Overfitting, Cross Validation, and Regularization

Statistical Machine Learning (ENGG\*6600\*02)

School of Engineering, University of Guelph, ON, Canada

Course Instructor: Benyamin Ghojogh Summer 2023 Measures for a Model

#### Learning Model

$$f: a_i \mapsto f(a_i) = f_i$$

• Assume we have a function f which gets the *i*-th input  $x_i$  and outputs  $f_i = f(x_i)$ . This figure shows this function and its input and output:



- We wish to know the function which we call it the <u>true model</u> but we do <u>not have access</u> to it as it is unknown.
- Also, the pure outputs (true observations), f<sub>i</sub>'s, are not available. The output may be corrupted with an additive noise ε<sub>i</sub>:

$$y_i = f_i + \varepsilon_i, \qquad (1)$$

where the noise is  $\varepsilon_i \sim \mathcal{N}(0, \sigma^2)$ .

## Learning Model

We have:

$$\operatorname{Var}(\varepsilon_{i}) = \operatorname{E}(\varepsilon_{i}^{2}) - \left(\operatorname{E}(\varepsilon_{i})\right)^{2} \mathbb{Z}$$

$$\underbrace{y_i = f_i + \varepsilon_i, \varepsilon_i \sim \mathcal{N}(0, \sigma^2)}_{\text{The true observation } f_i \text{ is not random, thus:}} \xrightarrow{\mathbb{E}(\varepsilon_i) = 0,}_{\mathbb{E}(\varepsilon_i) = 0,} \underbrace{\mathbb{E}(\varepsilon_i) = \mathbb{E}(\varepsilon_i) + \mathbb{E}(\varepsilon_i)^2}_{\mathbb{E}(f_i) = f_i.} \xrightarrow{\mathbb{E}(\varepsilon_i) = 0,}_{\mathbb{E}(\varepsilon_i) = 0,} \underbrace{\mathbb{E}(\varepsilon_i) = \mathbb{E}(\varepsilon_i) + \mathbb{E}(\varepsilon_i)^2}_{\mathbb{E}(\varepsilon_i) = 0,} \underbrace{\mathbb{E}(\varepsilon_i) = 0,}_{\mathbb{E}(\varepsilon_i) = 0,} \underbrace{\mathbb{E}(\varepsilon_i) = 0$$

The input training data {x<sub>i</sub>}<sup>n</sup><sub>i=1</sub> and their corrupted observations {y<sub>i</sub>}<sup>n</sup><sub>i=1</sub> are available to us. We would like to approximate (estimate) the true model by a model f in order to estimate the observations {y<sub>i</sub>}<sup>n</sup><sub>i=1</sub> from the input {x<sub>i</sub>}<sup>n</sup><sub>i=1</sub>.

- Calling the estimated observations by  $\{\hat{y}_i\}_{i=1}^n$ , we want the  $\{\hat{y}_i\}_{i=1}^n$  to be as close as possible to  $\{y_i\}_{i=1}^n$  for the training input data  $\{\mathbf{x}_i\}_{i=1}^n$ .
  - We train the model using the training data in order to estimate the true model.
  - After training the model, it can be used to estimate the output of the model for both the training input  $\{\mathbf{x}_i\}_{i=1}^n$  and the unseen test input  $\{\mathbf{x}_i\}_{i=1}^m$  to have the estimates  $\{\widehat{y}_i\}_{i=1}^n$  and  $\{\widehat{v}_i\}_{i=1}^m$ , respectively.



#### Learning Model

- Here, we denote the estimation of the observation of the *i*-th instance with either  $\hat{y}_i$  or  $\hat{f}_i$ .
- The model can be a <u>regression (prediction)</u> or <u>classification</u> model. In regression, the model's estimation is <u>continuous</u> while in classification, the estimation is a member of a discrete <u>set of possible</u> observations.
- The definitions of variance, bias, and MSE can also be used for the estimation  $\hat{f}_i$  of the true model  $f_i$ .

$$\underbrace{\mathsf{MSE}(\widehat{f}) = \mathbb{Var}(\widehat{f}) + (\mathbb{Bias}(\widehat{f}))^2}_{\mathbb{Bias}(\widehat{f})} \Longrightarrow \left( \sqrt{\mathsf{MSE}(\widehat{f})} \right)^2 = \left( \sqrt{\mathbb{Var}(\widehat{f})} \right)^2 + (\mathbb{Bias}(\widehat{f}))^2.$$
(4)

