

Linear Regression

Cost Function

$$L(w) = \frac{1}{2N} \sum_{i=1}^N (y_i - x_i^T w)^2 = \frac{1}{2N} \|y - X^T w\|_2^2$$

- Need to optimize the cost function: $w = \arg \min_w L(w)$
- Solution for cost function: $\frac{\nabla L(w)}{\nabla w} = \frac{1}{N} X(X^T w - y) = 0 \Leftrightarrow XX^T w = Xy$
 - If XX^T is invertible, $w = (XX^T)^{-1} Xy$
 - o.w. we use pseudo-inverse, $w = (XX^T)^+ Xy$

Bias Trick

- To make $y = x^T w + b$ become $y = x^T w$ $\begin{cases} x = [x_1 & x_2 & \dots & x_n]^T \rightarrow x = [x_1 & x_2 & \dots & x_n & 1]^T \\ w = [w_0 & w_1 & \dots & w_n]^T \rightarrow w = [w_0 & w_1 & \dots & w_n & b]^T \end{cases}$

Drawbacks

- Linear regression is very **sensitive to noise**, so we need to
 - Pre-processing data before applying linear regression
 - Change the cost function to make it **robust to noise**: **Huber Regression**
- Linear Regression cannot represent the high order function

Ridge Regression

- Aside from the purpose of avoiding overfitting, we can apply ridge when XX^T is not invertible and we don't want to apply pseudo-inverse

$$L_2(w) = \frac{1}{2N} (\|y - X^T w\|_2^2 + \lambda \|w\|_2^2)$$

$$\frac{\nabla L_2(w)}{\nabla w} = 0 \Leftrightarrow \frac{1}{N} (X(X^T w - y) + \lambda w) = 0 \Leftrightarrow (XX^T + \lambda I)w = Xy$$

- Prove $A = XX^T + \lambda I$ is positive definite, so A is always invertible

$$w^T A w = w^T (XX^T + \lambda I) w = w^T XX^T w + \lambda w^T w = \|X^T w\|_2^2 + \lambda \|w\|_2^2 > 0$$

- So, $\frac{\nabla L_2(w)}{\nabla w} = 0$ always have 1 solution: $w = (XX^T + \lambda I)^{-1} Xy$

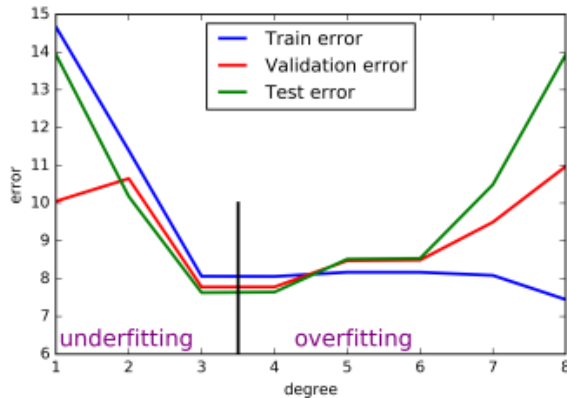
Overfitting

- If our model is too fit with the data, the model is more likely to lose its generalization
- We can use the metrics like MSE, Cross entropy, Hinge Loss (Multi class SVM) to measure the error of training and test set
- If the metric of error in training set < test set, we say it's overfitting

Validation

- Validation set is small dataset extracted from training set
- Its analogy: Suppose we are reviewing before taking exam. Exam questions in previous years is considered training set and exam questions in the exam we are gonna take is test set. When reviewing, we separate the training set into 2 parts: first part has questions and answer, second part just has questions

- So now, our model will have 3 errors, so our model needs to have small *training error* and *validation error*, so it's more likely to have small *test error*
- **Based on this method, we can tune the model or find different model on the training set and apply them on validation set to find small validation error**
- E.g. In regression, we increase the degree, the training error definitely decreases and validation error also decreases. But when we increase the degree to some extent, the validation error will stop decreasing and start to increase, so we will stop tune the degree at this point



Hình 8.2: Lựa chọn mô hình dựa trên validation error.

- Note: Test set and Validation set need to have something in common like distribution

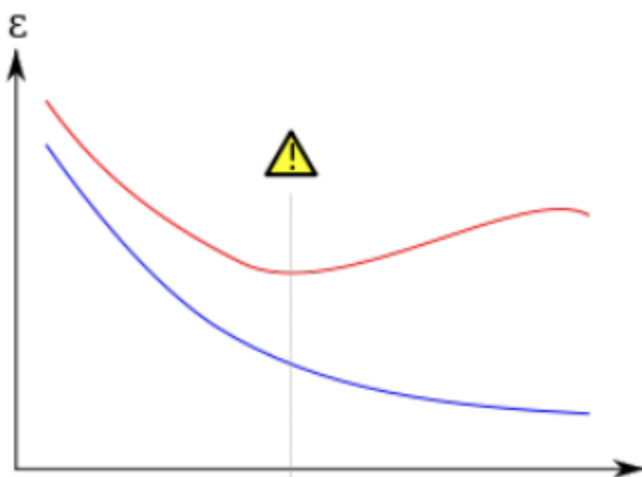
Cross Validation

- Suppose we just have not so big training set to build the model
- If we extract too much training data to make validation set, training set will not have enough data to build the model. So now validation set needs to be small. But when validation set is small, overfitting may occur in validation set
- Cross Validation will separate training set into k subsets. We run the model k times, each time, one of k subset will be considered as validation set, the model will be built based on $k-1$ subsets as training set
- Final model will be the mean of training error and mean of validation error
- **Leave-one-out:** $k = \text{number of data points}$, extract 1 data point to make validation set

Regularization

- Big disadvantage of cross validation is building the model k times
- In the regression above, we just need to tune the degree, but in the other models, we have to tune other parameters and the range of parameter is large, so cross validation is too complex in such models like this

Early Stopping



Hình 8.3: Early stopping. Đường màu xanh là training error, màu đỏ là validation error. Trục x thể hiện số lượng vòng lặp, trục y là giá trị error. Thuật toán huấn luyện dừng lại tại vòng lặp mà validation error đạt giá trị nhỏ nhất (Nguồn: [Overfitting – Wikipedia](#)).

Add to Cost function

$$L_{reg}(\theta) = L(\theta) + \lambda R(\theta)$$

- LASSO Regression: $R(\theta) = \|\mathbf{w}\|_1$, parameter \mathbf{w} is sparse vector. Data points which value in $\mathbf{w} \neq 0$ means this data point is important, contribute to predict output, otherwise this data point is not important. So LASSO regression is considered as FEATURE SELECTION and robust to noise

- But the LASSO regression has disadvantage: it's hard to find derivatives of $\|\mathbf{w}\|_1$