

DDA5001 Machine Learning

Overfitting (Part I)

Xiao Li

School of Data Science
The Chinese University of Hong Kong, Shenzhen



Recap: Gradient Descent with Momentum

GD:

$$\theta_{k+1} = \theta_k - \mu_k \nabla \mathcal{L}(\theta_k)$$

A quite popular technique for accelerating gradient descent method is the **momentum** technique:

$$\theta_{k+1} = \theta_k - \mu_k \nabla \mathcal{L}(\theta_k) + \underbrace{\beta_k (\theta_k - \theta_{k-1})}_{\text{momentum}}$$

An equivalent form:

$$\begin{aligned}\theta_{k+1} &= \theta_k - m_k \\ m_k &= \mu_k \nabla \mathcal{L}(\theta_k) + \beta_k m_{k-1}\end{aligned}$$

- ▶ Each iteration takes nearly the same time cost as GD.
- ▶ Widely used in practice, notably in **Adam** algorithm.

Recap: Nesterov's Acceleration

GD:

$$\theta_{k+1} = \theta_k - \mu_k \nabla \mathcal{L}(\theta_k)$$

Another very useful technique to accelerate GD is **Nesterov's accelerated gradient descent (AGD)**:

$$\begin{aligned}\theta_{k+1} &= w_k - \mu_k \nabla \mathcal{L}(w_k) \\ w_k &= \theta_k + \frac{k-1}{k+2}(\theta_k - \theta_{k-1})\end{aligned}$$

- ▶ Each iteration takes nearly the same time cost as GD.
- ▶ Widely used in practice.

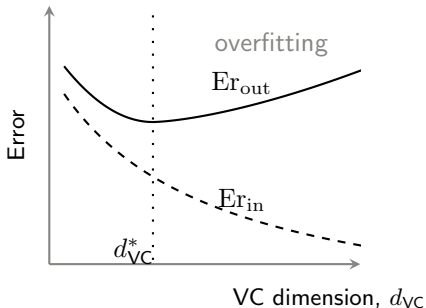
Overfitting

Validation

What is Overfitting?

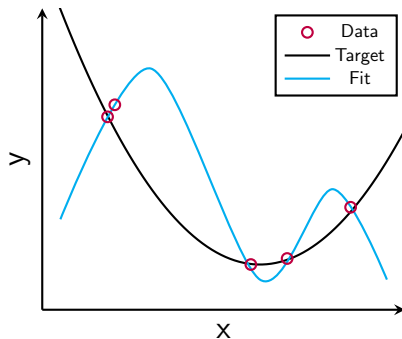
Overfitting

Fitting the data more than is needed



- ▶ One possible reason: \mathcal{H} is more complex than is needed.
- ▶ It means fitting the observed data well no longer indicates that we will get a small out-of-sample error, and may actually lead to the opposite effect. The main case: Small training error (small Er_{in}) but bad generalization (large Er_{out}).

Overfitting: Example



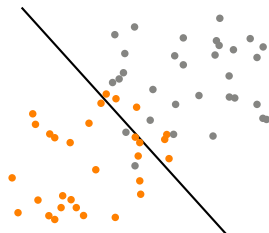
- ▶ Quadratic target function g .
- ▶ 5 data points with noise.
- ▶ 4-th order polynomial hypothesis set \mathcal{H}

$$Er_{in} = 0, \quad Er_{out} \text{ is huge}$$

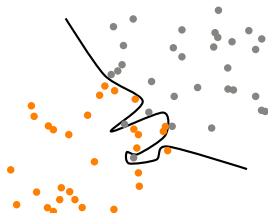
- ▶ The model uses its additional degrees of freedom to fit non-desirable pattern in the data (for example, noise), yielding a final learned model that is inferior.