Mean Squared Error of the Estimation of Observations

#### Mean Squared Error of the Estimation of Observations

• Suppose we have an instance 
$$(x_0, y_0)$$
. This instance can be either a training or  
test/validation instance. We will cover both cases.  
• According to Eq. (1), the observation  $y_0$  is:  $\mathcal{E}_{o} = \mathcal{Y}_{o} - \hat{f}_{o}$   $(\mathcal{Y}_{o}, \mathcal{Y}_{o})$   $\hat{f}_{o} + \mathcal{E}_{o} = \mathcal{I}_{o}$   
 $(a-b)^2 = a^2 + b^2 - 2ab$   $(5)$   
• Assume the model's estimation of  $y_0$  is  $\hat{f}_0$ . The MSE of the estimation is:  
 $M = \mathcal{I}$   $(\hat{f}_0 - \hat{f}_0)^2$   $\stackrel{(5)}{=} \mathbb{E}((\hat{f}_0 - f_0 - \hat{\epsilon}_0)^2) = \mathbb{E}((\hat{f}_0 - f_0)^2 + \hat{\epsilon}_0^2 - 2\hat{\epsilon}_0(\hat{f}_0 - f_0))$   
 $= \mathbb{E}((\hat{f}_0 - f_0)^2) + \mathbb{E}(\hat{\epsilon}_0^2) - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}((\hat{f}_0 - f_0)^2) + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}((\hat{f}_0 - f_0)^2) + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0)))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_0(\hat{f}_0 - f_0))$   $\stackrel{(6)}{=} \mathbb{E}(\hat{f}_0 - f_0)^2 + \hat{\sigma}^2 - 2\mathbb{E}(\hat{\epsilon}_$ 

$$\mathbb{E}\left(\widehat{f_0}-\widehat{f_0}\right) \stackrel{(5)}{=} \mathbb{E}\left((y_0-f_0)(\widehat{f_0}-f_0)\right). \quad \bigstar$$

$$(7)$$

For calculation of this term, we have two cases: (I) whether the instance (x<sub>0</sub>, y<sub>0</sub>) is in the training set or (II) not in the training set. In other words, whether the instance was used to train the model (estimator) or not.

- Assume the instance (x<sub>0</sub>, y<sub>0</sub>) was not in the training set, i.e., it was not used for training the model. In other words, we have y<sub>0</sub> ∉ T.
- This means that the estimation  $\hat{f_0}$  is independent of the observation  $y_0$  because the observation was not used to train the model but the estimation is obtained from the model. Therefore:

where (a) is because  $(y_0 - f_0) \perp (\hat{f_0} - f_0)$  and (b) is because:

$$\mathbb{E}((y_0 - f_0)) = \mathbb{E}(y_0) - \mathbb{E}(f_0) \stackrel{(c)}{=} f_0 - f_0 = 0,$$

where (c) is because of Eq. (3) and:

t

$$y_{s} = f_{0+} \varepsilon_{s} \implies \varepsilon_{(y_{0})} \stackrel{(5)}{=} \varepsilon_{(f_{0})} + \varepsilon_{(\varepsilon_{0})} = f_{0} + 0 = f_{0}.$$

Therefore, in this case, the last term in Eq. (6) is zero. Thus:

$$\mathsf{MSE} \longrightarrow \mathbb{E}((\widehat{f}_0 - y_0)^2) = \mathbb{E}((\widehat{f}_0 - f_0)^2) + \sigma^2$$
(8)

XILY => E(XY)=E(X+i)

 $E(x) \tilde{E}(y)$ 

1

# Case I: Instance not in the Training Set $(\mathcal{E}(X) \stackrel{\bullet}{\rightarrowtail} \stackrel{1}{\underset{j=1}{\overset{m}{\longrightarrow}}} \stackrel{\mathcal{R}_{t}}{\underset{j=1}{\overset{m}{\longrightarrow}}} \mathcal{R}_{t}$ • We found: $\mathcal{E}(g(X)) \stackrel{\bullet}{\rightharpoonup} \stackrel{1}{\underset{j=1}{\overset{m}{\longrightarrow}}} \stackrel{\mathcal{R}_{t}}{\underset{m}{\overset{m}{\longrightarrow}}} \mathcal{R}_{t}$

t

$$\bigstar \quad \underbrace{\mathbb{E}((\widehat{f}_0 - y_0)^2)}_{f_0} = \underbrace{\mathbb{E}((\widehat{f}_0 - f_0)^2)}_{f_0} + \sigma^2$$

Suppose the number of instances which are not in the training set is <u>m</u>. By Monte Carlo approximation of the expectation terms, we have:

$$\underbrace{\frac{1}{m}\sum_{i=1}^{m}(\widehat{f_{i}}-y_{i})^{2}}_{i=1} = \underbrace{\frac{1}{m}\sum_{i=1}^{m}(\widehat{f_{i}}-f_{i})^{2}}_{i=1} + \sigma^{2} \xrightarrow{\text{MN}}_{i=1} \underbrace{\sum_{i=1}^{m}(\widehat{f_{i}}-y_{i})^{2}}_{i=1} = \underbrace{\sum_{i=1}^{m}(\widehat{f_{i}}-f_{i})^{2}}_{i=1} + \underline{m\sigma^{2}}.$$
 (9)

- The term  $\sum_{i=1}^{m} (\hat{i}_i y_i)^2$  is the error between the predicted output and the <u>label</u> in the dataset. So, it is the <u>empirical error</u>, denoted by <u>err</u>.
- The term  $\sum_{i=1}^{m} (\hat{f}_i f_i)_i^2$  is the error between the predicted output and true unknown label. This error is referred to as **true error**, denoted by **Err**. Therefore:

• The term  $m\sigma^2$  is a constant and can be ignored. Hence, in this case, the empirical error is a good estimation of the true error. Thus, we can minimize the empirical error in order to properly minimize the true error.

