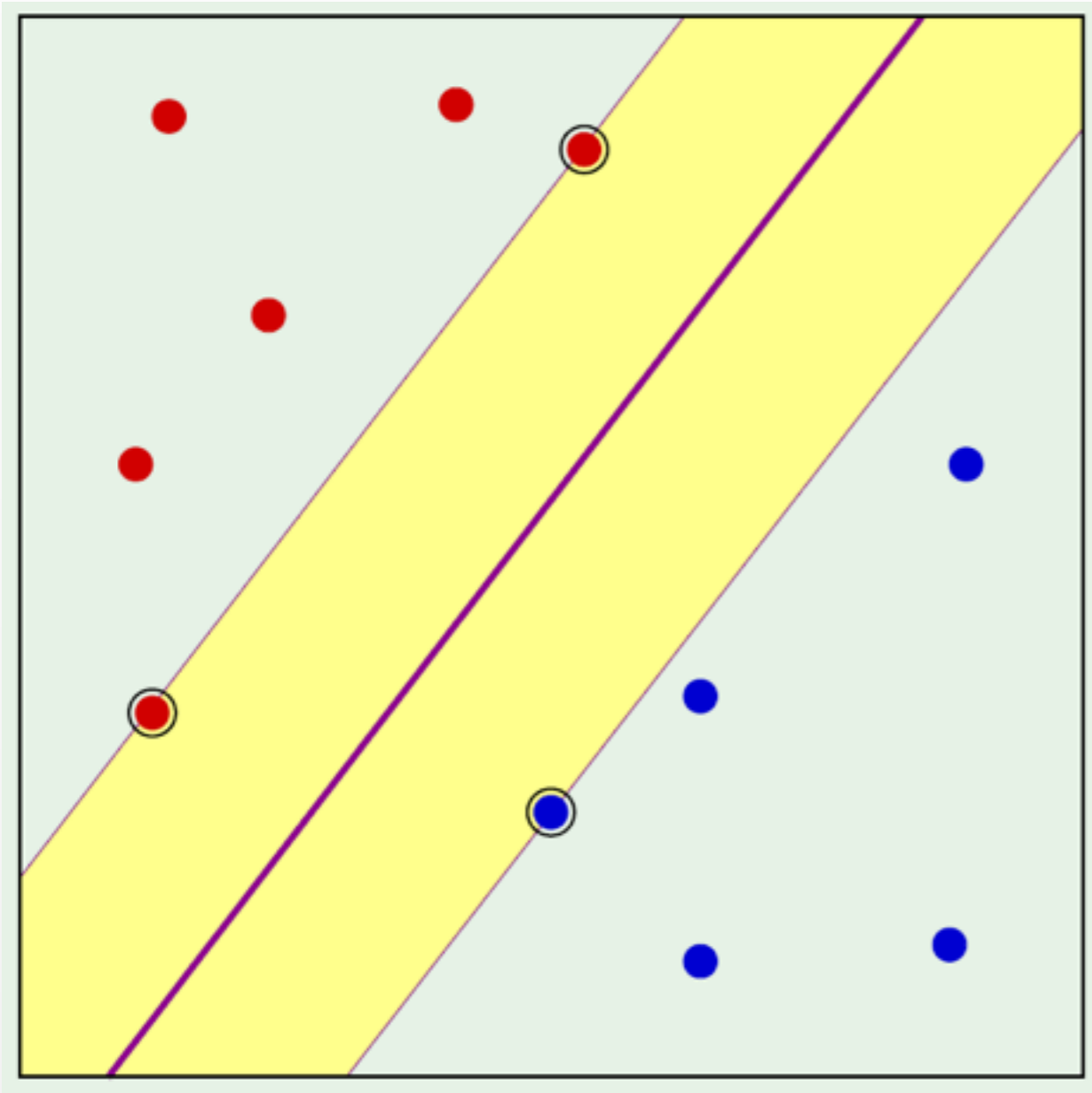
$$w = \sum_{i=1}^m \alpha_i y^{(i)} x^{(i)}.$$

Compute b

$$b^* = -\frac{\max_{i:y^{(i)}=-1} w^{*T} x^{(i)} + \min_{i:y^{(i)}=1} w^{*T} x^{(i)}}{2}.$$

