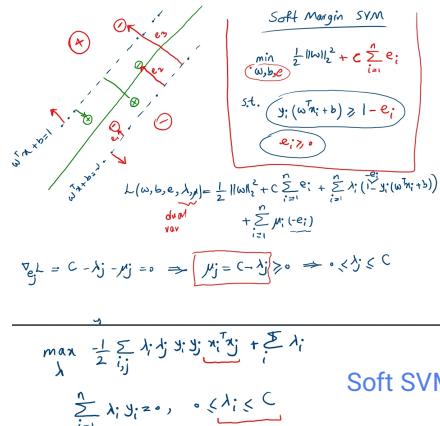
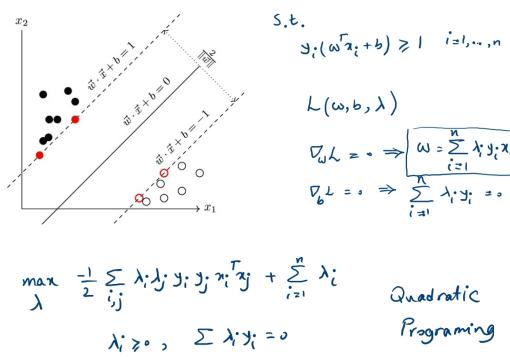


Hard SVM



Techniques for Constructing New Kernels.

Given valid kernels $k_1(x, x')$ and $k_2(x, x')$, the following new kernels will also be valid:

$$k(x, x') = ck_1(x, x') \quad (6.13)$$

$$k(x, x') = f(x)k_1(x, x')f(x') \quad (6.14)$$

$$k(x, x') = q(k_1(x, x')) \quad (6.15)$$

$$k(x, x') = \exp(k_1(x, x')) \quad (6.16)$$

$$k(x, x') = k_1(x, x') + k_2(x, x') \quad (6.17)$$

$$k(x, x') = k_1(x, x')k_2(x, x') \quad (6.18)$$

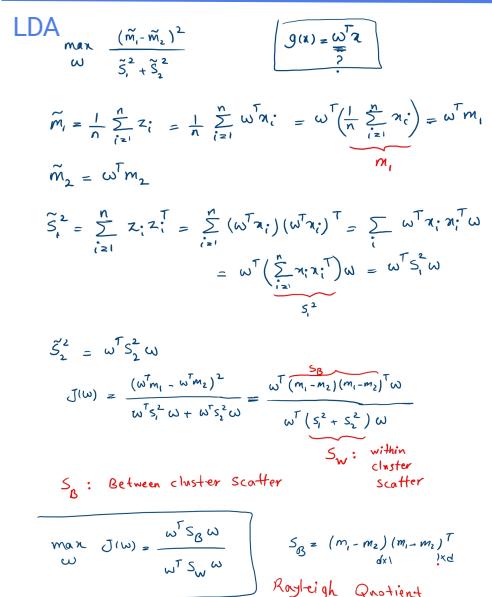
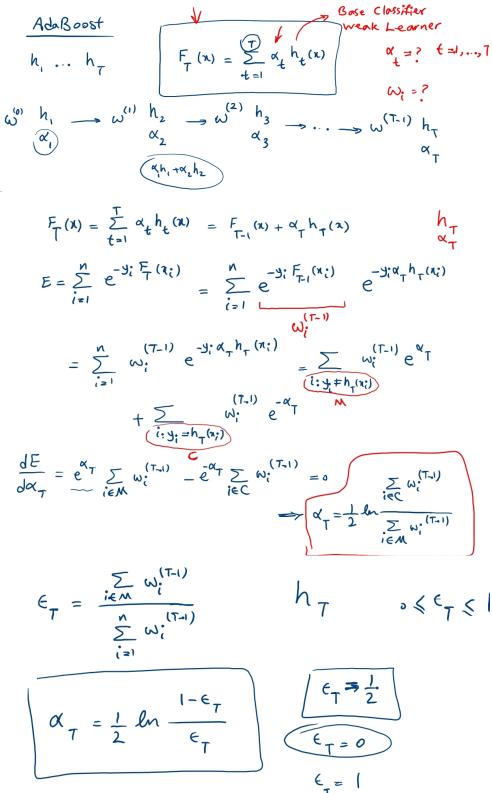
$$k(x, x') = k_3(\phi(x), \phi(x')) \quad (6.19)$$

$$k(x, x') = x^T A x' \quad (6.20)$$

$$k(x, x') = k_a(x_a, x_a') + k_b(x_b, x_b') \quad (6.21)$$

$$k(x, x') = k_a(x_a, x_a')k_b(x_b, x_b') \quad (6.22)$$

where $c > 0$ is a constant, $f(\cdot)$ is any function, $q(\cdot)$ is a polynomial with nonnegative coefficients, $\phi(x)$ is a function from x to \mathbb{R}^M , $k_3(\cdot, \cdot)$ is a valid kernel in \mathbb{R}^M . A is a symmetric positive semidefinite matrix, x_a and x_b are variables (not necessarily disjoint) with $x = (x_a, x_b)$, and k_a and k_b are valid kernel functions over their respective spaces.



