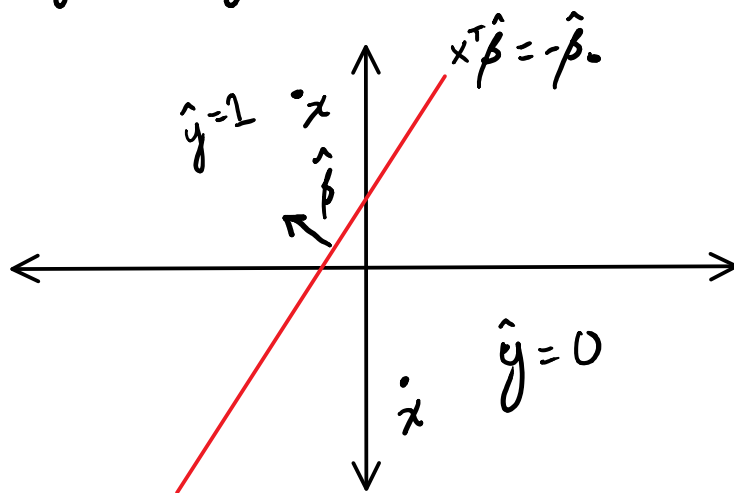


# Margin Based Methods

Tuesday, April 25, 2017

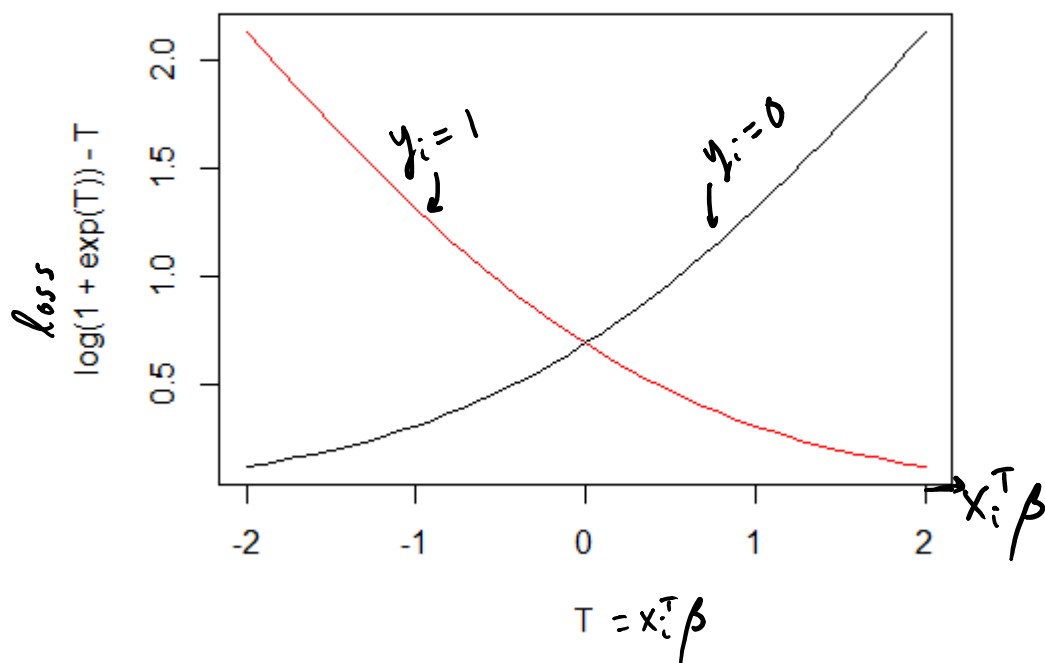
4:22 PM

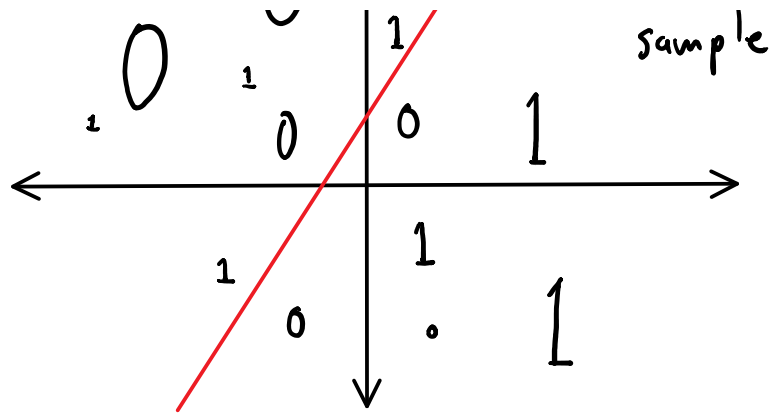
Predict for logistic regression and LDA :  $\hat{y} = 1 \{x^T \hat{\beta} + \hat{\beta}_0 \geq 0\}$



Empirical Risk Minimization

$$\min_{\beta} \frac{1}{n} \sum_{i=1}^n \ell(y_i, x_i, \beta) \rightarrow -y_i x_i^T \beta + \log(1 + e^{x_i^T \beta})$$





