

# Soft Margin SVM

Thursday 2 May 2024 12:39 PM

→ The problem with hard margin SVM:  $\arg \max_{A, B, C} \left( \frac{2}{\sqrt{A^2 + B^2}} \right)$  s.t.  $\left\{ \begin{array}{l} y_i (A x_{i1} + B x_{i2} + C) \geq 1 \\ \text{for all } x_i \end{array} \right\}$

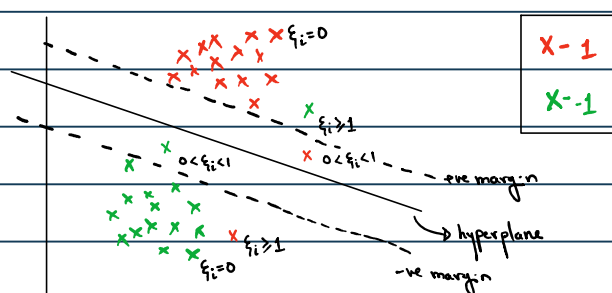
This constraint makes hard margin SVM inapplicable in most scenarios.

→ We need to 'soften' the constraint

→ The concept of slack variable is used in the formulation of the soft margin SVM to handle cases where data is not linearly separable, or when one allows for some degree of error in classification.

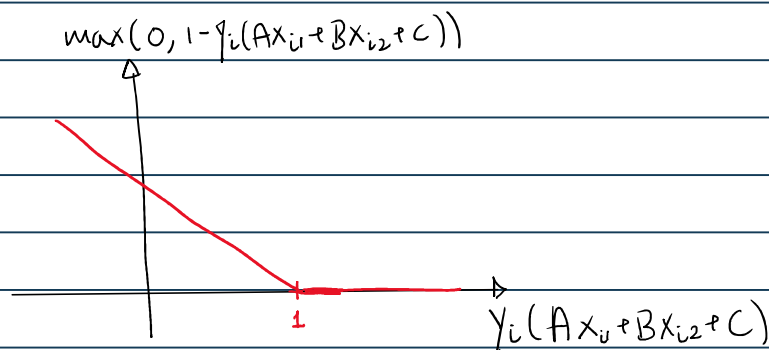
→ For each data point, a slack variable  $\xi_i$  is introduced.

→  $\xi_i$  measures the degree of misclassification of the data point  $x_i$



→  $\xi_i$  'misclassification score'  
↳ Hinge loss

$$\text{Hinge loss} = \max(0, 1 - y_i (A x_{i1} + B x_{i2} + C))$$



modified problem (overcorrection)

$$\arg \min_{A, B, C} \left( \frac{\sqrt{A^2 + B^2}}{2} \right) \text{ s.t. } \left\{ \begin{array}{l} y_i (A x_{i1} + B x_{i2} + C) \geq 1 - \xi_i \\ \text{for all } x_i \\ \text{s.t. } \xi_i \geq 0 \end{array} \right\}$$

→ it is allowing all points / All true condition / "no more a constraint"

→ In hard margin SVM the support vectors helped in deciding the margins. But now we have modified the constraint in a way that any point can lie anywhere, there is nothing stopping the margins from increasing infinitely. There is no criteria to stoppage for the margins in this modified problem.

→ addition of this term acts as constraint

$$\arg \min_{A, B, C} \left( \underbrace{\frac{\sqrt{A^2 + B^2}}{2}}_{\text{maximize margin}} + \underbrace{\frac{1}{n} \sum_{i=1}^n \xi_i}_{\text{reduce misclassification}} \right) \text{ s.t. } \left\{ \begin{array}{l} y_i (Ax_{i1} + Bx_{i2} + C) \geq 1 - \xi_i \\ \text{for all } x_i \\ \text{s.t. } \xi_i \geq 0 \end{array} \right.$$

addition of this term acts as constraint  
 opposite forces → goal is to find the balance

$$\arg \min_{A, B, C} \left( \frac{\sqrt{A^2 + B^2}}{2} + \underbrace{C \frac{1}{n} \sum_{i=1}^n \xi_i}_{\text{reduce misclassification}} \right) \text{ s.t. } \left\{ \begin{array}{l} y_i (Ax_{i1} + Bx_{i2} + C) \geq 1 - \xi_i \\ \text{for all } x_i \\ \text{s.t. } \xi_i \geq 0 \end{array} \right.$$

$C > 0$ , hyperparameter, allows us to decide which "force" to prioritize more  
 $C \uparrow$  more focus to reduce misclassification i.e.  $\downarrow$  margin vice versa.  
 opposite forces → goal is to find the balance

### Bias Variance tradeoff

$\uparrow C \rightarrow$  overfitting (low bias high variance)

$\downarrow C \rightarrow$  underfitting (high bias low variance)

### Relationship with Logistic Regression

$$\arg \min_{\beta_0, \beta_1, \beta_2} \left( \underbrace{\frac{\sqrt{\beta_1^2 + \beta_2^2}}{2}}_{\text{log loss}} + \underbrace{C \frac{1}{n} \sum_{i=1}^n \xi_i}_{\text{hinge loss}} \right)$$

$$\arg \min_{\beta_0, \beta_1, \beta_2} \left( \text{log loss} + \lambda \left( \sqrt{\beta_1^2 + \beta_2^2} \right) \right)$$

$\lambda \propto \frac{1}{C}$  i.e.  $C \propto \frac{1}{\lambda}$   
 $\| \beta \|^2$

→ This also creates linear decision boundaries only.