Overfitting: Example

Linear: nearly appropriate



Highly nonlinear: overfit



- ▶ We have linearly separable data plus noise, it becomes slightly non-linearly separable.
- ▶ Linear classifier

$E_{\text{in}} \neq 0$, E_{out} is appropriate

- ▶ **Highly nonlinear** hypothesis set \mathcal{H}

$E_{\text{in}} = 0$, E_{out} is larger

Overfitting: Catalysts

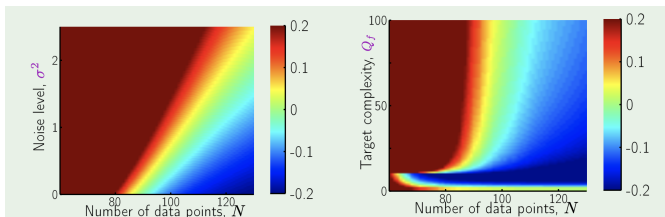


Figure: Color means the overfitting measure $\text{Er}_{\text{out}}(f_{10}) - \text{Er}_{\text{out}}(f_2)$.

Overfitting: Catalysts

- ▶ **Number of training samples** increase, overfitting decreases.
- ▶ **Noise** in data increase, overfitting increases.
- ▶ **Target model complexity** increases, overfitting increases.

⇒ Our HW2 will study the Catalysts for overfitting.

- ▶ One important approach to avoid overfitting is validation.
- ▶ Another one is regularization.

Overfitting

Validation

Validation as A Cure for Overfitting

From VC generalization analysis, we have

$$\text{Er}_{\text{out}}(f) \leq \text{Er}_{\text{in}}(f) + \text{overfit penalty.}$$

- ▶ Generalization error is large in the overfitting regime, hence it becomes the “overfit penalty”.
- ▶ The ideal case is to minimize Er_{out} directly, which is, however, not available.
- ▶ Can we turn to estimate the Er_{out} ?

Validation technique tries to estimate the out-of-sample error for eliminating overfitting:

$$\underbrace{\text{Er}_{\text{out}}(f)}_{\text{validation estimates this quantity}} \leq \text{Er}_{\text{in}}(f) + \text{overfit penalty.}$$

Where we have estimated Er_{out} ? Q5 in HW1. Why is it possible? Will justify later in the understanding of validation error.

Applications of Validation

Hyper-parameters: Something that are **NOT** automatically determined by the learning algorithm.

- ▶ Complexity of \mathcal{H} .
- ▶ Number of iterations in learning algorithm, i.e., stopping criterion.
- ▶ Learning rate.
- ▶ Regularization parameter λ (in our later lecture).
- ▶ ...

In many cases, we have one (or more) hyper-parameters. The value chosen for hyper-parameters has a significant impact on the algorithm's output.

Model selection:

- ▶ The problem of selecting values for hyper-parameters is called **model selection**.
- ▶ The **validation** is used for model selection.

The Validation Dilemma

- ▶ What we have? **Training data**.
- ▶ Can we use the training data to estimate the Er_{out} and then select hyper-parameters?
- ▶ These hyper-parameters usually control the balance between **underfitting** and **overfitting**.
- ▶ Thus, it is not appropriate to let the **training data** influence the selection of hyper-parameters, as **this almost always leads to overfitting**.

Example: If we let the training data determine the degree in polynomial regression (\mathcal{H} is the set of all polynomial functions with degree less than d), we will just end up choosing the maximum degree and doing interpolation (overfitting).

\leadsto We need some other data for estimating Er_{out} .

The Idea of Validation: Split Training Data

The idea

Split the training set to another 'training set' and a validation set.

Validation

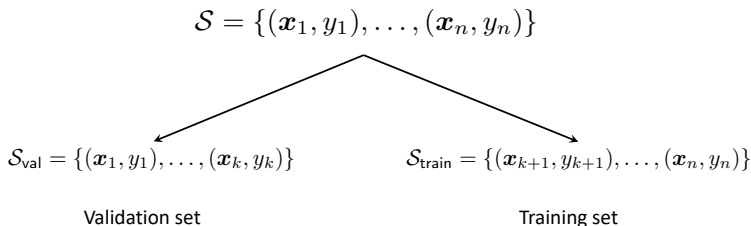
After learning $f_{\theta'} \in \mathcal{H}$ based on the new 'training set', use validation set to estimate $\text{Er}_{\text{out}}(f_{\theta'})$.

Why estimate $\text{Er}_{\text{out}}(f_{\theta'})$? We can see if the learned $f_{\theta'} \in \mathcal{H}$ is good or not in terms of (the estimated) out-of-sample error.

