

Residual Sum of Squares (RSS): $RSS(w) = (y - Xw)^T(y - Xw)$, $\widehat{w}_{LS} = (X^T X)^{-1} X^T y$ (Least squares optimal solution)

Handling intercept: (1) Show that $\frac{1}{n} \sum_i x_i = 0(*)$ (2) "Demmean" linear regression model so * holds

Interpreting the result: When other features are fixed, the estimated change in \hat{y} for one-unit change of x

Reason for minimizing RSS: if $y_i = wx_i + \epsilon_i$ and $\epsilon \sim N(0, \sigma^2)$ which means $y \sim (X^T w, \sigma^2)$, the MLE for w produces minimum RSS.

Feature map: in polynomial regression, transform each data point x_i to higher-dimensional feature vector $h(x_i)$

$h_j(x)$: j^{th} feature associated with input x (j^{th} basis function)

Supervised Learning: learning a function that maps inputs to outputs based on labelled examples

For linear regression, the loss function is RSS.

Handling an intercept in linear regression: model is $y = wx + b$, $\hat{b} = \frac{1}{n} \sum_{i=1}^n y_i$, w still the same.

In the case of feature map, replace x with h

Training Error: $e_{train} = \frac{1}{n} \sum_{i=1}^n \text{loss}(y_i, \hat{f}_w(x_i))$, dependent on how loss function is defined. Training error is overly optimistic on training data (unless training set includes everything you might ever see).

Generalization error: $\int_{x,y} (y - f_w(x))^2 P(x,y) dx dy$

As model complexity goes up, training error goes down; generalization error goes down, then goes up.

Testing Error: same as training error, but on a labelled testing data set.

Three sources of error:

Data inherited noise: $y = wx + \epsilon$, irreducible.

Bias: Let $\widehat{f}_w(x) = E_{train}[f_{w_{train}}(x)]$, then $\text{Bias}_{x_i} = f_{w_{true}}(x_i) - \widehat{f}_w(x_i) \rightarrow$ is our model flexible enough to catch true w ?

Variance: How much do specific fit varies from expected fit? $\text{Var}_x = E_{train} \left[(f_{w_{train}}(x_i) - \widehat{f}_w(x_i))^2 \right]$

Low complexity: high bias, low variance

High complexity: high variance, low bias

If $y = wX + \epsilon$ and $\epsilon \sim N(0, \sigma)$, average prediction error_i = $\sigma^2 + \text{Bias}_i^2 + \text{Variance}_i$

Mean Squared Error (MSE) = $E_{train} \left[(f_{w_{true}}(x_t) - \widehat{f}_{w_{train}}(x_t))^2 \right] = \text{Bias}^2 + \text{Var}$

Regularization: imposes "simpler" solutions by a complexity penalty

Rigid Regression: L_2 penalized least-squares regression. Trade of model complexity with training error

$\widehat{w}_{ridge} = \text{argmin}_w \left(\sum_{i=1}^n (y_i - x_i^T w)^2 \right) + \lambda \|w\|_2^2$, solves to $\widehat{w}_{ridge} = (X^T X + \lambda I)^{-1} X^T y$

Shrinkage Property: assume $X^T X = I$ (a special case), $\widehat{w}_{ridge} = \frac{n}{n+\lambda} \widehat{w}_{LS}$

Assume $X^T X = nI$ and $y = Xw + \epsilon$, $\epsilon \sim N(0, \sigma^2)$, we have $\widehat{w}_{ridge} = \frac{n}{n+\lambda} w + \frac{1}{n+\lambda} X^T \epsilon$

predicted average error = $\sigma^2 + \frac{\lambda^2}{(n+\lambda)^2} (w^T x)^2 + \frac{d\sigma^2 n}{(n+\lambda)^2} \|x\|_2^2$

$\lim_{\lambda \rightarrow 0} \widehat{w}_{ridge} = \widehat{w}_{LS}$, $\lim_{\lambda \rightarrow \infty} \widehat{w}_{ridge} = 0$.

Leave-One-Out (LOO) cross validation: $D \setminus j$ is the set of training data with data point j removed

$$\text{error}_{LOO} = \frac{1}{n} \sum_{j=1}^n \left(y_j - f_{D \setminus j}(x_j) \right)^2$$

error_{LOO} is usually pessimistic because learning with less data points gives worse answer.

