

# LN12. ML Debugging, Over or Underfitting

Kun He (何琨)

Data Mining and Machine Learning Lab  
(John Hopcroft Lab)  
Huazhong University of Science & Technology

*brooklet60@hust.edu.cn*

2022年5月



# Table of Contents

- 1 Make Sure Your Model is Optimally Tuned
- 2 Overfitting and Underfitting
- 3 Identify the Sweet Spot
  - Cross-Validation
  - Telescopic Search
- 4 Early Stopping

# Table of Contents

## 1 Make Sure Your Model is Optimally Tuned

## 2 Overfitting and Underfitting

## 3 Identify the Sweet Spot

- Cross-Validation
- Telescopic Search

## 4 Early Stopping

# Make Sure Your Model is Optimally Tuned

Motivating example:

- Remember Empirical Risk Minimization (ERM):

$$\min_{\mathbf{w}} \underbrace{\frac{1}{n} \sum_{i=1}^n l_{(s)}(h_{\mathbf{w}}(\mathbf{x}_i), y_i)}_{\text{Loss}} + \underbrace{\lambda r(\mathbf{w})}_{\text{Regularizer}}$$

- How should we set  $\lambda$  (regulates the bias/variance)?

# Table of Contents

1 Make Sure Your Model is Optimally Tuned

2 Overfitting and Underfitting

3 Identify the Sweet Spot

- Cross-Validation
- Telescopic Search

4 Early Stopping

# Overfitting and Underfitting

- There are two equally problematic cases which can arise when learning a classifier on a data set: **underfitting and overfitting**, each of which relate to the degree to which the data in the training set is extrapolated to apply to unknown data:

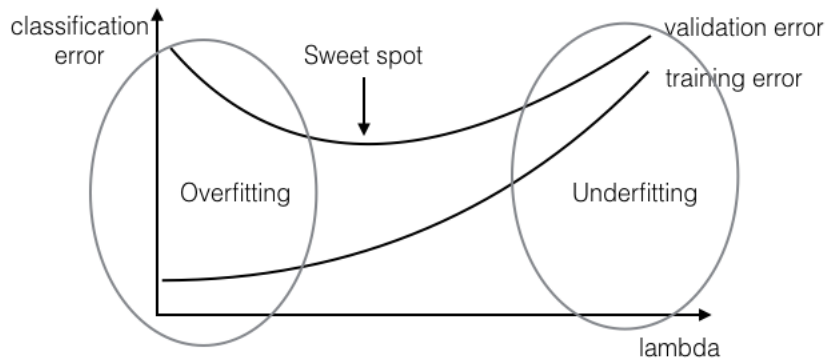
## Underfitting

- The classifier learned on the training set is not expressive enough to even account for the data provided.
- In this case, both the training error and the test error will be high, as the classifier does not account for relevant information present in the training set.

## Overfitting

- The classifier learned on the training set is too specific, and cannot be used to accurately infer anything about unseen data.
- Although training error continues to decrease over time, test error will begin to increase again as the classifier begins to make decisions based on patterns which exist only in the training set and not in the broader distribution.

# Overfitting and Underfitting



# Table of Contents

1 Make Sure Your Model is Optimally Tuned

2 Overfitting and Underfitting

3 Identify the Sweet Spot

- Cross-Validation
- Telescopic Search

4 Early Stopping



## Cross-Validation

Divide data into training and validation portions. Train your algorithm on the "training" split and evaluate it on the "validation" split, for various value of  $\lambda$  (Typical values:  $10^{-5}$   $10^{-4}$   $10^{-3}$   $10^{-2}$   $10^{-1}$   $10^0$   $10^1$   $10^2$  ...).

## k-Fold Cross-Validation

Divide your training data into  $k$  partitions. Train on  $k-1$  of them and leave one out as validation set. Do this  $k$  times (i.e. leave out every partition exactly once) and average the validation error across runs. This gives you a good estimate of the validation error (even with standard deviation).



## K-Fold Validation Error

### k-fold validation error

K-fold validation error is defined as:

$$\text{K-fold validation error} : \frac{1}{n} \sum_{k=1}^K \sum_{i \in D_k} \mathbf{1}(h(\mathbf{x}_i; D - D_k) \neq y_i)$$

where  $D - D_k$  is the dataset minus  $k^{\text{th}}$  fold and  $h(\mathbf{x}_i; D - D_k)$  is the classifier learned on  $D - D_k$ .

# Leave-One-Out Cross-Validation

In the extreme case, you can have  $k=n$ , i.e. you only leave a single data point out (this is often referred to as LOOCV- Leave One Out Cross Validation). LOOCV is important if your data set is small and cannot afford to leave out many data points for evaluation .

## Leave-one-out error

Leave-one-out error is defined as:

$$\text{LOO error} : \frac{1}{n} \sum_i \mathbf{1}(h(\mathbf{x}_i; D_{-i}) \neq y_i)$$

where  $D_{-i}$  is the dataset minus  $i^{\text{th}}$  example and  $h(\mathbf{x}_i; D_{-i})$  is the classifier learned on  $D_{-i}$ .

## Telescopic Search

Do two searches:

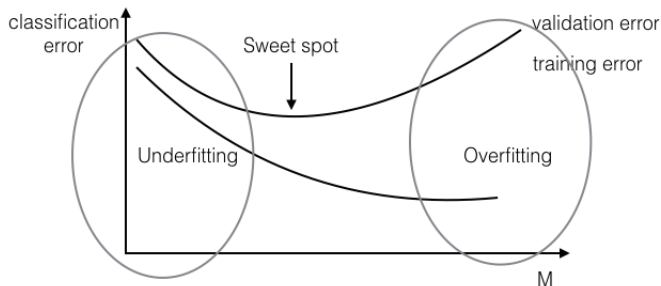
- 1st, find the best order of magnitude for  $\lambda$ ;
  - 2nd, do a more fine-grained search around the best  $\lambda$  found so far.
- 
- For example, first you try  $\lambda = 0.01, 0.1, 1, 10, 100$ . It turns out 10 is the best performing value.
  - Then you try out  $\lambda = 5, 10, 15, 20, 25, \dots, 95$  to test values "around" 10.

# Table of Contents

- 1 Make Sure Your Model is Optimally Tuned
- 2 Overfitting and Underfitting
- 3 Identify the Sweet Spot
  - Cross-Validation
  - Telescopic Search
- 4 Early Stopping

## Early Stopping

Stop your optimization after  $M$  ( $\geq 0$ ) number of gradient steps, even if optimization has not converged yet.



# The End