

Recall

1. Logistic Regression, (glm).
2. Introduced how to evaluate a classifier.

Training-Testing Split:

$$\text{data } D = \{(x_i, y_i) : i=1, \dots, n\}.$$

Let \mathcal{A} be an algorithm s.t. $\mathcal{A} : \mathcal{P}(D) \rightarrow \{g : \mathbb{R}^p \rightarrow \{0, 1\}\}$.

$$\mathcal{P}(D) = \{V : V \subseteq D\}. \quad \text{set of all subsets}$$

$$\{g : \mathbb{R}^p \rightarrow \{0, 1\}\}. \quad \text{set of all classifiers.}$$

Given data D , we can train a classifier $g = \mathcal{A}(D)$.

$$D = D_{\text{train}} \cup D_{\text{test}}.$$



$g = \mathcal{A}(D_{\text{train}})$ is our classifier.

$$\text{def. } \text{err}_{\text{train}} = \frac{1}{|D_{\text{train}}|} \sum_{i: x_i \in D_{\text{train}}} \mathbb{I}(y_i \neq g(x_i)).$$

$$\text{def } \text{err}_{\text{test}} = \frac{1}{|D_{\text{test}}|} \sum_{i: x_i \in D_{\text{test}}} \mathbb{I}(y_i \neq g(x_i)).$$

We argued that we cannot evaluate how "good" a classifier is solely based on training error.

If we take $g(x_i) = y_i$, then, $\text{err}_{\text{train}} = 0$.

It may do a good job on training set, but it cannot generalize to testing set.

If $\text{err}_{\text{test}} > \text{err}_{\text{train}} \Rightarrow$ overfitting.

$\text{err}_{\text{test}} \ll \text{err}_{\text{train}} \Rightarrow$ underfitting.

Drawback: We lose some efficiency.

K-fold cross-validation, (K-fold CV)

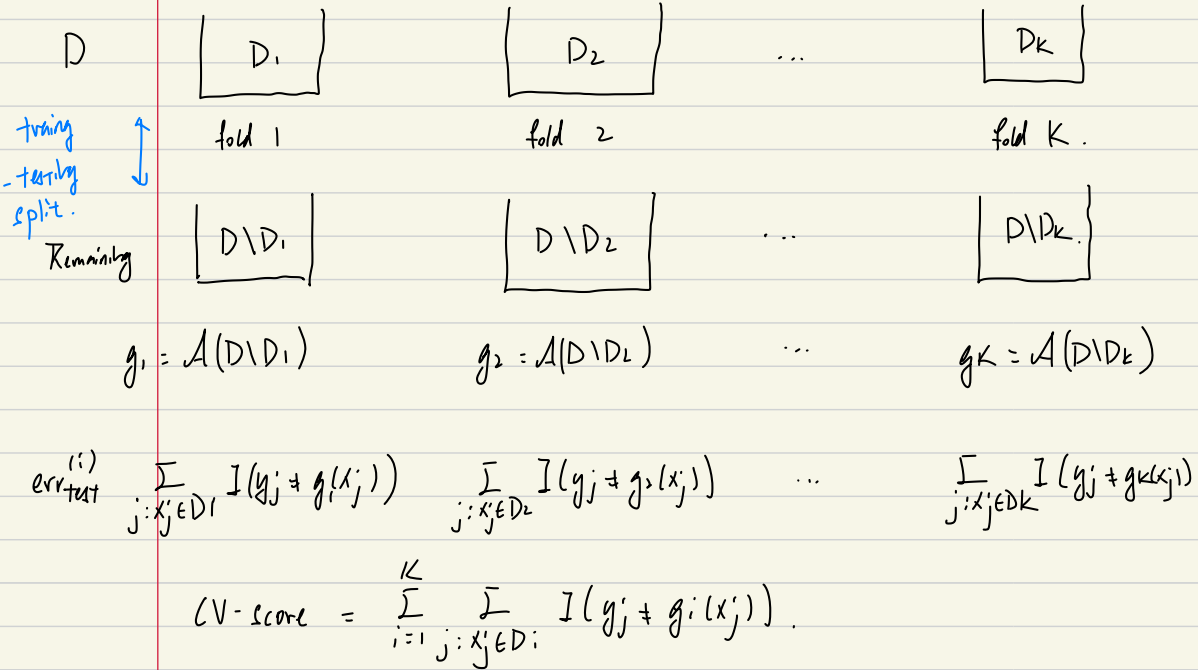
$D = \bigcup_{i=1}^K D_i$ each D_i is a "sub-test" set

