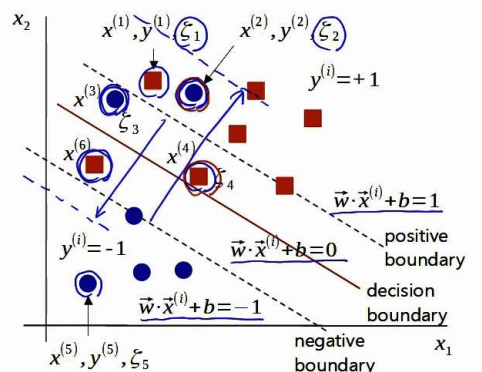# Soft Margin Revisited

Thursday 23 May 2024        8:22 AM

- **Linear Soft Margin – Slack variable**
  - If the classes +1 and -1 are mixed and cannot be completely separated linearly, they can be linearly separated by allowing for some errors. It is called linear soft margin SVM.
  - We assign a slack variable (ξ) to each data point to allow for misclassification.
  - Our goal is to find a decision boundary that minimizes the total error = $\sum \xi_i$, (ξ ≥ 0) while maximizing the margin.

  - **Constraints**



$$\vec{w}\cdot\vec{x}^{(i)}+b \geq 1-\zeta_i \quad if \ y^{(i)}=+1$$
$$\vec{w}\cdot\vec{x}^{(i)}+b \leq -1+\zeta_i \quad if \ y^{(i)}=-1$$
$$y^{(i)}\left(\vec{w}\cdot\vec{x}^{(i)}+b\right) \geq 1-\zeta_i$$
$$\zeta_i \geq 1-y^{(i)}\left(\vec{w}\cdot\vec{x}^{(i)}+b\right)$$

For correctly classified data, x(1), x(4), x(5), ξ is greater than or equal to 0. For misclassified data, x(2), x(3), x(6), ξ is greater than or equal to 1. The data point x(4) that is classified correctly but falls between the positive and negative boundaries will have a ξ value between 0 and 1. Therefore ξ can be considered a measure of misclassification.

$$\boxed{\zeta_i = max\left(0, 1-y^{(i)}\left(\vec{w}\cdot\vec{x}^{(i)}+b\right)\right)}$$

$$\zeta_1 = max\left(0, 1-(\vec{w}\cdot\vec{x}^{(1)}+b)\right) = 0$$
$$\quad (>1)$$
$$\zeta_2 = max\left(0, 1+(\vec{w}\cdot\vec{x}^{(2)}+b)\right) > 2$$
$$\quad (>1)$$
$$\zeta_3 = max\left(0, 1+(\vec{w}\cdot\vec{x}^{(3)}+b)\right) = 1\sim2$$
$$\quad 0\sim1$$
$$\zeta_4 = max\left(0, 1-(\vec{w}\cdot\vec{x}^{(4)}+b)\right) = 0\sim1$$
$$\quad 0\sim1$$
$$\zeta_5 = max\left(0, 1+(\vec{w}\cdot\vec{x}^{(5)}+b)\right) = 0$$
$$\quad (\leq-1)$$
$$\zeta_6 = max\left(0, 1-(\vec{w}\cdot\vec{x}^{(6)}+b)\right) = 1\sim2$$
$$\quad (-1\sim0)$$

- **Objective function**
  - Soft margin is a method that allows for some errors. Two objective functions must be considered together. The first is to maximize margin, and the second is to minimize errors. The two goals are trade-offs and must be properly balanced.
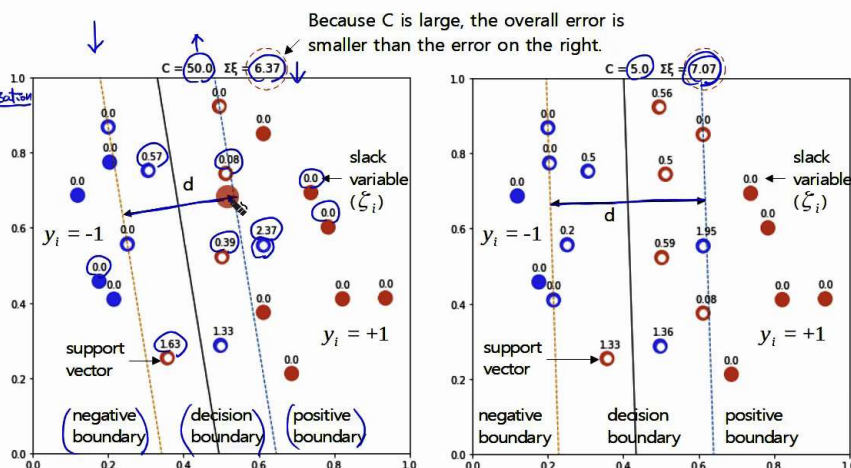  - The two objective functions are combined into one by the weight, C.