Recall 0-1 loss  $l(y_i, x_i, \beta) = \mathbb{1}\{y_i \neq \hat{y}_i\}$  and re-encode  $y_i \leftarrow 2y_i - 1$  so that  $y_i \in \{-1, 1\}$ , and  $\hat{y}_i = \text{sign}(x_i^T \beta)$

Error if  $y_i \cdot \hat{y}_i \neq 1 \Leftrightarrow y_i \cdot x_i^T \beta < 0$  so

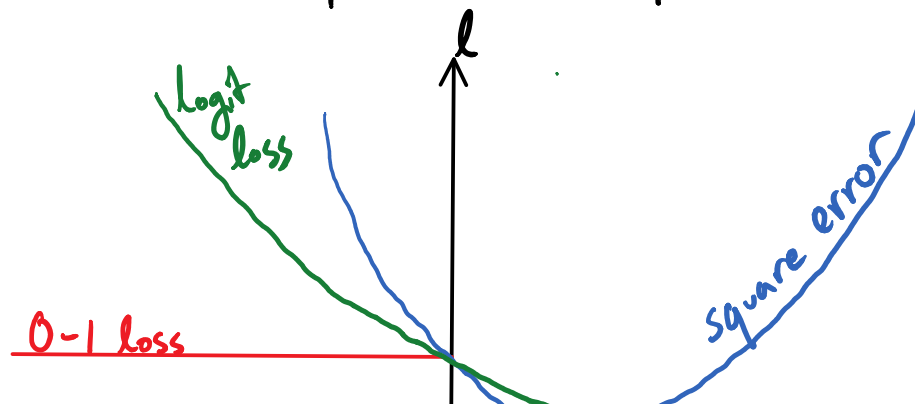
$$l_{01}(y_i, x_i, \beta) = \mathbb{1}\{y_i \cdot x_i^T \beta \leq 0\}$$

