

Learning Reflection 2

* 2022/4/4 (Mon.)

* Cross Validation

choosing complexity \rightarrow need lower train error, lower test error.

2 ways to pick the model complexity without training test set:

1. Validation set

2. Cross Validation

| Train | val: | Test |



fast



sacrifice training data

Cross Validation :

use small validation sets w/o losing too much data.

Slow, usually use $K=5$ to 10 .

train, validation, test = split_data(dataset)

chunk1, ..., chunkK, test = split_data(dataset)

for each model complexity p :

model = train_model(model-p, train)

val_err = error(model, validation)

Keep track of p with smallest val_err

return best p + error(model, test).

for each model complexity p:

for i in [1, K]:

model = train_model(model-p, chunks-i)

val_err = error(model, chunk=i)

avg_val_err = average val_err over chunks.

Keep track of p with smallest avg_val_err

return model trained on train with best

p + error(model, test)

* Coefficients and overfitting

Overfitting - large estimated parameters \hat{w} , low bias, high variance.

Underfitting - small parameters \hat{w} , high bias, low variance.

* Overfitting is not limited to polynomial regression of large degree.

It can also happen if you use a large # of features.

* prevent overfitting. Slide 17.

Regularization.

Quality metric (minimized loss) : $\hat{w} = \min_w L(w)$

$\hat{w} = \min_w L(w) + \pi R(w)$ measure the magnitude of coefficients,
measure of fit

Say, $w = [w_0, w_1, \dots, w_D]$,

Sum: $R(w) = \sum_{j=0}^D w_j$

Doesn't work, $w = [10000, -10000]$
 $R(w) = 0$

Sum of abs.: $R(w) = \sum_{j=0}^D |w_j| \triangleq \|w\|_1$, L1 norm Lasso

Sum of sqrs: $R(w) = \sum_{j=0}^D w_j^2 \triangleq \|w\|_2^2$ L2 norm Ridge.

Ridge Regression (L2 Norm)

minimize $\hat{w} = \min_w RSS(w) + \pi \|w\|_2^2$

π is a tuning parameter that changes how much the model cares about the regularization term.

$\pi=0$, $\hat{w} = \min_w RSS(w) = \hat{w}_{OLS}$ (OLS) ordinary least squares.

$\pi=\infty$, $\hat{w} = \vec{0}$

If any $w_j \neq 0$, $RSS(w) + \pi \|w\|_2^2 = \infty$

If all $w_j = 0$, $RSS(w) + \pi \|w\|_2^2 = RSS(w) < \infty$

π in between

$$0 \leq \|w\|_2^2 \leq \|\hat{w}_{OLS}\|_2^2$$



Complex model.

$N=0 \rightarrow$ model have many parameters, low bias, high variance

$N=\infty \rightarrow$ model have few parameters, low variance, high bias.
Simple model.

* Code to train a Ridge Model.

```
from sklearn.linear_model import Ridge.
```

```
model = Ridge(alpha=1.5)
```

* Intercept.

Two ways dealing with Intercept Biases.

1. Change the measure of overfitting to not including the intercept.

$$\min_{w_0, w_{rest}} \text{RSS}(w_0, w_{rest}) + \gamma \|w_{rest}\|_2^2$$

2. Center the y values so they have mean 0.

* Scaling Features.

Scale of coefficients.

$$x = 100 \text{ m}^2 = \frac{100}{(1000)^2} \text{ km}^2$$

This makes the input value to be smaller, and the w_i will be larger.

This accidentally penalize features

for having large coefficients due to having small value inputs.

Fix this by normalizing the feature

$$x_j(x_i) = \frac{h_j(x_i) - \mu_j(x_1, \dots, x_N)}{\sigma_j(x_1, \dots, x_N)}$$

where the mean of feature j : $\mu_j(x_1, \dots, x_N) = \frac{1}{N} \sum_{i=1}^N h_j(x_i)$

standard deviation of feature j : $\sigma_j(x_1, \dots, x_N) = \sqrt{\frac{1}{N} \sum_{i=1}^N (h_j(x_i) - \mu_j(x_1, \dots, x_N))^2}$

* Must scale the test data & all feature data using the means & std. of training set.

2022, 4.6(3)

- * Overfitting :
 - | polynomial regression of large degree.
 - | large number of features.

- * How to choose the best value of γ ?

After we train each model with a certain γ_i in order to find

$$\hat{w}_i = \operatorname{argmin}_w \text{MSE}(w) + \gamma_i \|w\|_2^2$$

→ Pick the γ_i that has the smallest $\text{MSE}(\hat{w}_i)$ on the validation set.

We want to choose the γ that will do best on future data.

= Need to minimize validation error.

- * for γ in γ s:

Train a model using Gradient Descent

$$\hat{w}_{\text{ridge}}(\gamma) = \underset{w}{\operatorname{min}} \text{MSE}_{\text{train}}(w) + \gamma \|w_{2:D}\|_2^2$$

Compute Validation Error.

$$\text{Validation error} = \text{MSE}_{\text{val}}(\hat{w}_{\text{ridge}}(\gamma))$$

Train γ with smallest validation error.

Return γ^* + estimated future error $\text{RSS}_{\text{test}}(\hat{w}_{\text{ridge}}(\gamma))$

Feature selection

- ① Complexity, adjust to appropriate
- ② Interpretability
- ③ Efficiency

If w is sparse, only need to look at the non-zero coefficients.

$$\hat{y} = \sum_{w_j \neq 0} w_j h_j(x)$$

How many linear regression models in total do we have to evaluate for a dataset with d features in order to find the global optimum?

$$\boxed{2^d}$$

Greedy Algo.

| Forward stepwise
| Backward stepwise
| Combining.

Runtime: $O(K^2)$

$$\text{min_val} = \infty$$

for $i \leftarrow 1 \dots k$:

Find feature f_i not in S_{i-1} , that when we

- combined w/ S_{i-1} , minimizes the validation loss the most.

$$S_i \leftarrow S_{i-1} \cup \{f_i\}$$

if $\text{val-loss}(S_i) > \text{min_val}$:

break.

Sparcity

In regression, means many of the learned coefficients are 0.

Hyper-parameter

A parameter you specify for the model that influences which param. (e.g. coefficient) are learned by Mr. alg.

LASSO

Adds bias to the Least Squares solution (and Variance \rightarrow overfitting)

If want to remove bias from LASSO soln.:

1. Run LASSO to select features
2. Run regular OLS on dataset with only those features.

LASSO select highly correlated features arbitrarily amongst them.

\rightarrow ElasticNet: $\hat{w}_{\text{elastic}} = \min_w \text{RSS}(w) + \gamma_1 \|w\|_1 + \gamma_2 \|w\|_2^2$

Summary:

1. For feature selection,
adjusting feature size to pick the best many having lowest $MSE_{val}(\hat{w})$.
2. Three options choosing feature numbers:
a Validation set; Cross Validation;
other metrics for penalizing complexity like BIC. Bayesian Information Criterion.
3. L₁ penalty causes more sparse than L₂ penalty.
4. LASSO \rightarrow feature selection.
5. L₁ (LASSO)
more sparsity
feature selection.

L₂(RIDGE)

Makes weights small (not 0)

More sensitive to outliers (due to square)

usual works better in practice

Uncertainties:

1. If Elastic Net allows us to combine L₁ & L₂ together to do the regularization, why not usily elastic net instead of using L₁ or L₂ separately?