$$g_i = A(D \setminus D_i)$$

$$\text{err}_{\text{test}}^{(i)} = \sum_{j: x_j \in D_i} I(y_j \neq g_i(x_j)).$$

$$\text{CV-score} = \sum_{j=1}^K \text{err}_{\text{test}}^{(i)}.$$

(Total error)



Special case of CV: Leave-One-Out CV (LOOCV).
 LOOCV is n -fold CV.

thm. For Linear Regression, LOOCV score is given by

$$CV\text{-score} = \sum_{i=1}^n \left(\frac{y_i - \hat{y}_i}{1 - h_{ii}} \right)^2, \quad h_{ii} = [X(X^T X)^{-1} X^T]_{ii}$$

$$\hat{y} = X\hat{\beta} = \frac{X(X^T X)^{-1} X^T y}{H} = Hy, \quad H = \begin{bmatrix} h_{11} & & & \\ & h_{22} & & \\ & & \ddots & \\ & & & h_{nn} \end{bmatrix}$$

proof.

i -th fold, Test set: $\{x_i, y_i\} = D_i$ Training set $\{(x_j, y_j)_{j \neq i}\} = D \setminus D_i$

$$\begin{aligned} \tilde{y}_i &= g(x_i), \quad g = \mathcal{A}(D \setminus D_i) \\ &= x_i^T (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T y_{(-i)}. \end{aligned}$$

$$D \setminus D_i = (X_{(-i)}, y_{(-i)}), \quad X_{(-i)} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_{i-1} \\ \vdots & \vdots \\ 1 & x_{i+1} \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \quad y_{(-i)} = \begin{pmatrix} y_1 \\ \vdots \\ y_{i-1} \\ \vdots \\ y_{i+1} \\ \vdots \\ y_n \end{pmatrix}.$$

thm, it is clear $\hat{\beta}_{(-i)} = (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T y_{(-i)} = ((X^T X)^{-1} X^T y)$.

$$\text{hence, } \tilde{y}_i = x_i^T \hat{\beta}_{(-i)} = x_i^T (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T y_{(-i)}.$$

$$X_{(-i)}^T X_{(-i)} = \sum_{\substack{j=1 \\ j \neq i}}^n x_j x_j^T = \sum_{j=1}^n x_j x_j^T - x_i x_i^T = X^T X - x_i x_i^T.$$

$$X_{(-i)}^T y_{(-i)} = \sum_{\substack{j=1 \\ j \neq i}}^n x_j y_j = \sum_{j=1}^n x_j y_j - x_i y_i.$$

$$(X^T X - \lambda_i X_i^T)^{-1} = (X^T X)^{-1} + \frac{(X^T X)^{-1} X_i X_i^T (X^T X)^{-1}}{1 - X_i^T (X^T X)^{-1} X_i} \quad \text{Sherman-Morrison.}$$

$$X_{(-i)}^T y_{(-i)} = \sum_{j=1}^n X_j^T y_j = \sum_{j=1}^n X_j^T y_j - X_i^T y_i = X^T y - X_i^T y_i.$$

$$\tilde{y}_i = X_i^T \hat{\beta}_{(-i)} = X_i^T (X_{(-i)}^T X_{(-i)})^{-1} X_{(-i)}^T y_{(-i)}.$$

