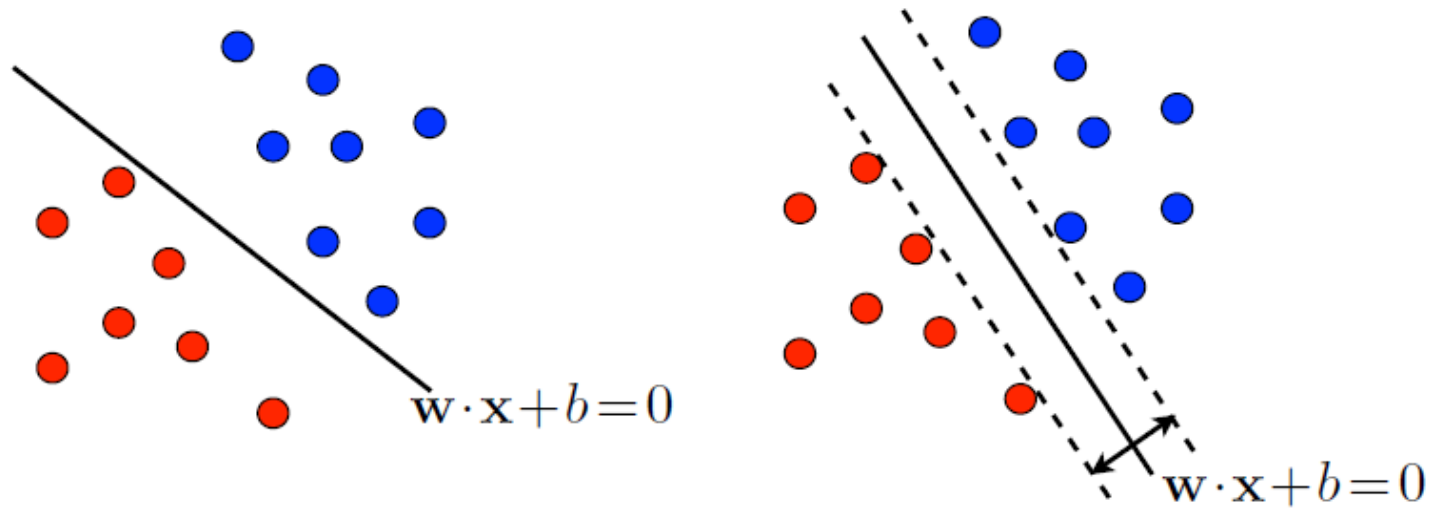


Maximum Margin Clustering

NIPS 2004

SVM Classification



- Maximum Margin Classification

SVM Classification

$$\begin{aligned}\gamma^{*-2} &= \min_{\mathbf{w}, b, \epsilon} \|\mathbf{w}\|^2 + C\epsilon^T \mathbf{e} \quad \text{subject to} \quad y^i(\mathbf{w}^T \phi(\mathbf{x}^i) + b) \geq 1 - \epsilon_i, \forall_{i=1}^N, \epsilon \geq 0 \\ &= \max_{\lambda} 2\lambda^T \mathbf{e} - \langle K \circ \lambda \lambda^T, \mathbf{y} \mathbf{y}^T \rangle \quad \text{subject to} \quad 0 \leq \lambda \leq C, \lambda^T \mathbf{y} = 0 \quad (3)\end{aligned}$$

$$\langle a, b \rangle = \sum_{ij} a_{ij} b_{ij}$$

$$\lambda^T (K \circ \mathbf{y} \mathbf{y}^T) \lambda = \langle K \circ \mathbf{y} \mathbf{y}^T, \lambda \lambda^T \rangle = \langle K \circ \lambda \lambda^T, \mathbf{y} \mathbf{y}^T \rangle$$

- Supervised maximum margin training
- This is an standard QP problem:
 - active sets methods
 - interior point method

Problem Definition

- Given data x_1, x_2, \dots, x_n , we wish to assign the data points to two classes $\{-1, +1\}$ such that separation between the two classes is as wide as possible

$$\min_{y \in \{-1, +1\}^N} \max_{\lambda} 2\lambda^T e - \langle K \circ \lambda \lambda^T, yy^T \rangle \quad \text{subject to} \quad 0 \leq \lambda \leq C \quad \lambda^T y = 0$$

- Integer programming
- May leads to trivial solution, highly unbalanced clusters
- Not a convex function of y !