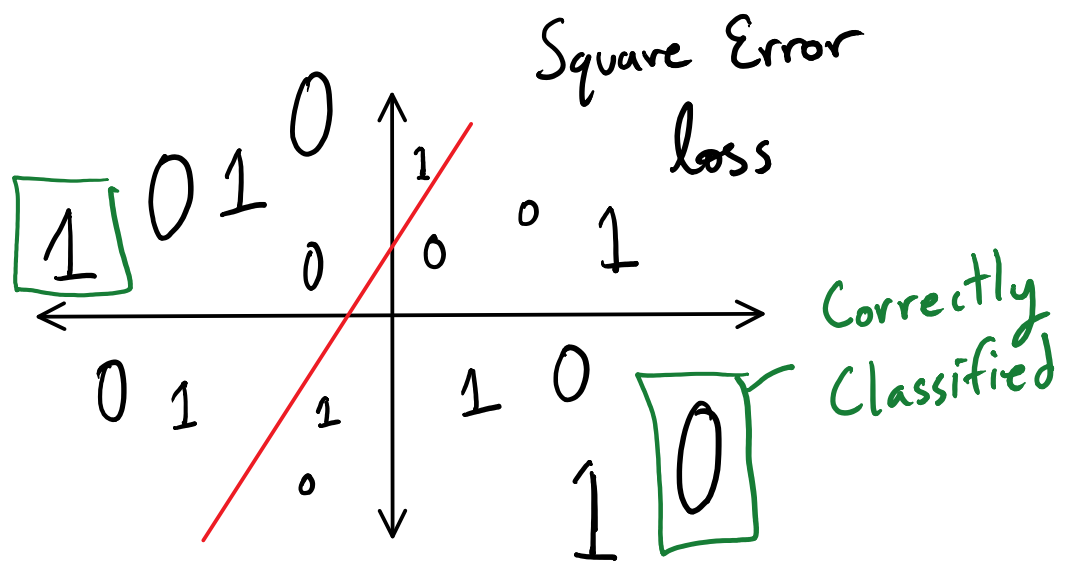
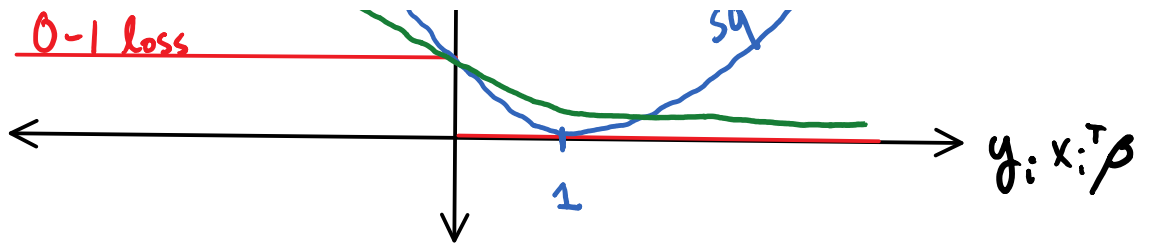
$$\text{and } l_{\text{logit}}(y_i, x_i, \beta) = \begin{cases} \log(1 + e^{x_i^T \beta}) & , y_i = -1 \\ \log(1 + e^{-x_i^T \beta}) & , y_i = 1 \end{cases}$$

$$= \log(1 + e^{-y_i x_i^T \beta})$$

You can also use square error loss,

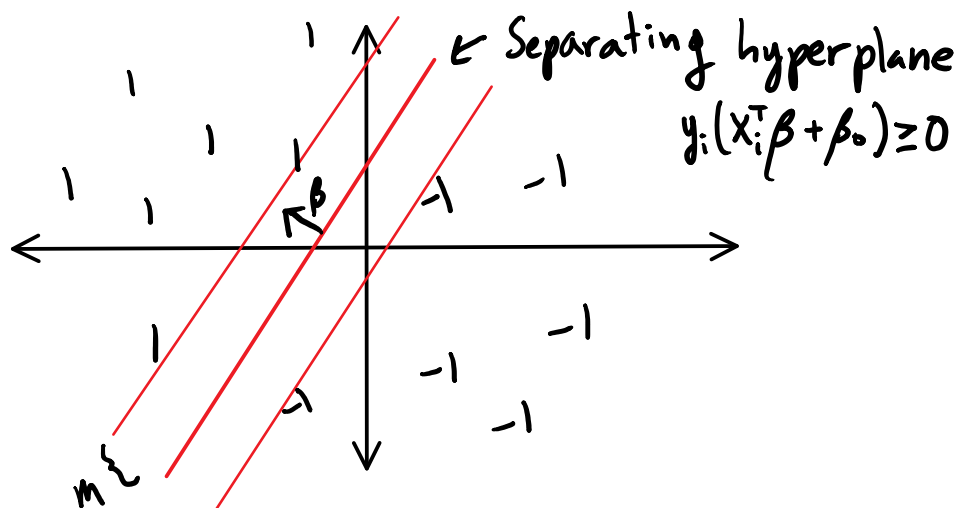
$$l_2(y_i, x_i, \beta) = (y_i - x_i^T \beta)^2 = (1 - y_i x_i^T \beta)^2$$





# Support Vector Machines

Wednesday, April 26, 2017 12:04 PM



Max-margin separating hyperplane

$$\max_{\beta, \beta_0: \|\beta\|=1} m \quad \text{s.t.} \quad y_i(x_i^T \beta + \beta_0) \geq m \quad \forall i$$