$$= X_i^T (X^T X)^{-1} X^T y - X_i^T (X^T X)^{-1} X_i^T y_i.$$

$$+ \frac{X_i^T (X^T X)^{-1} X_i X_i^T (X^T X)^{-1} X_i^T y_i}{1 - X_i^T (X^T X)^{-1} X_i} - \frac{X_i^T (X^T X)^{-1} X_i X_i^T (X^T X)^{-1} X_i^T y_i}{1 - X_i^T (X^T X)^{-1} X_i}.$$

$$\text{Note: } X (X^T X)^{-1} X^T = H.$$

$$X_i^T (X^T X)^{-1} X_i = h_{ii}.$$

$$X_i^T (X^T X)^{-1} X^T y = \hat{y}_i.$$

$$= \hat{y}_i - h_{ii} y_i + \frac{h_{ii} \hat{y}_i}{1 - h_{ii}} - \frac{h_{ii}^2 y_i}{1 - h_{ii}}.$$

$$= \hat{y}_i \left\{ 1 + \frac{h_{ii}}{1 - h_{ii}} \right\} - h_{ii} y_i \left\{ 1 + \frac{h_{ii}}{1 - h_{ii}} \right\}.$$

$$\tilde{y}_i = \hat{y}_i \frac{1}{1 - h_{ii}} - h_{ii} y_i \frac{1}{1 - h_{ii}} = \frac{1}{1 - h_{ii}} (\hat{y}_i - h_{ii} y_i).$$

\tilde{y}_i = predicted from $X_{(-i)}$, $y_{(-i)}$.

\hat{y}_i = predicted from X , y .

$$y_i - \tilde{y}_i = y_i - \frac{\hat{y}_i}{1 - h_{ii}} + \frac{h_{ii}}{1 - h_{ii}} y_i = \frac{1}{1 - h_{ii}} y_i - \frac{1}{1 - h_{ii}} \hat{y}_i = \frac{y_i - \hat{y}_i}{1 - h_{ii}}.$$

$$\Rightarrow \text{CV-score} = \sum_{j=1}^n (y_j - \tilde{y}_j)^2 = \sum_{j=1}^n \left(\frac{y_j - \hat{y}_j}{1 - h_{jj}} \right)^2.$$

thm. If $\hat{y} = S y$, where S is a matrix, then, (S linear smoother)

$$CV\text{-Score} = \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{1 - h_{ii}}$$

If $h_{ii} \approx 1$, then, this is not numerically stable.

$$\text{Generalize CV (GCV)} = \sum_{i=1}^n \frac{(y_i - \hat{y}_i)^2}{(1 - \text{tr}(S)/n)^2} \quad \text{tr}(S) = \sum_{i=1}^n h_{ii}$$

Unfortunately, logistic regression is not a linear smoother.

Smoothing Technique. (Spline).
Kernel Trick (KDE). } \Rightarrow CV closed form available.

Use of CV.

1. To evaluate the performance of classifier.
2. Hyperparameter selection.

HW2 Q2. Regularized Logistic Regression.

Logistic Regression: $l(\beta) = \sum_{i=1}^n \log_2 k(\beta)$.

Regularized Logistic Regression. $p_l(\beta, \lambda) = \sum_{i=1}^n \log_2 k(\beta) + \lambda \|\beta\|_2^2$
 \downarrow
 hyperparameter.

Ridge Regression. $\Rightarrow l(\beta) + \lambda \|\beta\|_2^2 \quad \|\beta\|_2^2 = \sum_{j=1}^p \beta_j^2$

LASSO Regression. $\Rightarrow l(\beta) + \lambda \|\beta\|_1 \quad \|\beta\|_1 = \sum_{j=1}^p |\beta_j|$

$$CV(\lambda) = \sum_{j=1}^K \sum_{i: x_i \in D_j} I(y_i \neq g_j(x_i, \lambda))$$

$$\hat{\lambda} = \underset{\lambda}{\text{argmin}} CV(\lambda)$$