  1) Goal-1: $max(d)=min\frac{1}{2}\|\vec{w}\|^2$     2) Goal-2: $min\sum_{i=1}^{N}\zeta_i = min\sum_{i=1}^{N}max\left(0, 1-y^{(i)}\left(\vec{w}\cdot\vec{x}^{(i)}+b\right)\right)$ → hinge loss

  **Objective function:** $min\left(\frac{1}{2}\|\vec{w}\|^2\right) + C\sum_{i=1}^{N}\zeta_i^{(k)}$
  loss      regularization
  Ridge regularization    loss

  - If C is small, the margin is made large even if the error is large. If C is large, the error is made small even if the margin is small.
  - Typically k=1, 2 is used.
    (k = 1 : hinge loss, k = 2 : squared hinge loss).
  - The first term can be viewed as a loss and the second term as a regularized term (penalty term).
  - Conversely, the second term can be viewed as a loss (hinge loss) and the first term as a regularized term (L2, Ridge).

  Because C is large, the overall error is smaller than the error on the right.



- **Hinge loss**
  - Slack variables, ξ, can be considered errors.
  - In regression, the further a data point is from the regression curve, above or below, the larger the error. However, in classification, the error increases the further a data point is from the boundary in the wrong direction, but the error is zero no matter how far the data point is in the right direction. This is called hinge loss.
  - The typical form of hinge loss is something like max(0, something).

  $$\hat{y}^{(i)}=\vec{w}\cdot\vec{x}^{(i)}+b \quad \leftarrow decision \ function$$
  $$\zeta_i = max\left(0, 1-y^{(i)}(\vec{w}\cdot\vec{x}^{(i)}+b)\right)=max\left(0, 1-y^{(i)}\cdot\hat{y}^{(i)}\right) \quad typical \ error$$
  $$\quad (1-\hat{y})$$
  $$y^{(i)}=+1 \ \rightarrow \ \zeta_i = max\left(0, y^{(i)}-\hat{y}^{(i)}\right)=max\left(0, actual-predict\right)$$
  $$y^{(i)}=-1 \ \rightarrow \ \zeta_i = max\left(0, \hat{y}^{(i)}-y^{(i)}\right)=max\left(0, predict-actual\right)$$

  Correctly classified data has all errors of zero.

  - **Hinge & squared hinge loss for (+) sample**



  - **Quadratic loss**

  $$y^{(i)}=\pm1 \ \rightarrow \ \zeta_i^2=\left(\hat{y}^{(i)}-y^{(i)}\right)^2$$

  The squared (quadratic) loss shows the following characteristics. The farther left or right you go from the boundary, the larger the error becomes.
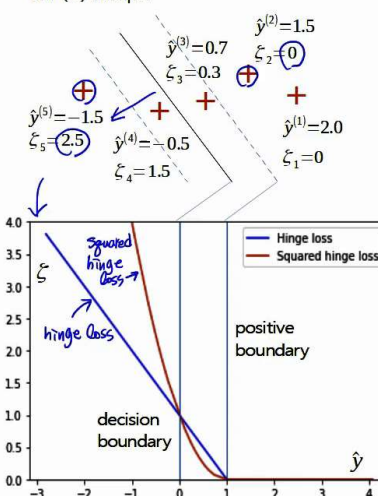


- **Objective function**

  $$min\left[\frac{1}{2}\|\vec{w}\|^2 + C\sum_{i=1}^{N}\zeta_i\right] \leftarrow \text{L2-regularized (or penalty) \& hinge loss}$$

  $$min\left[\frac{1}{2}\|\vec{w}\|^2 + C\sum_{i=1}^{N}\zeta_i^2\right] \leftarrow \begin{array}{l}\text{L2-regularized \& squared hinge loss}\\ \text{It is more sensitive to errors far from the}\\ \text{boundary.}\end{array}$$