- In case 2, the instance is in the training set. For this case, we need to use a mathematical formula named <u>SURE</u>, introduced in the following.
- Consider a multivariate random variable ℝ<sup>d</sup> ∋ z = [z<sub>1</sub>,..., z<sub>d</sub>] whose components are independent random variables with normal distribution, i.e., z<sub>i</sub> ~ N(μ<sub>i</sub>, σ).
- Take  $\mathbb{R}^d \ni \boldsymbol{\mu} = [\mu_1, \dots, \mu_d]^\top$  and let  $\mathbb{R}^d \ni \boldsymbol{g}(\boldsymbol{z}) = [g_1, \dots, g_d]^\top$  be a function of the random variable  $\boldsymbol{z}$  with  $\boldsymbol{g}(\boldsymbol{z}) : \mathbb{R}^d \to \mathbb{R}^d$ .
- There exists a lemma, named Stein's Lemma, which states:

$$\bigstar \left[ \mathbb{E}((\boldsymbol{z} - \boldsymbol{\mu})^{\top} \boldsymbol{g}(\boldsymbol{z})) = \sigma^2 \sum_{i=1}^{d} \mathbb{E}(\frac{\partial \boldsymbol{g}_i}{\partial \boldsymbol{z}_i}), \right]$$
(11)

which is used in <u>Stein's Unbiased Risk Estimate (SURE)</u> [1]. See our <u>tutorial</u> [2] for the proof of Eq. (11).

• If the random variable is a <u>univariate variable</u>, the Stein's lemma becomes:

$$\bigstar \left[ \mathbb{E}((z-\mu)g(z)) = \sigma^2 \mathbb{E}(\frac{\partial g(z)}{\partial z}) \right]$$
(12)

$$y_{j} = f_{0} + \varepsilon_{0}$$

0

• <u>SURE</u> for univariate variable:

$$\bigstar \quad \widetilde{\mathbb{E}((\underline{z}-\mu)g(z))} = \sigma^2 \mathbb{E}(\frac{\partial g(z)}{\partial z}).$$

• In the SURE formula for univariate variable, we take  $\varepsilon_0$ , 0, and  $\hat{f}_0 - f_0$  as the z,  $\mu$ , and  $\underline{g(z)}$ , respectively. We do this to make Eq. (7).

• Using Eq. (12), the last term in Eq. (6) is:

$$\underbrace{\mathbb{E}((\varepsilon_{0}-0)(\widehat{f_{0}}-f_{0}))}_{3\mathcal{H}} = \sigma^{2} \mathbb{E}(\frac{\partial(\widehat{f_{0}}-f_{0})}{\partial\varepsilon_{0}}) = \sigma^{2} \mathbb{E}(\frac{\partial\widehat{f_{0}}}{\partial\varepsilon_{0}} - \frac{\partial\widehat{f_{0}}}{\partial\varepsilon_{0}}) \stackrel{(a)}{=} \sigma^{2} \mathbb{E}(\frac{\partial\widehat{f_{0}}}{\partial\varepsilon_{0}}) \stackrel{(a)}{=} \sigma^{2} \mathbb{E}(\frac{\partial\widehat{f_{0}}}{\partial\varepsilon_{0}}),$$

where (a) is because the true model f is not dependent on the noise, (b) is because of the chain rule in derivative, and (c) is because:

$$\underbrace{y_0 \stackrel{(5)}{=} f_0 + \varepsilon_0}_{y_0 \stackrel{(5)}{=} f_0 + \varepsilon_0} \implies \frac{\partial y_0}{\partial \varepsilon_0} = 1.$$

• Therefore, in this case, the Eq. (6) is:

$$\mathsf{MSE} \longrightarrow \mathbb{E}((\widehat{f_0} - y_0)^2) = \mathbb{E}((\widehat{f_0} - f_0)^2) + \sigma^2 - 2\sigma^2 \mathbb{E}(\frac{\partial \widehat{f_0}}{\partial y_0}).$$
(13)

We had:

• Suppose the number of training instances is <u>n</u>. By <u>Monte Carlo approximation</u> of the expectation terms, we have:

$$\underbrace{\frac{1}{n}\sum_{i=1}^{n}(\widehat{f_{i}}-y_{i})^{2}}_{\sum_{i=1}^{n}(\widehat{f_{i}}-f_{i})^{2}+n\sigma^{2}-2\sigma^{2}\frac{1}{n}\sum_{i=1}^{n}\frac{\partial\widehat{f_{i}}}{\partial y_{i}} \stackrel{\textbf{Xh}}{\Longrightarrow} \sum_{i=1}^{n}(\widehat{f_{i}}-y_{i})^{2} = \sum_{i=1}^{n}(\widehat{f_{i}}-f_{i})^{2}+n\sigma^{2}-2\sigma^{2}\sum_{i=1}^{n}\frac{\partial\widehat{f_{i}}}{\partial y_{i}}.$$
(14)

- The term ∑<sub>i=1</sub><sup>m</sup> (*i*<sub>i</sub> − y<sub>i</sub>)<sup>2</sup> is the error between the predicted output and the label in the dataset. So, it is the empirical error, denoted by err.
- The term  $\sum_{i=1}^{m} (\hat{f}_i \bar{f}_i)^2$  is the error between the predicted output and true unknown label. This error is referred to as **true error**, denoted by **<u>Err</u>**. Therefore:

$$\operatorname{err} = \operatorname{Err} + n\sigma^{2} - 2\sigma^{2}\sum_{i=1}^{n} \frac{\partial \widehat{f}_{i}}{\partial y_{i}} \Longrightarrow \operatorname{Err} = \operatorname{err} - n\sigma^{2} + 2\sigma^{2}\sum_{i=1}^{n} \frac{\partial \widehat{f}_{i}}{\partial y_{i}}.$$
(15)

We had:

$$\mathbf{Err} = \mathbf{err} - n\,\sigma^2 + 2\,\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{f_i}}{\partial y_i},$$

- The last term in this equation is a <u>measure of complexity</u> (or <u>overfitting</u>) of the model. Note that ∂f<sub>i</sub>/∂y<sub>i</sub> means if we move the <u>i-th training instance</u>, how much the model's estimation of that instance will change? This shows how much the model is complex or overfitted.
- For better understanding, suppose a line regressing a training set via least squares problem. If we change a point, the line will not change significantly because the model is not complex (is <u>underfitted</u>). On the other hand, consider a regression model passing through "all" the points. If we move a training point, the regressing curve changes noticeably which is because the model is very complex (<u>overfitted</u>).







- According to this equation, in the case where the instance is in the training set, the empirical error is not a good estimation of the true error.
- The reason is that minimization of **err** usually increases the complexity of the model, cancelling out the minimization of **Err** after some level of training.

Overfitting, Underfitting, and Generalization

#### Overfitting, Underfitting, and Generalization



- If the model is trained in an <u>extremely simple</u> way so that its estimation has <u>low variance</u> but <u>high bias</u>, we have <u>underfitting</u>. Note that underfitting is also referred to as <u>over-generalization</u>.
- On the other hand, if the model is trained in an <u>extremely complex</u> way so that its estimation has high <u>variance</u> but low bias, we have <u>overfitting</u>.
- To summarize:
  - ▶ in underfitting: low variance, high bias, and low complexity.
  - in overfitting: high variance, low bias, high complexity.



#### Overfitting, Underfitting, and Generalization

- An example for underfitting, good fit, and overfitting is illustrated in this figure.
- As this figure shows, in both underfitting and overfitting, the estimation of a test instance might be very weak while in a good fit, the test instance, which was not seen in the training phase, is estimated well enough with smaller error.
- The ability of the model to estimate the <u>unseen test</u> (out-of-sample) data is referred to as <u>generalization</u>.
- The lack of generalization is the reason why both overfitting and underfitting, especially overfitting, is not acceptable. In overfitting, the training error, i.e., err, is very small while the test (true) error, i.e., Err, is usually awful!



**Cross Validation** 

#### **Cross Validation**

- In order to either (I) find out until which complexity we should train the model or (II) tune the parameters of the model, we should use cross validation [3].
- In cross validation, we divide the dataset  $\mathcal{D}$  into two partitions, i.e., **training set** denoted by  $\mathcal{T}$  and **test set** denoted by  $\mathcal{R}$  where the union of these two subsets is the whole dataset and the intersection of them is the empty set:

$$\begin{array}{c}
\mathcal{T} \cup \mathcal{R} = \mathcal{D}, \\
\mathcal{T} \cap \mathcal{R} = \varnothing.
\end{array}$$
(16)
(17)

- The  $\mathcal{T}$  is used for training the model. After the model is trained, the  $\mathcal{R}$  is used for testing the performance of the model.
- We have different methods for cross validation. Two of the most well-known methods for cross validation are <u>K-fold cross validation</u> and <u>Leave-One-Out Cross Validation</u> (LOOCV).



where |.| denoted the cardinality of set.

• Sometimes, the dataset D is shuffled before the cross validation for better randomization.