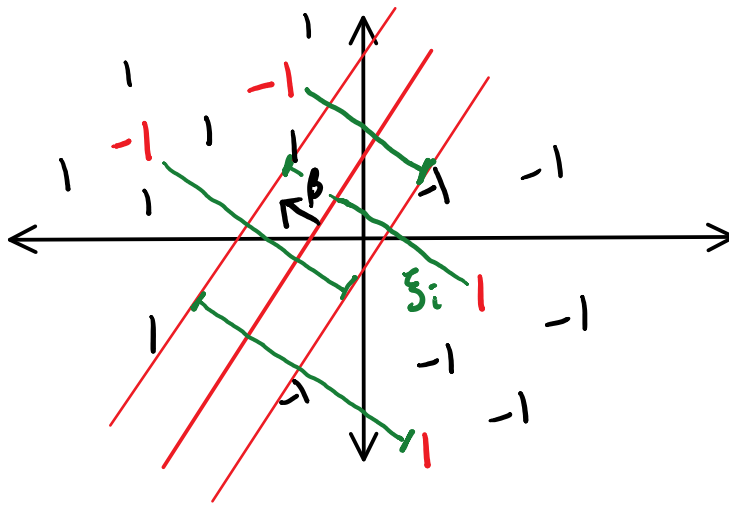
$$\equiv \max_{\beta, \beta_0} m \quad \text{s.t.} \quad \frac{1}{\|\beta\|} y_i(x_i^T \beta + \beta_0) \geq m \quad \forall i$$

$\hookrightarrow$  can scale  $\beta$  arbitrarily so set  
 $\|\beta\| = \frac{1}{m}$

$$\equiv \min_{\beta, \beta_0} \frac{1}{2} \|\beta\|_2^2 \quad \text{s.t.} \quad y_i(x_i^T \beta + \beta_0) \geq 1 \quad \forall i$$

$\hookrightarrow$  solution only dependent on "support vectors"

!!! / Feasible only if linearly separable !!!



if not linearly separable then add "slack variable"  $\xi_i$

$$\min \|\beta\| \text{ s.t. } y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad \forall i$$

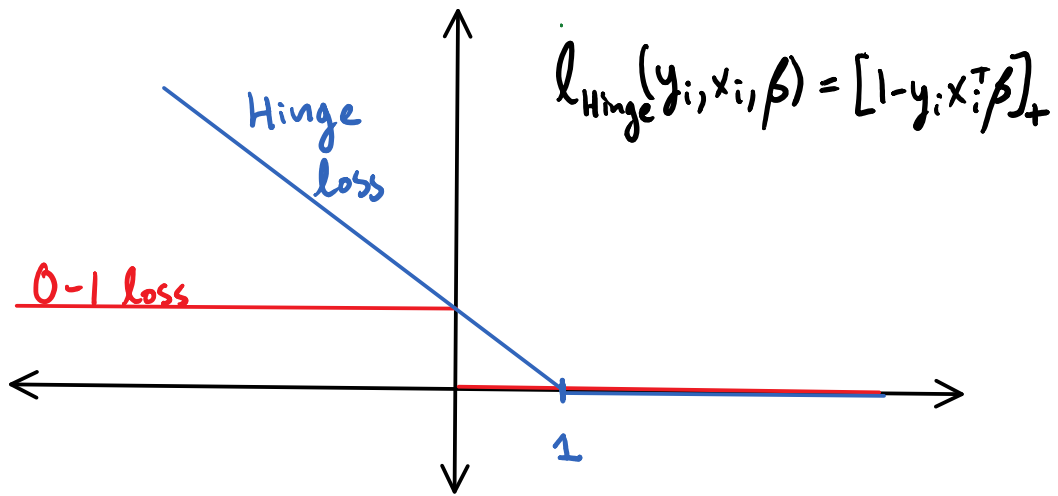
$$\xi_i \geq 0 \quad \sum_i \xi_i \leq C$$