■ Optimization: Lagrange primal function

▪ objective
$$min[\frac{1}{2}\|\vec{w}\|^2 + C\sum_{i=1}^{N}\zeta_i]$$

▪ constraints
$$1-\zeta_i-y^{(i)}(\vec{w}\cdot\vec{x}^{(i)}+b)\leq 0 \qquad y^{(i)}=\begin{cases} 1 & for\ "+" \\ -1 & for\ "-" \end{cases} \qquad -\zeta\leq 0$$

▪ inequality constrained optimization problem
$$\min_{x} f(x), \quad s.t \quad h(x)\leq 0$$

▪ Lagrangian primal function
$$L_p(x,\lambda)=f(x)+\lambda h(x) \quad \rightarrow \quad \lambda\geq 0$$

$$\min_{x} L_p(x,\lambda), \quad s.t \quad h(x)\leq 0$$

$$\frac{\partial L_p}{\partial x} = 0$$

$$L_p=\frac{1}{2}\|\vec{w}\|^2+C\sum_{i=1}^{N}\zeta_i+\sum_{i=1}^{N}\lambda_i\{1-\zeta-y^{(i)}(\vec{w}\cdot\vec{x}^{(i)}+b)\}+\sum_{i=1}^{N}\{-\mu_i\zeta_i\} \qquad (\zeta\geq 0,\ \lambda_i\geq 0,\ \mu_i\geq 0)$$

$$L_p=\frac{1}{2}\|\vec{w}\|^2+C\sum_{i=1}^{N}\zeta_i-\sum_{i=1}^{N}\lambda_i\{y^{(i)}(\vec{w}\cdot\vec{x}^{(i)}+b)-1+\zeta\}-\sum_{i=1}^{N}\mu_i\zeta_i$$

$$\frac{\partial L_p}{\partial \vec{w}} = 0 \rightarrow \vec{w}=\sum_{i=1}^{N}\lambda_i y^{(i)}\vec{x}^{(i)} \qquad\qquad \frac{\partial L_p}{\partial b} = 0 \rightarrow \sum_{i=1}^{N}\lambda_i y^{(i)} = 0$$

$$\frac{\partial L_p}{\partial \zeta_i} = C-\lambda_i-\mu_i=0 \rightarrow C=\lambda_i+\mu_i \qquad \lambda_i=C-\mu_i \rightarrow 0\leq\lambda_i\leq C$$

---

■ Optimization: Lagrange dual function

▪ The Lagrange dual function for the soft margin is obtained in the same way as for the hard margin.
▪ The dual function is the same as that of the hard margin, but has different constraints.

Primal function
$$\begin{cases} L_p=\frac{1}{2}\|\vec{w}\|^2+C\sum_{i=1}^{N}\zeta_i-\sum_{i=1}^{N}\lambda_i\{y^{(i)}(\vec{w}\cdot\vec{x}^{(i)}+b)-1+\zeta\}-\sum_{i=1}^{N}\mu_i\zeta_i \qquad (\zeta\geq 0,\ \lambda_i\geq 0,\ \mu_i\geq 0) \\ \vec{w}=\sum_{i=1}^{N}\lambda_i y^{(i)}\vec{x}^{(i)} \qquad \sum_{i=1}^{N}\lambda_i y^{(i)} = 0 \qquad C-\lambda_i-\mu_i=0 \rightarrow C=\lambda_i+\mu_i \rightarrow 0\leq\lambda_i\leq C \end{cases}$$

$$L_D=\frac{1}{2}(\sum_{i=1}^{N}\lambda_i y^{(i)}\vec{x}^{(i)})(\sum_{j=1}^{N}\lambda_j y^{(j)}\vec{x}^{(j)}) + C\sum_{i=1}^{N}\zeta_i-\sum_{i=1}^{N}[\lambda_i y^{(i)}(\sum_{i=1}^{N}\lambda_i y^{(i)}\vec{x}^{(i)})\cdot\vec{x}^{(i)}+\lambda_i y^{(i)}b-\lambda_i+\lambda_i\zeta_i]-\sum_{i=1}^{N}\mu_i\zeta_i$$

$$L_D=\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\lambda_i\lambda_j y^{(i)}y^{(j)}\vec{x}^{(i)}\cdot\vec{x}^{(j)}+\sum_{i=1}^{N}(C-\lambda_i-\mu_i)\zeta_i+\sum_{i=1}^{N}\lambda_i-\sum_{i=1}^{N}\sum_{j=1}^{N}\lambda_i\lambda_j y^{(i)}y^{(j)}\vec{x}^{(i)}\cdot\vec{x}^{(j)}+b\sum_{i=1}^{N}\lambda_i y^{(i)}$$