- Moreover, both simple random sampling without replacement and stratified sampling [4, 5] can be used for this splitting.
- The <u>K-fold cross validation</u> includes <u>K iterations</u>, where in each of them, one of the partitions is used as the test set and the rest of data is used for training. The overall estimation error is the average test error of iterations.
- We usually have K = 2, 5, 10 in the literature but K = 10 is the most common.

#### K-fold Cross Validation

• The algorithm of K-fold cross validation is shown in the following.



Algorithm 1: K-fold Cross Validation

#### Leave-One-Out Cross Validation



- In Leave-One-Out Cross Validation (LOOCV), we iterate for |D| = N times and in each iteration, we take one instance as the R (so that |R| = 1) and the rest of instances as the training set.
- The overall estimation error is the average test error of iterations.
- Usually, when the size of dataset is <u>small</u>, LOOCV is used in order to use the most of dataset for training and then test the model properly.
- The algorithm of LOOCV is shown in the following.

1 for k from 1 to 
$$|\underline{\mathcal{D}}| = N$$
 do  
2  $|\underline{\mathcal{R}} \leftarrow \text{Take the } \underline{k \cdot \text{th}} \text{ instance from } \mathcal{D}.$   
3  $|\underline{\mathcal{T}} \leftarrow \mathcal{D} \setminus \mathcal{R}.$   
4 Use  $\underline{\mathcal{T}}$  to train the model.  
5  $|\underline{\text{Err}}_{k} \leftarrow \text{Use the trained model to predict } \underline{\mathcal{R}}.$   
6  $|\underline{\text{Err}}_{k} \leftarrow \underline{1}_{k=1}^{|\mathcal{D}|} \sum_{k=1}^{|\mathcal{D}|} \underline{\text{Err}}_{k}$ 

Algorithm 2: Leave-One-Out Cross Validation

#### Cheating #1 in Machine Learning

- The test set and the training set should be disjoint, i.e., T ∩ R = Ø; otherwise, we are introducing the whole or a part of the test instances to the model to learn them.
- Of course, in that way, the model will learn to estimate the test instances easier and better; however, in the real-world applications, the test data is not available at the time of training. Therefore, if we mistakenly have *T* ∩ *R* ≠ Ø, it is referred to as cheating in machine learning (we call it cheating #1 here).

#### Cheating #2 in Machine Learning

• In cross validation with validation set, we have:

$$\star \qquad \underbrace{\mathcal{T} \cap \mathcal{R} = \emptyset}_{\mathbf{X}}, \quad \underbrace{\mathcal{T} \cap \mathcal{V} = \emptyset}_{\mathbf{X}}, \quad \underbrace{\mathcal{V} \cap \mathcal{R} = \emptyset}_{\mathbf{X}}.$$

- The <u>validation and test sets</u> should be disjoint because the parameters of the model should not be optimized by testing on the test set. In other words, in real-world applications, the training and validation sets are available but the test set is not available yet. If we mistakenly have <u>V ∩ R ≠ Ø</u>, it is referred to as <u>cheating</u> in machine learning (we call it cheating #2 here).
- This kind of mistake is very common in the <u>literature unfortunately</u>, where some people optimize the parameters by testing on the test set without having a validation set.
- Moreover, the training and test sets should be disjoint as explained beforehand; otherwise, that would be another kind of cheating in machine learning (introduced before as cheating #1).
- On the other hand, the <u>training and validation sets should be disjoint</u>. Although having *T* ∩ *V* ≠ Ø is not cheating but it should not be done for the reason which will be explained later in this section.
   *T* ∩ *V* = Ø
- To have validation set in cross validation, we usually first split the dataset  $\mathcal{D}$  into  $\underline{\mathcal{T}}'$  and  $\underline{\mathcal{R}}$  where  $\mathcal{T}' \cup \mathcal{R} = \underline{\mathcal{D}}$  and  $\underline{\mathcal{T}}' \cap \underline{\mathcal{R}} = \underline{\varnothing}$ . Then, we split the set  $\underline{\mathcal{T}}'$  into the training and validation sets, i.e.,  $\overline{\mathcal{T}} \cup \mathcal{V} = \overline{\mathcal{T}}'$  and  $\overline{\mathcal{T}} \cap \mathcal{V} = \emptyset$  and usually  $|\overline{\mathcal{T}}| > |\mathcal{V}|$ .
- The algorithms of *K*-fold cross validation and LOOCV can be modified accordingly to include the validation set. In LOOCV, we usually have  $|\mathcal{V}| = 1$ .

(21)

Justification of Overfitting



- When the instance is not in the training set, the true error, Err, and the test error, err, behave differently as shown in Fig. (a).
- At the first stages of training, the **err** and **Err** both decrease; however, after some training, the model becomes more complex and goes toward overfitting. In that stage, the **Err** starts to increase.
- We should end the training when the **Err** starts to increase because that stage is the good fit. Usually, in order to find out when to stop training, we train the model for one stage (e.g., iteration) and then test the trained model on the validation set where the error is named **Err**. This is commonly used in training neural networks [6] where **Err** is measured after every epoch for example.