Constrained Optimization

$$\min_{\omega} f(\omega)$$
s.t.
 $y_i(\omega) \leq 0 \quad i=1, \dots, m$
 $h_i(\omega) = 0 \quad i=1, \dots, p$

KKT conditions

$$\text{Lagrangian: } L(\omega, \lambda, \mu) = f(\omega) + \sum_{i=1}^m \lambda_i y_i(\omega) + \sum_{i=1}^p \mu_i h_i(\omega)$$

$$\lambda_i \geq 0 \quad i=1, \dots, m$$

KKT:

$$(1) \quad \nabla_\omega L(\omega, \lambda, \mu) = 0$$

$$(2) \quad h_i(\omega) = 0 \quad (y_i(\omega) \leq 0)$$

$$(3) \quad \lambda_i \geq 0$$

$$(4) \quad \lambda_i y_i(\omega) = 0 \rightarrow \text{Complementary Slackness}$$

$$g(\mu, \lambda) = \min_{\omega} L(\omega, \mu, \lambda)$$

$$g(\mu, \lambda) \leq f(\omega^*)$$

Dual solution \leq primal solution

$$g(\omega^*, \mu^*) \leq f(\omega^*)$$

weak duality

Strong Duality:

$$g(\omega^*, \mu^*) = f(\omega^*)$$

$$\begin{cases} (1) \text{ convex} \\ (2) \exists \omega \quad g_i(\omega) < 0 \end{cases}$$

$$\min_{\omega, b} \frac{1}{2} \|\omega\|_2^2$$
s.t.
 $y_i(\omega^T x_i + b) \leq 1 \quad i=1, \dots, n$

$$L(\omega, b, \lambda) = \frac{1}{2} \|\omega\|_2^2 + \sum_{i=1}^n \lambda_i (1 - y_i(\omega^T x_i + b))$$

$$\frac{\partial L}{\partial \omega} = \omega + \sum_{i=1}^n \lambda_i y_i x_i = 0 \Rightarrow \boxed{\omega = \sum_{i=1}^n \lambda_i y_i x_i}$$

$$\frac{\partial L}{\partial b} = -\sum_{i=1}^n \lambda_i y_i = 0 \Rightarrow \boxed{\sum_{i=1}^n \lambda_i y_i = 0}$$

$$g(\lambda) = \frac{1}{2} \left(\sum_i \lambda_i y_i x_i \right)^T \left(\sum_i \lambda_i y_i x_i \right) + \sum_{i=1}^n \lambda_i$$

$$- \sum_i \lambda_i y_i \left(\sum_j \lambda_j y_j x_j \right)^T x_i - \sum_i \lambda_i y_i$$

$$= \frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j - \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_i \lambda_i$$

$$\Rightarrow g(\lambda) = -\frac{1}{2} \sum_i \sum_j \lambda_i \lambda_j y_i y_j x_i^T x_j + \sum_i \lambda_i$$

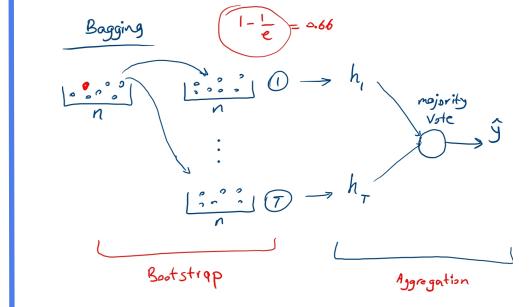
$$\max_{\lambda} g(\lambda)$$
s.t.
 $\lambda_i \geq 0$

$$\Rightarrow \lambda^* = \begin{bmatrix} \lambda_1^* \\ \vdots \\ \lambda_n^* \end{bmatrix}$$