$$L_D=\sum_{i=1}^{N}\lambda_i-\frac{1}{2}\sum_{i=1}^{N}\sum_{j=1}^{N}\lambda_i\lambda_j y^{(i)}y^{(j)}\vec{x}^{(i)}\cdot\vec{x}^{(j)}$$

constraints: $0\leq\lambda_i\leq C,\ \sum_{i=1}^{N}\lambda_i y^{(i)} = 0$

▪ These are the Lagrange dual function and the constraints for the linear soft margin SVM.
▪ As in the hard margin SVM, the λ that maximizes this equation can be obtained through quadratic programming.

---

■ Decision function

▪ We can find the lambda from the dual function, and then use the lambda to find w. And we can use the support vectors to find b.
▪ We use the w and b to obtain the decision function. And we use the decision function to predict the class of the test data.



▪ Use λ to find w
$$\vec{w}=\sum_{i=1}^{N}\lambda_i y^{(i)}\vec{x}^{(i)}$$

▪ Use the support vectors (SV) to find b. SVs lie between the positive and negative boundaries.

$$\begin{cases} \vec{w}\cdot\vec{x}_i^{(+)}+b\geq 1-\zeta_i & \longleftarrow \text{"+" sample i satisfies this inequality.} \\ \vec{w}\cdot\vec{x}_j^{(-)}+b\leq -1+\zeta_j & \longleftarrow \text{"-" sample i satisfies this inequality.} \end{cases}$$

The SV(+) sample with the largest wx is on the positive boundary, and the SV(-) sample with the smallest wx is on the negative boundary (ξi = ξj = 0). Calculate b using these two samples.

$$b=-\frac{\left(max(\vec{w}\cdot\vec{x}_{SV}^{(+)}) + min(\vec{w}\cdot\vec{x}_{SV}^{(-)})\right)}{2}$$

▪ Decision function: We use this function to predict the class of the test data.

$$\hat{y}= w_1 x_1 + w_2 x_2 + ... + b$$

- ■ Quadratic Programming for linear soft margin SVM

  - It is the same as for hard margin SVM. However, only the constraints on λ are different.
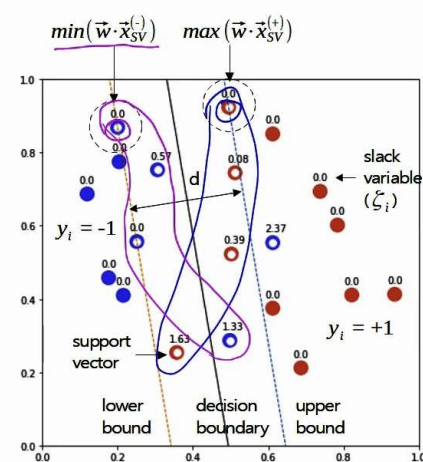
$$L_D = \sum_{i=1}^{N} \lambda_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \lambda_i \lambda_j y^{(i)} y^{(j)} \vec{x}^{(i)} \cdot \vec{x}^{(j)}$$

$$H_{i,j} = y^{(i)} y^{(j)} \vec{x}^{(i)} \cdot \vec{x}^{(j)} \quad \longleftarrow \quad \text{definition}$$

$$H = \begin{bmatrix} y^{(1)} y^{(1)} & y^{(1)} y^{(2)} \\ y^{(2)} y^{(1)} & y^{(2)} y^{(2)} \end{bmatrix} \times \begin{bmatrix} \vec{x}^{(1)} \cdot \vec{x}^{(1)} & \vec{x}^{(1)} \cdot \vec{x}^{(2)} \\ \vec{x}^{(2)} \cdot \vec{x}^{(1)} & \vec{x}^{(2)} \cdot \vec{x}^{(2)} \end{bmatrix} \quad \longleftarrow \quad \text{For } \underline{N=2} \text{ element wise product}$$

H = np.outer(y, y) * np.dot(x, x.T)   ← Python code

$$L_D = \sum_{i=1}^{N} \lambda_i - \frac{1}{2} \sum_{i=1}^{N} \sum_{j=1}^{N} \lambda_i \lambda_j H_{i,j}$$

$$L_D = \sum_{i=1}^{N} \lambda_i - \frac{1}{2} [\lambda_1 \ \lambda_2] \cdot H \cdot \begin{bmatrix} \lambda_1 \\ \lambda_2 \end{bmatrix}$$

$$L_D = \sum_{i=1}^{N} \lambda_i - \frac{1}{2} \lambda^T \cdot H \cdot \lambda$$

$$\text{s.t } -\lambda_i \le 0, \quad \lambda_i \le C$$

$$\sum_{i=1}^{N} \lambda_i y^{(i)} = 0$$

This constraint is added to the hard margin SVM.