Recall the goal of machine learning is to make the out-of-sample error small. If $\text{Er}_{\text{out}}(f_{\theta'})$ is small, it means we are nearly done. If not, it means we have a bad learned model, either underfitting or overfitting.

The Validation Set

Split training data set:



- ▶ Data set: $\mathcal{S} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$, size n .
- ▶ Validation set: $\mathcal{S}_{\text{val}} = \{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_k, y_k)\}$, size k .
- ▶ New training set: $\mathcal{S}_{\text{train}} = \{(\mathbf{x}_{k+1}, y_{k+1}), \dots, (\mathbf{x}_n, y_n)\}$, size $n - k$.

Validation Error

Use the validation set to form an estimate (we use f' to represent $f_{\theta'}$ to ease notation)

$$\text{Er}_{\text{val}}(f') = \frac{1}{k} \sum_{i \in \mathcal{S}_{\text{val}}} e(f'(\mathbf{x}_i), y_i)$$

How well $\text{Er}_{\text{val}}(f')$ approximates $\text{Er}_{\text{out}}(f')$?

► In expectation

$$\mathbb{E} [\text{Er}_{\text{val}}(f')] = \frac{1}{k} \sum_{i \in \mathcal{S}_{\text{val}}} \mathbb{E} [e(f'(\mathbf{x}_i), y_i)] = \text{Er}_{\text{out}}(f')$$

This implies that the validation error is an **unbiased** estimation of the out-of-sample error of f' .

However, unbiased estimation is a too weak guarantee. Can we have more?

Validation Error Approximates Out-of-Sample Error

Mimicking the derivation of the **generalization result for a single fixed hypothesis**, we can utilize the Hoeffding's inequality to show that

$$\text{Er}_{\text{out}}(f') \leq \text{Er}_{\text{val}}(f') + \mathcal{O}\left(\frac{1}{\sqrt{k}}\right).$$

- ▶ The key of this result is that we only need to derive the bound for **a single fixed hypothesis f'** , thus the derivation is exactly the same as the fixed f generalization result. No union bounds needed.
- ▶ One key feature is that **$\text{Er}_{\text{val}}(f')$ is computable**.
- ▶ This means that once we have enough validation data, i.e., **k is large**, we get good estimation of $\text{Er}_{\text{out}}(f')$ by $\text{Er}_{\text{val}}(f')$. Then, we will know whether the model f' is good or not.
- ▶ This result applies to **test data** as well, i.e.,

$$\text{Er}_{\text{out}}(f') \leq \text{Er}_{\text{test}}(f') + \mathcal{O}(1/\sqrt{k})$$

if we have k test data points. The reasoning is exactly the same.

Any remaining issues?

Trade-off in Validation

- Where is k from? Where is f' from?

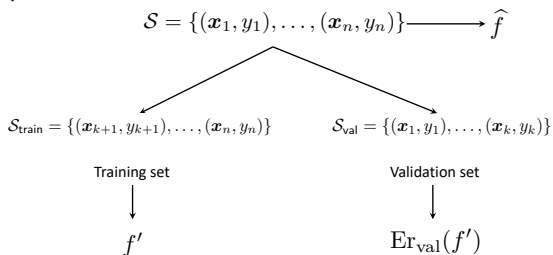
$$\underbrace{n}_{\text{number of data}} = \underbrace{k}_{\text{number of validation data}} + \underbrace{n - k}_{\text{number of new training data}}$$

- $\text{Er}_{\text{val}}(f')$ is $\mathcal{O}(1/\sqrt{k})$ -close to $\text{Er}_{\text{out}}(f')$, while f' is learned by using $n - k$ new training data.
- Small k : Bad estimate of $\text{Er}_{\text{out}}(f')$.
- Large k : Good estimate of $\text{Er}_{\text{out}}(f')$, but bad f' (since few training samples).

Rule of thumb: $k = n/5$.