$$\omega^* = \sum_{i=1}^n \lambda_i^* y_i x_i$$

Complementary Slackness:
 $\lambda_i^* y_i(\omega^*) = 0$

$$\lambda_k^* (y_k(\omega^T x_k + b)) = 0$$



$$\max_{\omega} \text{var}(z)$$
s.t.
 $\|\omega\|_2^2 = 1$

$E[z] = 0 \quad \checkmark$

$$\text{PCA}$$

$$z = u^T x \Rightarrow E[z] = E[u^T x] = u^T E[x] = 0$$

$$\text{var}(z) = E[z^2] - E[z]^2 \simeq \frac{1}{n} \sum_{i=1}^n z_i^2$$

$$\max_{\omega} \frac{1}{n} \sum_{i=1}^n z_i^2 = \frac{1}{n} \sum_{i=1}^n (x_i^T \omega)^2 = \frac{1}{n} \sum_{i=1}^n x_i^T \omega \omega^T x_i$$
s.t.
 $\|\omega\|_2^2 = 1$

$$(u^T \omega)(\omega^T u)$$

$$= u^T \left(\frac{1}{n} \sum_{i=1}^n x_i x_i^T \right) u$$

$S: \text{Sample cov. Matrix}$

$$= u^T S u$$

$$\max_{\omega} u^T S u$$
s.t.
 $\|\omega\|_2^2 = 1$

$$L(\omega) = u^T S u - \lambda (\|\omega\|_2^2 - 1)$$

$$\nabla_\omega L = 0 \Rightarrow 2 S u - \lambda 2 \omega = 0$$

$$\Rightarrow S u = \lambda \omega$$

eigenvector eigenvalue

$$\max_{\omega} u^T \omega = \max_{\omega} \frac{\omega^T u}{\|\omega\|_2^2} = \max_{\omega} \frac{1}{\|\omega\|_2^2}$$

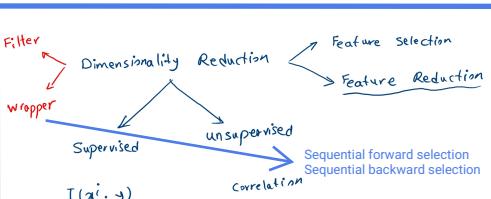
$$\max_{\omega} (\omega^T S_B \omega)$$
s.t.
 $\omega^T S_W \omega = k$

$$\max_{\omega} \omega^T S_B \omega - \lambda (\omega^T S_W \omega - k)$$

$$\lambda \gg 0$$

$$2 S_B \omega - 2 \lambda S_W \omega = 0 \Rightarrow S_B \omega = \lambda S_W \omega$$

$$\boxed{S_B^{-1} S_W \omega = \lambda \omega}$$



$$\lambda \omega = S_W^{-1} S_B \omega$$

$$S_B = \frac{(m_1 - m_2)(m_1 - m_2)^T}{n}$$

$$\lambda \omega = S_W^{-1} (m_1 - m_2)(m_1 - m_2)^T \omega \Rightarrow \omega = \frac{\lambda}{n} S_W^{-1} (m_1 - m_2)$$

$$\boxed{\omega = \frac{\lambda}{n} S_W^{-1} (m_1 - m_2)}$$

Types of Clustering



K-means Recap ...



What is K-means optimizing? (Objective Function)

- Randomly initialize k centers

$$\bullet \mu^{(0)} = \mu_1^{(0)}, \mu_2^{(0)}, \mu_3^{(0)}, \dots, \mu_k^{(0)}$$

- Interate t=0, 1, 2, ...

- **Classify:** Assign each point $j \in \{1, 2, \dots, m\}$ to nearest center

$$C^{(t)}(j) \leftarrow \arg \min_{i=1, \dots, k} \|\mu_i^{(t)} - x_j\|^2$$

- **Recenter:** μ_i becomes centroid if its points:

$$\mu_i^{(t+1)} \leftarrow \arg \min_{\mu} \sum_{j: C^{(t)}(j)=i} \|\mu - x_j\|^2 \quad i \in \{1, \dots, k\}$$

- Equivalent to μ_i average of its points!