Dual complementarity

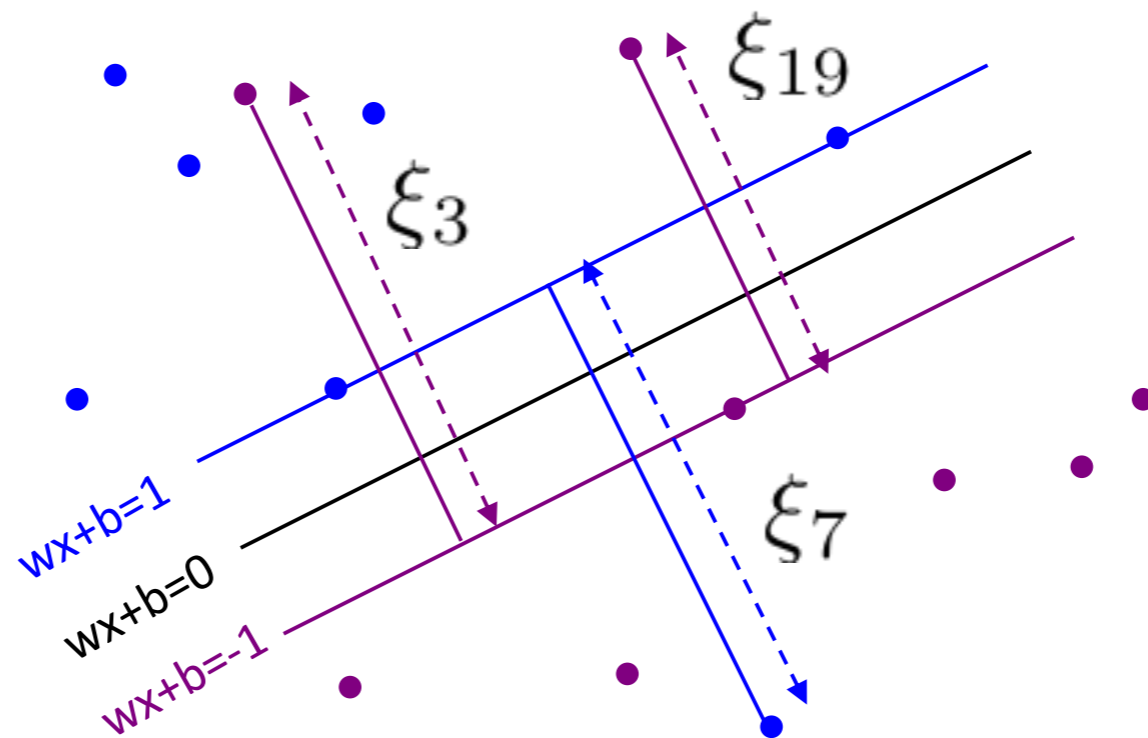
- $\alpha_i = 0$ when $g_i(w) < 0$
- $g_i(w) = 0$ when $\alpha_i > 0$



SVM with Non-Separable Data

$$\arg \min_{\mathbf{w}, b, \xi \geq 0} \quad \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i$$

s.t. $y_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i, \quad i = 1, \dots, m$



SVM with Non-Separable Data

Generalized Lagrangian

$$L(\mathbf{w}, b, \boldsymbol{\xi}, \boldsymbol{\alpha}, \boldsymbol{\lambda}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} + C \sum_{i=1}^m \xi_i + \sum_{i=1}^m \alpha_i (1 - \xi_i - y_i (\mathbf{w}^T \mathbf{x}_i + b)) - \sum_{i=1}^m \lambda_i \xi_i$$

Solve for Dual Form

$$\frac{\partial L}{\partial b} = \sum_{i=1}^m y_i \alpha_i = 0$$

$$\frac{\partial L}{\partial \mathbf{w}} = \mathbf{w} - \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i \Rightarrow \mathbf{w} = \sum_{i=1}^m \alpha_i y_i \mathbf{x}_i$$

$$\frac{\partial L}{\partial \xi} = C - \alpha_i - \lambda_i = 0$$

SVM with Non-Separable Data

Dual Optimization Problem

$$\begin{aligned} \arg \max_{\alpha \geq 0} \quad & \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m \alpha_i \alpha_j y_i y_j \mathbf{x}_i^T \mathbf{x}_j \\ \text{s.t.} \quad & \sum_{i=1}^m \alpha_i y_i = 0 \\ & 0 \leq \alpha_i \leq C, i = 1, \dots, m \end{aligned}$$