- Standard form of QP

$$\underset{x}{argmin} \ \frac{1}{2} x^T P x + q^T x$$

$$\text{s.t } Gx \le h$$
$$Ax = b$$

- Lagrange dual function for linear soft margin SVM

$$\underset{\lambda}{argmin} \ \frac{1}{2} \lambda^T H \lambda - 1^T \lambda_i \quad \longleftarrow \quad \underset{\lambda}{argmax} \ L_D \rightarrow \underset{\lambda}{argmin} \ (-L_D)$$

$$\text{s.t } -\lambda_i \le 0$$
$$\lambda_i \le C$$
$$\sum_{i=1}^{N} \lambda_i y^{(i)} = 0$$

← Expressed in the same format as the standard form of QP.

$$P := H \quad \text{size} = N \times N$$
$$q := -\vec{1} = [[-1], [-1], ...] - \text{size} = N \times 1$$
$$A := y \quad \text{size} = N \times 1$$
$$b := 0 \quad \text{scalar}$$

$$G \cdot \lambda \le h \longrightarrow \begin{bmatrix} -1 & 0 & 0 & 0 \\ 0 & -1 & 0 & 0 \\ 0 & 0 & -1 & 0 \\ 0 & 0 & 0 & -1 \\ 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix} \cdot \begin{bmatrix} \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} = \begin{bmatrix} -\lambda_1 \\ -\lambda_2 \\ -\lambda_3 \\ -\lambda_4 \\ \lambda_1 \\ \lambda_2 \\ \lambda_3 \\ \lambda_4 \end{bmatrix} \le \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ C \\ C \\ C \\ C \end{bmatrix}$$
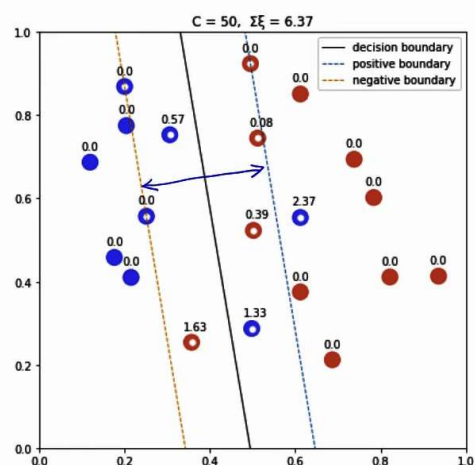
$$-\lambda_i \le 0$$
$$\lambda_i \le C$$

Since the constraints have changed, G and h are different from those of the hard margin SVM.

G     h

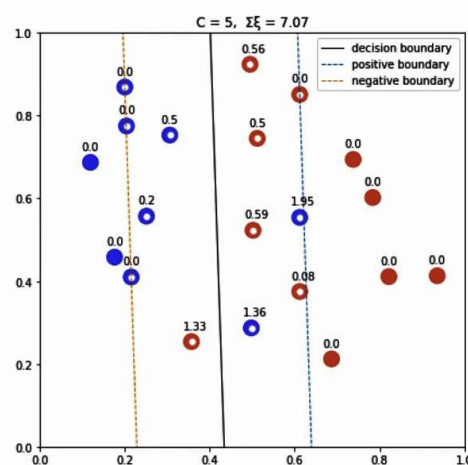- ■ Observation of changes in decision boundary according to changes in C

  - As C gets smaller, the margin and Σξ get bigger. This is because it focuses more on the goal of maximizing the margin. Conversely, as C increases, it focuses more on the goal of minimizing Σξ, so Σξ decreases but the margin also becomes smaller.

$$min\left[ \frac{1}{2} \|\vec{w}\|^2 + C \sum_{i=1}^{N} \zeta_i \right]$$

  - C = 50,  Σξ = 6.37
  - C = 5,  Σξ = 7.07
  - C = 1,  Σξ = 10.58