Lagrangian

$$\min_{\beta, \beta_0, \xi_i \geq 0} \|\beta\|_2^2 + \lambda \sum_i \xi_i \text{ s.t. } y_i(x_i^T \beta + \beta_0) \geq 1 - \xi_i \quad \forall i$$

$$\equiv \min_{\beta, \beta_0} \sum_i [1 - y_i(x_i^T \beta + \beta_0)]_+ + \lambda \|\beta\|_2^2$$

$$\text{where } a_+ = \begin{cases} a, & a \geq 0 \\ 0, & a < 0 \end{cases}$$



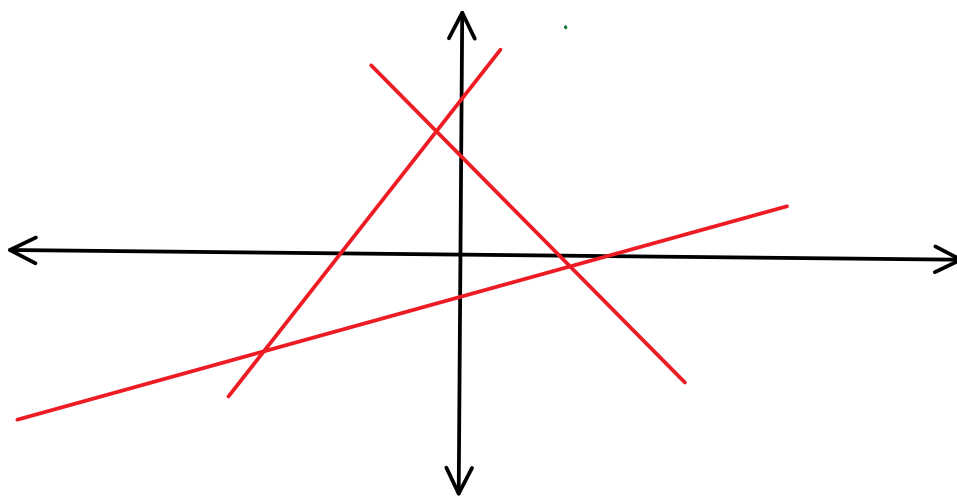
# Multiclass Classification

Wednesday, April 26, 2017 5:15 PM

Encode  $K$  classes:  $y_i \in \{0, 1\}^K$

$$\ell_{0,1}(y_i, \hat{y}_i) = 1 - y_i^T \hat{y}_i$$

Multiple linear separators:  $\{\beta_k\}_{k=1}^K$



$$\hat{y}_j = \begin{cases} 1, & j = \arg\max_k \beta_k^T x \\ 0, & \text{otherwise} \end{cases}$$

Soft-argmax

$z \in \mathbb{R}^K$  is vector of scores

$$s(z) = \left( \frac{e^{z_1}}{\sum_k e^{z_k}}, \frac{e^{z_2}}{\sum_k e^{z_k}}, \dots, \frac{e^{z_K}}{\sum_k e^{z_k}} \right)$$

Replace  $\hat{y}_i = s(x_i^T \hat{\beta}_1, \dots, x_i^T \hat{\beta}_K)$  then

$$u_i^T \hat{u}_i = e^{x_i^T \hat{\beta}_i} \quad \text{for } i = 1, \dots, K$$

$$\begin{aligned}
 y_i^T \hat{y}_i &= \frac{e^{x_i^T \hat{\beta}_j}}{\sum_k e^{x_i^T \hat{\beta}_k}} \quad \text{for } y_{ij} = 1. \\
 &\quad \underbrace{\hspace{1.5cm}}_{\text{rename } \hat{\beta}_j} \\
 &= \frac{e^{x_i^T (\hat{\beta}_j - \hat{\beta}_K)}}{1 + \sum_{k=1}^{K-1} e^{x_i^T (\hat{\beta}_k - \hat{\beta}_K)}} = P\{\text{Class} = j | x\} \\
 &\quad \text{for logistic model}
 \end{aligned}$$

$$\log \frac{P\{\text{Class } j | x\}}{P\{\text{Class } K | x\}} = x^T \hat{\beta}_j.$$

- ▷ multiclass SVM predicts class with max margin / smallest slack var.
- ▷ Confusion matrix is  $K \times K$ .