## Justification of Overfitting



• Recall the Eqs. (10) and (15) where the true error for the test (not in training set) and training instance are related to the training error, respectively:

**Err** = err - m \sigma^2, (instance in training set)  
**Err** = err - m \sigma^2, 
$$2\sigma^2 \sum_{i=1}^n \frac{\partial \widehat{f_i}}{\partial y_i}$$
. (instance not in training set)

• The reason why Err increases after a while of training is according to Eq. (15). Dropping the constant nσ<sup>2</sup> from that expression, we have: Err = err + 2 σ<sup>2</sup> Σ<sub>i=1</sub><sup>n</sup> ∂f<sub>i</sub>/∂y<sub>i</sub> where the term 2 σ<sup>2</sup> Σ<sub>i=1</sub><sup>n</sup> ∂f<sub>i</sub>/∂y<sub>i</sub> shows the model complexity. See Fig. (b) where both err and the model complexity are illustrated as a function of training stages (iterations). According to Eq. (15), the Err is the summation of these two curves which clarifies the reason of its behavior.

## Justification of Overfitting



- That is why we should not train a lot on the training set because the model will get too
  much fitted on the training set and will lose its ability to generalize to new unseen data.
- The Fig. (a) and Eq. (15) show that it is better to have *T* ∩ *V* = Ø. Otherwise, for example if we have *T* = *V*, the Err will be equivalent to err and thus it will go down even in overfitting stages. This is harmful to our training because we will not notice overfitting properly.
- The Eq. (10) also explains that the error on validation or test set is a good measure for the true error. That is why we can use test or validation error in order to know until what stage we can train the model without overfitting.

#### Discussion of Cheating in a Nut Shell

• If we have only training and test sets without validation set:

- $\mathsf{T} \cap \mathcal{R} \neq \varnothing \implies \mathsf{cheating} \ \#1$
- If we have training, test, and validation sets without validation set:

$$\mathcal{T} \cap \mathcal{R} \neq \varnothing \implies \text{cheating } \#1$$

$$\mathcal{V} \cap \mathcal{R} \neq \varnothing \implies$$
 cheating #2

- **\***  $\mathcal{T} \cap \mathcal{V} \neq \emptyset \implies$  harmful to training (not noticing overfitting properly)
- The first two items are <u>advantageous to the model's performance on test data</u> but that is <u>cheating</u> and also it may be disadvantageous to future new test data.
- The third item is <u>disadvantageous to the model's performance on test data</u> because we may not find out overfitting or we may find it out late and the generalization error will become worse; therefore, it is better not to do it.

#### Regularization

#### Regularization: Definition

We can minimize the true error, Err, using optimization. According to Eq. (15), we have:

mhimize 
$$(Err) \implies \mininimize \left[ err - \int_{1}^{n} \sigma^{2} + 2\sigma^{2} \sum_{i=1}^{n} \frac{\partial \widehat{f}_{i}}{\partial y_{i}} \right] \rightarrow Complexity (22)$$

- As the term  $n\sigma^2$  is a constant, we can drop it.
- Calculation of <u>∂f<sub>i</sub>/∂y<sub>i</sub></u> is usually very difficult; therefore, we usually use a penalty term in place of it where the penalty increases as the complexity of the model increases in order to imitate the behavior of ∂f<sub>i</sub>/∂y<sub>i</sub>.
- Therefore, the optimization can be written as a regularized optimization problem:

minimize 
$$\widetilde{J}(\mathbf{x};\theta) := J(\mathbf{x};\theta) + \alpha \, \Omega(\mathbf{x}),$$
 (23)

where  $\theta$  is the parameter(s) of the cost function, J(.) is the objective err to be minimized,  $\Omega(.)$  is the penalty function representing the complexity of model,  $\alpha > 0$  is the regularization parameter, and  $\tilde{J}(.)$  is the regularized objective function.

#### Regularization: Definition

$$\begin{bmatrix} \text{minimize} & \operatorname{Err} := \operatorname{err} - n \sigma^{2} + 2 \sigma^{2} \sum_{i=1}^{n} \frac{\partial \widehat{f_{i}}}{\partial y_{i}}, \\ & & & & \\ & & & & \\$$

- The penalty function can be different things such as log norm [7], l1 norm [8, 9], l2,1 norm [10], etc.
- The <u>l<sub>1</sub> and <u>l<sub>2,1</sub> norms</u> are useful for having <u>sparsity</u> [11, 12]. The sparsity is very effective because of the "bet on <u>sparsity</u>" principal: "Use a procedure that does well in sparse problems, since no procedure does well in dense problems [7, 13]."
  </u>
- The effectiveness of the sparsity can also be explained by <u>Occam's razor</u> [14] stating that "simpler solutions are more likely to be correct than complex ones" or "simplicity is a goal in itself".