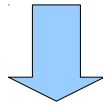
Problems & solutions

$$\min_{\mathbf{y} \in \{-1, +1\}^N} \quad \max_{\lambda} 2\lambda^T \mathbf{e} - \langle K \circ \lambda \lambda^T, \mathbf{y} \mathbf{y}^T \rangle \quad \text{subject to} \quad 0 \leq \lambda \leq C \quad \lambda^T \mathbf{y} = 0$$

- Unbalanced cluster problem
 - Impose a constraint on $-\ell \leq \mathbf{e}^T \mathbf{y} \leq \ell$
- Integer programming
 - Soft clustering
- Non-convexity
 - Set $b=0$ to drop the constraint $\lambda^T \mathbf{y} = 0$
 - Centering the data at the origin and more...

Re-express optimization problem

$$\gamma^{*-2}(y) = \max_{\lambda} 2\lambda^T e - \langle K \circ \lambda\lambda^T, yy^T \rangle \quad \text{subject to} \quad 0 \leq \lambda \leq C$$

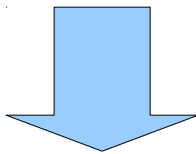


$$\gamma^{*-2}(M) = \max_{\lambda} 2\lambda^T e - \langle K \circ \lambda\lambda^T, M \rangle \quad \text{subject to} \quad 0 \leq \lambda \leq C \quad M = yy^T$$

- Objective function is linear 😊
- $M=yy'$ and M is $[-1,+1]^n$ is not convex 😞

Indirectly enforce $M=yy^T$

- M encodes equivalence relation
 - transitive, reflexive and symmetric
- M has two equivalence classes

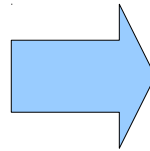


$$\mathcal{L}_1: m_{ii} = 1; m_{ij} = m_{ji}; m_{ik} \geq m_{ij} + m_{jk} - 1; \forall_{ijk}$$

$$\mathcal{L}_2: m_{jk} \geq -m_{ij} - m_{ik} - 1; \forall_{ijk}$$

$$\mathcal{L}_3: \sum_i m_{ij} \leq N - 2; \forall_j$$

$$\mathcal{L}_4: -\ell \leq \sum_i m_{ij} \leq \ell; \forall_j$$



$$\min_{M \in [-1, +1]^{N \times N}} \max_{\lambda} 2\lambda^T e - \langle K \circ \lambda \lambda^T, M \rangle$$

subject to $0 \leq \lambda \leq C, \mathcal{L}_1, \mathcal{L}_2, \mathcal{L}_4, M \succeq 0$

$$y = \sqrt{\lambda_1} v_1 \quad \text{😊}$$

Semi-supervised learning

- Combine both unlabeled data and labeled data to produce a more accurate classification result

$$\mathcal{S}_1: m_{ij} = y_i y_j \quad \text{for labeled examples } i, j \in \{1, \dots, n\}$$

$$\mathcal{S}_2: \sum_{i=1}^n m_{ij} \geq 2 - n \quad \text{for unlabeled examples } j \in \{n + 1, \dots, N\}$$

Experiment Results

| | Gaussians | Circles | A I | Joined Circles | Digits | Faces |
|---------------------|-----------|---------|------|----------------|--------|-------|
| Maximum Margin | 1.25 | 0 | 0 | 1 | 3 | 0 |
| Spectral Clustering | 1.25 | 0 | 0 | 24 | 6 | 16.7 |
| K-means | 5 | 50 | 38.5 | 50 | 7 | 24.4 |

Table 1: Percentage misclassification errors of the various clustering algorithms on the various data sets.

| | HWD 1-7 | HWD 2-3 | UCI Austra. | UCI Flare | UCI Vote | UCI Diabet. |
|------------|---------|---------|-------------|-----------|----------|-------------|
| Max Marg | 3.3 | 4.7 | 32 | 34 | 14 | 35.55 |
| Spec Clust | 4.2 | 6.4 | 48.7 | 40.7 | 13.8 | 44.67 |
| TSVM | 4.6 | 5.4 | 38.7 | 33.3 | 17.5 | 35.89 |
| SVM | 4.5 | 10.9 | 37.5 | 37 | 20.4 | 39.44 |

Table 2: Percentage misclassification errors of the various semisupervised learning algorithms on the various data sets. SVM uses no unlabeled data. TSVM is due to [8].

My Comments

- Class balance problem
- For pure clustering problem, soft margin is not necessary. It leads to C value can be arbitrary, which is consistent with experiment report by the author.