Restoring: k is Put Back to n

- Validation is to estimate $\text{Er}_{\text{out}}(f')$ using $\text{Er}_{\text{val}}(f')$, must we output f' as the learned model?
- No. It is mainly used to know how to choose **hyper-parameters** such as the best possible hypothesis and learning rate.
- Hence, after we have used the validation set to estimate the out-of-sample error, **re-train on the whole data set**.



By using the previous bound, we have

$$\text{Er}_{\text{out}}(\hat{f}) \leq \text{Er}_{\text{out}}(f') \leq \text{Er}_{\text{val}}(f') + \mathcal{O}\left(\frac{1}{\sqrt{k}}\right).$$

The gray part is a **reasonable guess** (not proof) from VC analysis.

▪

Model Selection

Validation for Model Selection

We have discussed how to estimate $\text{Er}_{\text{out}}(f')$ using validation. It can be important to estimate how good our f' will perform on the test data.

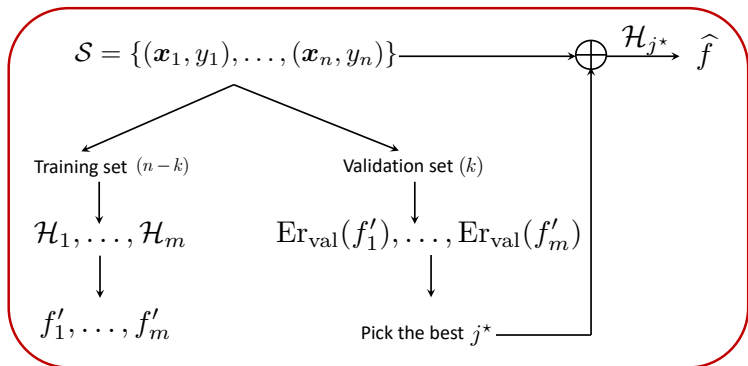
- ▶ Indeed, validation can be used to **guide the learning process** systematically, by the way of **model selection**.
- ▶ Recall that model selection is to select **hyper-parameters**. These could be the choice between a linear model and a nonlinear model, the choice of order of polynomials in a model, the choice of learning rate, or any other choice that affects the learning process.
- ▶ **Setup:** Suppose we have m candidate hypothesis spaces (can also be m different learning rate choices, etc)

$$\mathcal{H}_1, \dots, \mathcal{H}_m.$$

We can use the validation set to estimate the out-of-sample error by using $\text{Er}_{\text{val}}(f'_j)$ for each f'_j learned from these model spaces.

- ▶ **Selection:** Choose j^* such that $\text{Er}_{\text{val}}(f'_{j^*}) \leq \text{Er}_{\text{val}}(f'_j)$ for all j .
- ▶ **Restoring:** Train f on the whole set using model space \mathcal{H}_{j^*} , get \hat{f} .

Validation for Model Selection: illustration



Generalization Error of Model Selection

The model selection process gives a new hypothesis space consists of

$$\mathcal{H}_{\text{val}} = \{f'_1, f'_2, \dots, f'_m\}.$$

Model selection chooses **one** from \mathcal{H}_{val} that achieves the smallest validation error.

- ▶ This process is equivalent to learn a model from \mathcal{H}_{val} using the validation set, where Er_{val} is the “in-sample” error.
- ▶ This setting allows us to apply the generalization analysis for **finite hypotheses space** case, since \mathcal{H}_{val} only consists of m hypotheses, i.e., $|\mathcal{H}_{\text{val}}| = m$.
- ▶ This gives

$$\text{Er}_{\text{out}}(\hat{f}) \leq \text{Er}_{\text{out}}(f'_{j^*}) \leq \text{Er}_{\text{val}}(f'_{j^*}) + \mathcal{O}\left(\sqrt{\frac{\log m}{k}}\right).$$

Again, the gray part is a reasonable guess from VC analysis.

Validation vs. Testing

- ▶ We call this “validation”, but how is it different from “testing”?
- ▶ Typically, validation is used to **make learning choices**, i.e., choosing hyper-parameters to avoid overfitting.

However,

The test data can never influence the training phase in any way.

If it impacts the learning process, i.e., which final $\hat{f} \in \mathcal{H}$ we choose, then **it is no longer a test set,**

it becomes a validation set.