\*

**¥** • We are minimizing the Err (i.e.,  $\tilde{J}(x;\theta)$ ) and not err (i.e.,  $J(x;\theta)$ ). As discussed before, minimizing err results in overfitting. Therefore, regularization helps avoid overfitting.

#### $\ell_2$ Regularization

#### Theory for $\ell_2$ Norm Regularization

• The 
$$\ell_2$$
 norm regularization [7]:  
minimize  $\widetilde{J}(\mathbf{x};\theta) := \underbrace{J(\mathbf{x};\theta)}_{\mathcal{X}} + \underbrace{\beta_1(\mathbf{x})}_{\mathcal{Y}} + \underbrace{\beta_2(\mathbf{x})}_{\mathcal{Y}} + \underbrace$ 

• The  $\ell_2$  norm regularization is also referred to as ridge regression or Tikhonov regularization [6].

• Suppose 
$$x^*$$
 is minimizer of the  $J(x; \theta)$ , i.e.:

$$\checkmark \checkmark \left( \nabla J(\boldsymbol{x}^*; \boldsymbol{\theta}) = \boldsymbol{0}. \right)$$

The Taylor series expansion of  $J(\mathbf{x}; \theta)$  up to the second derivative at  $\mathbf{x}^*$  gives:

$$\underbrace{\widehat{J}(\boldsymbol{x};\theta)}_{i} \approx \underbrace{J(\boldsymbol{x}^{*};\theta)}_{i} + \underbrace{\nabla J(\boldsymbol{x}^{*};\theta)}_{i} + \underbrace{\frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}^{*})^{\top} \boldsymbol{H}(\boldsymbol{x} - \boldsymbol{x}^{*})}_{i} = \underbrace{J(\boldsymbol{x}^{*};\theta)}_{i} + \underbrace{\frac{1}{2}(\boldsymbol{x} - \boldsymbol{x}^{*})^{\top} \boldsymbol{H}(\boldsymbol{x} - \boldsymbol{x}^{*})}_{i}, \qquad (26)$$

where  $\boldsymbol{H} \in \mathbb{R}^{d \times d}$  is the Hessian.

(25)

#### Theory for $\ell_2$ Norm Regularization



## Theory for $\ell_2$ Norm Regularization $(AB)^{-1} = B^{-1}A^{-1}$

• If we apply eigenvalue decomposition on the Hessian matrix, we will have:

• Using this decomposition in Eq. (27),  $\mathbf{x}^{\dagger} = (\mathbf{H} + \alpha \mathbf{I})^{-1} \mathbf{H} \mathbf{X}^{*}$ , gives us:

$$\mathbf{x}^{\dagger} = (\underbrace{\mathbf{U}} \underbrace{\mathbf{N}} \underbrace{\mathbf{U}}^{\top} + \underbrace{\mathbf{U}} \alpha \mathbf{I})^{-1} \underbrace{\mathbf{U}} \underbrace{\mathbf{N}} \underbrace{\mathbf{U}}^{\top} \mathbf{x}^{*} \stackrel{(a)}{=} (\mathbf{U} \underbrace{\mathbf{N}} \underbrace{\mathbf{U}}^{\top} + \underbrace{\mathbf{U}} \underbrace{\mathbf{U}} \underbrace{\mathbf{U}} \underbrace{\mathbf{U}}^{\top} \mathbf{I})^{-1} \underbrace{\mathbf{U}} \underbrace{\mathbf{N}} \underbrace{\mathbf{U}}^{\top} \mathbf{x}^{*} = (\underbrace{\mathbf{U}} \underbrace{(\mathbf{N}} + \alpha \mathbf{I}) \underbrace{\mathbf{U}}^{\top} \underbrace{\mathbf{U}} \underbrace{\mathbf{$$

where (a) and (c) are because U is an orthogonal matrix so we have  $U^{-1} = U^{\top}$  which yields to  $U^{\top}U = I$  and  $UU^{\top} = I$  (because U is not truncated). The (b) is because  $\alpha$  is a scalar and can move between the multiplication of matrices.

• The Eq. (29) means that we are rotating  $x^*$  by  $U^{\top}x^*$  but before rotating it back with  $UU^{\top}x^*$ , we manipulate it with the term  $(\Lambda + \alpha I)^{-1}\Lambda$ .

where (a) is because U is an orthogonal matrix and (b) is because U is a non-truncated orthogonal matrix. This means that if we do not have the penalty term, the minimizer of  $\tilde{J}(\mathbf{x};\theta)$  is the minimizer of  $J(\mathbf{x};\theta)$  as expected. In other words, we are rotating the solution  $\mathbf{x}^*$  by  $\mathbf{U}^\top$  and then rotate it back by  $\mathbf{U}$ .