LOO is almost unbiased, so it can be used for model selection (i.e. picking the right λ in regularization).

k-fold cross validation: Divide data to k equal parts $D_1, D_2, D_3, \dots, D_k$

Estimation error on set D_i : $\text{error}_{D_i} = \frac{1}{|D_i|} \sum_{x_i, y_i \in D_i} \left(y_i - f_{D \setminus D_i}(x_i) \right)^2$

k-fold error: average over all sets $\text{error}_{k\text{-fold}} = \frac{1}{k} \sum_{i=1}^k \text{error}_{D_i}$

Faster to compute and much more pessimistic than LOO (because model is trained with even less data points)

Usually choose $k = 10$

Given a set of input data: (1) split to train and test set (2) use k-fold cross validation on training set to train predictor and choose magic parameters such as λ (3) Assessment: use test set to check the accuracy of the model

Greedy Model/Feature Selection: (1) Start from an empty set of features $F_0 = \{\}$ (2) Run learning algorithm on current feature set to get weight for features in F_0 (3) Choose the best next feature h_1 which minimizes the training error when using $\{F_0, h_1\}$ (4) Recurse from step (1)

Theorem for Penalized Least Squares: For any $\lambda \geq 0$ for which \hat{w}_r achieves the minimum, there exists a $v \geq 0$ such that

$$\hat{w}_r = \operatorname{argmin}_w \sum_{i=1}^n (y_i - x_i^T w)^2 \text{ subject to } r(w) \leq v$$

LASSO Regression: $\hat{w}_{ridge} = \operatorname{argmin}_w \left(\sum_{i=1}^n (y_i - x_i^T w)^2 \right) + \lambda \|w\|_1$, keep $\sum_{i=1}^n y_i = 0$ and $\sum_{i=1}^n x_i = 0$ for not dealing with offset terms.

Optimizing LASSO objective one coordinate at a time: For $j = 1, 2, \dots, n$, $r_i^j = \sum_{k \neq j} x_{ik} w_k$ and $\hat{w}_j = \operatorname{argmin}_{w_j} \sum_{i=1}^n \left(r_i^j - x_{i,j} w_j \right)^2 + \lambda |w_j|$

Subgradient: $f(y) \geq f(x) + g^T(y - x)$, w is a minimum if 0 is a subgradient at w

Coordinate Descent for LASSO (Shooting Algorithm): Pick a random coordinate j , compute $a_j = \sum_{i=1}^n x_{ij}^2$, $c_j =$

$$2 \sum_{i=1}^n (y_i - \sum_{k \neq j} x_{ik} w_k) x_{i,j}$$

L_1 Regularization is one way to do variable selection (finding a sparse solution). LASSO is non-differentiable, but convex

Binary Classification: $\text{loss}(f(x), y) = 1\{f(x) \neq y\}$

$$\text{Sigmoid for Binary Class: } P(Y = 0|w, X) = \frac{1}{1 + \exp(w_0 + \sum_k w_k x_k)}$$

$$\text{Linear Decision Rule: } \log \left(\frac{P(Y = 1|w, X)}{P(Y = 0|w, X)} \right) = w_0 + \sum_k w_k x_k$$

Conditional Likelihood: if $y_i \in \{-1, 1\}$, $P(Y = y|x, w) = \frac{1}{1 + \exp(-y w^T x)}$, $\hat{w}_{MLE} = \operatorname{argmin}_w \sum_{i=1}^n \log(1 + \exp(-y_i x_i^T w)) = J(w)$, $J(w)$ is a convex function, no closed-form solution

To prevent overfitting: add a regularization term (i.e. L_2)

Gradient Descent: $w_{t+1} = w_t - \alpha \nabla w_t$

With Taylor Series Approximation: $w_{t+1} = w_t - \alpha \nabla f(w_t)$

$$\text{For } f(w) = \frac{1}{2} \|Xw - y\|_2^2, \nabla f(w_t) = (X^T X)^{-1} X^T y, (w_{t+1} - w_*) = (I - \alpha X^T X)(w^t - w_*) = (I - \alpha X^T X)^{t+1} (w_0 - w_*)$$

$$\text{Total Derivative: } \sum_{i=1}^n \frac{\delta l_i(w)}{\delta w_j} + 2\lambda w_j$$

Stochastic Gradient Descent: Use one data point to estimate the actual averaged gradient over all coordinates

Each data point contributes to $\frac{1}{N}$ of regularization.

$$\text{Estimated Derivative: } \frac{\delta l_i(w)}{\delta w_j} + \frac{2\lambda w_j}{N} \text{ for some random data point } i$$

SGD has noisy convergence.