Inner Products

Dual Optimization Problem

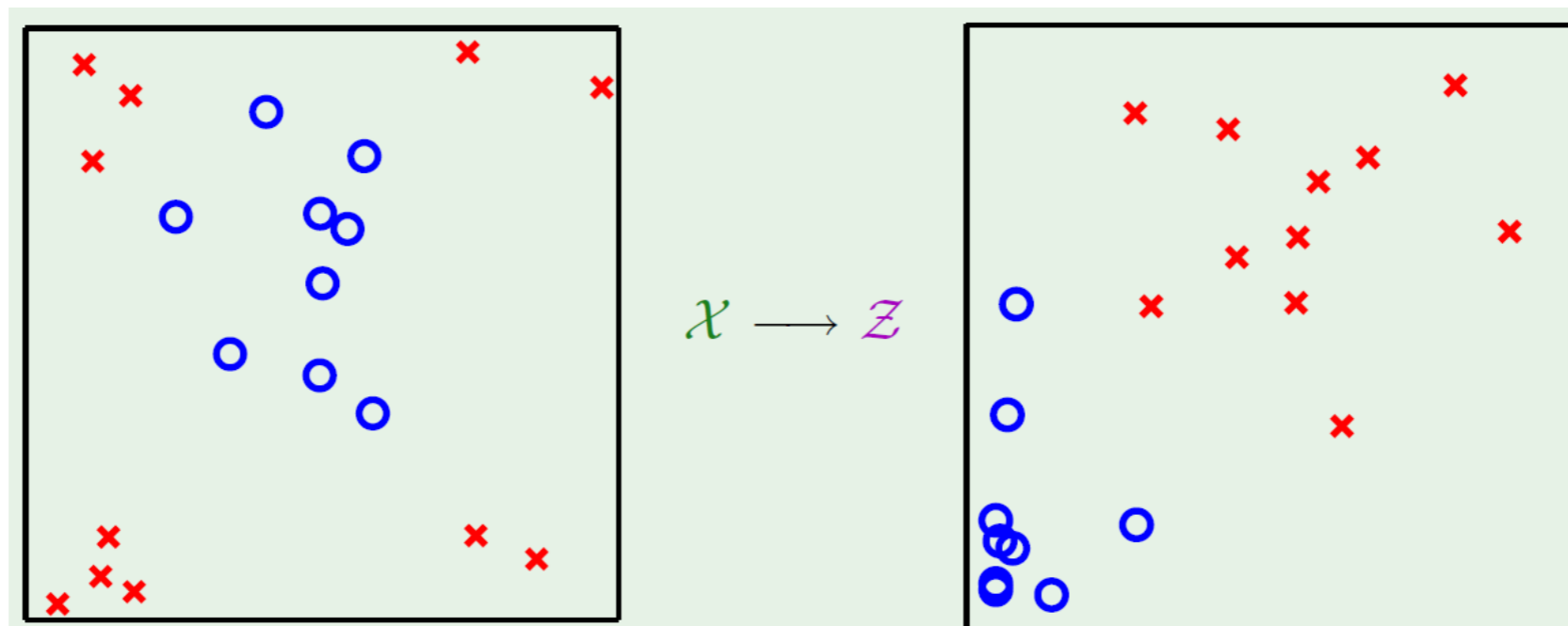
$$\begin{aligned} \max_{\alpha} \quad & W(\alpha) = \sum_{i=1}^m \alpha_i - \frac{1}{2} \sum_{i,j=1}^m y^{(i)} y^{(j)} \alpha_i \alpha_j \langle x^{(i)}, x^{(j)} \rangle. \\ \text{s.t.} \quad & \alpha_i \geq 0, \quad i = 1, \dots, m \\ & \sum_{i=1}^m \alpha_i y^{(i)} = 0, \end{aligned}$$

Prediction

$$\begin{aligned} w^T x + b &= \left(\sum_{i=1}^m \alpha_i y^{(i)} x^{(i)} \right)^T x + b \\ &= \sum_{i=1}^m \alpha_i y^{(i)} \langle x^{(i)}, x \rangle + b. \end{aligned}$$

Next Lecture: Nonlinear SVMs

Let $\mathbf{z} = \Phi(\mathbf{x})$ for some function Φ :



Apply SVM in the \mathbf{z} space by maximizing:

$$L(\alpha) = \sum_{i=1}^N \alpha_i - \frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N y_i y_j \alpha_i \alpha_j \mathbf{z}_i^T \mathbf{z}_j$$