- **Potential function** (objective function) $F(\mu, C)$ of centers μ and point allocation C :

$$\begin{aligned} F(\mu, C) &= \sum_{j=1}^m \|x_j - \mu_{C(j)}\|^2 \\ &= \sum_{i=1}^k \sum_{j: C(j)=i} \|\mu_i - x_j\|^2 \end{aligned}$$

- Optimal K-means:

$$\bullet \min_{\mu} \min_C F(\mu, C)$$

- A clustering is a process to find a set of clusters
- Important distinction between **hierarchical** and **partitional** sets of clusters
- **Partitional Clustering** \leftarrow K-means
 - A division data objects into non-overlapping subsets (clusters) such that each data object is in exactly one subset
- **Hierarchical clustering** \leftarrow
 - A set of nested clusters organized as a hierarchical tree
- **Density based clustering** \rightarrow EM
 - Discover clusters of **arbitrary** shape.
 - Clusters **dense regions** of objects separated by regions of low density

EM / GMM

$$\hat{\theta} = \arg \max_{\theta} E_Z \left[\sum_{i=1}^n \log P(x_i; z_i; \theta) \right]$$

complete Log Likehood $L(\theta)$

$$Q(\theta) = E_Z [L(\theta)]$$

E-step

$$\frac{\partial Q(\theta)}{\partial \theta} = 0 \quad \max_{\theta} Q(\theta) \quad M\text{-step}$$

$$P(z| \theta) = \sum_{k=1}^K \alpha_k N(z | \mu_k, \Sigma_k) \quad \sum_{k=1}^K \alpha_k = 1$$

$$\alpha_k \geq 0$$

$$D = \{x_1, x_2, x_3, \dots, x_n\}$$

$$\rightarrow H = \{z_1, z_2, z_3, \dots, z_n\}$$

$\boxed{z_i = \begin{bmatrix} \vdots \\ 1 \\ \vdots \end{bmatrix} \xrightarrow{k} k}$

$$P(x_i, z_i | \theta) = \prod_{k=1}^K (\alpha_k N(x_i | \mu_k, \Sigma_k))^{z_{ik}}$$

$$= \left(\sum_{k=1}^K z_{ik} \alpha_k N(x_i | \mu_k, \Sigma_k) \right)^{-1}$$

$$L(\theta) = \log P(D, H | \theta) = \log P(x_1, z_1, x_2, z_2, \dots, x_n, z_n | \theta) = \log \prod_{i=1}^n P(x_i, z_i | \theta) = \sum_{i=1}^n \log P(x_i, z_i | \theta)$$

$$= \sum_{i=1}^n \sum_{k=1}^K z_{ik} (\log \alpha_k + \log N(x_i | \mu_k, \Sigma_k))$$

$$Q(\theta) = E_{Z_i} [L(\theta)] = \sum_{i=1}^n \sum_{k=1}^K \underbrace{E[z_{ik}]}_{\gamma_{ik}^t} (\log \alpha_k + \log N(x_i | \mu_k, \Sigma_k))$$

$$\begin{aligned} E[z_{ik}] &= \Pr(z_{ik}=1 | x_i, \theta^t) \\ &= \frac{P(x_i | z_{ik}=1, \theta^t) P(z_{ik}=1 | \theta^t)}{P(x_i | \theta^t)} \\ &= \frac{N(x_i | \mu_k, \Sigma_k)}{\sum_{k=1}^K \alpha_k N(x_i | \mu_k, \Sigma_k)} = \gamma_{ik}^t \end{aligned}$$

$$\frac{\partial Q(\theta)}{\partial \alpha_k} = \hat{\alpha}_k = \frac{\sum_{i=1}^n \gamma_{ik}^t}{n}$$

$$\hat{\mu}_k = \frac{\sum_{i=1}^n \gamma_{ik}^t x_i}{\sum_{i=1}^n \gamma_{ik}^t}$$

$$Q(\theta) = \sum_{i=1}^n \left(\sum_{k=1}^K \gamma_{ik}^t (\log \alpha_k + \log N(x_i | \mu_k, \Sigma_k)) \right) - \lambda \left(\sum_{k=1}^K \alpha_k - 1 \right)$$

$$\frac{\partial Q(\theta)}{\partial \alpha_j} = \sum_{i=1}^n \gamma_{ij}^t \frac{1}{\alpha_j} - \lambda = 0 \Rightarrow \alpha_j = \frac{n \gamma_{ij}^t}{\lambda} \quad j = 1, \dots, K$$

$$1 = \sum_{j=1}^K \alpha_j = \frac{\sum_{j=1}^K \sum_{i=1}^n \gamma_{ij}^t}{\lambda}$$

$$\alpha_j = \frac{\sum_{i=1}^n \gamma_{ij}^t}{\sum_{j=1}^K \sum_{i=1}^n \gamma_{ij}^t} = \frac{\sum_{i=1}^n \gamma_{ij}^t}{n}$$