#### Theory for $\ell_2$ Norm Regularization

• If  $\alpha \neq 0$ , the term  $(\mathbf{\Lambda} + \alpha \mathbf{I})^{-1}\mathbf{\Lambda}$  is:



where  $\mathbf{\Lambda} = \operatorname{diag}([\lambda_1, \dots, \lambda_d]^\top)$ . Therefore, for the *j*-th direction of Hessian, we have  $\begin{bmatrix} \lambda_j \\ \overline{\lambda_j + \alpha} \end{bmatrix}$ .

If λ<sub>j</sub> ≫ α, we will have λ<sub>j</sub>/λ<sub>j+α</sub> ≈ 1 so for the j-th direction we have (Λ + αI)<sup>-1</sup>Λ ≈ I; therefore, x<sup>†</sup> ≈ x\*. This makes sense because λ<sub>i</sub> ≫ α means that the j-th direction of Hessian and thus the j-th direction of J(x; θ) is large enough to be effective. Therefore, the penalty is roughly ignored with respect to it.
 If λ<sub>j</sub> ≪ α, we will have λ<sub>j</sub>/λ<sub>j+α</sub> ≈ 0 so for the j-th direction we have (Λ + αI)<sup>-1</sup>Λ ≈ 0; therefore, x<sup>†</sup> ≈ 0. This makes sense because λ<sub>j</sub> ≪ α means that the j-th direction of Hessian and thus the j-th direction of J(x; θ) is small and not effective. Therefore, the penalty shrinks that direction to almost zero.

#### Theory for $\ell_2$ Norm Regularization

- Therefore, the l<sub>2</sub> norm regularization keeps the effective directions but shrinks the weak directions to close to zero.
- The following measure is referred to as <u>effective number of parameters</u> or <u>degree of</u> <u>freedom</u> [7]:

$$\bigstar \qquad \sum_{j=1}^{d} \frac{\lambda_j}{\lambda_j + \alpha}, \tag{30}$$

because it counts the number of effective directions as discussed above. Moreover, the term  $\lambda_j/(\lambda_j + \alpha)$  or  $(\Lambda + \alpha I)^{-1}\Lambda$  is called the shrinkage factor because it shrinks the weak directions.

#### $\ell_1$ Regularization

#### Theory for $\ell_1$ Norm Regularization

As explained before, sparsity is very useful and effective. If x = [x<sub>1</sub>,...,x<sub>d</sub>]<sup>⊤</sup>, for having sparsity, we should use subset selection for the regularization:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \widetilde{J}(\mathbf{x};\theta) := J(\mathbf{x};\theta) + \alpha ||\mathbf{x}||_{0}, \tag{31}$$

where:

$$\boxed{||\mathbf{x}||_{0}} := \sum_{j=1}^{d} \mathbb{I}(x_{j} \neq 0) = \begin{cases} 0 & \text{if } x_{j} = 0, \\ 1 & \text{if } x_{j} \neq 0, \end{cases}$$
(32)

is " $\ell_0$ " norm, which is not a norm (so we used "." for it) because it does not satisfy the norm properties [15]. The " $\ell_0$ " norm counts the number of non-zero elements so when we penalize it, it means that we want to have sparser solutions with many zero entries.

 According to [16], the convex relaxation of "l<sub>0</sub>" norm (subset selection) is l<sub>1</sub> norm. Therefore, we write the regularized optimization as:

$$\underset{\mathbf{x}}{\text{minimize}} \quad \widetilde{J}(\mathbf{x}; \theta) := J(\mathbf{x}; \theta) + \alpha \underbrace{||\mathbf{x}||_{1.}}$$
(33)

The l<sub>1</sub> regularization is also referred to as lasso (least absolute shrinkage and selection operator) regularization [8].

#### Theory for $\ell_1$ Norm Regularization

- Different methods exist for solving optimization having <u>l<sub>1</sub> norm</u>, such as <u>proximal</u> algorithm using soft thresholding [17] and <u>coordinate descent</u> [18, 19]. Here, we explain solving the optimization using the <u>coordinate descent</u> algorithm.
- The idea of <u>coordinate descent</u> algorithm is similar to the idea of <u>Gibbs sampling</u> [20] where we work on the <u>dimensions of the variable</u> one by one. Similar to what we did for obtaining Eq. (27), we have:  $\sum_{k=1}^{\infty} |\mathbf{x}_{k}|^{2} = |\mathbf{x}_{k}|_{F}$

$$\bigstar \quad \widetilde{J}(\mathbf{x};\theta) = \widehat{J}(\mathbf{x};\theta) + \alpha ||\mathbf{x}||_1 = \underbrace{J(\mathbf{x}^*;\theta)}_{2} + \frac{1}{2} (\mathbf{x} - \mathbf{x}^*)^\top H(\mathbf{x} - \mathbf{x}^*)_{+} \alpha ||\mathbf{x}||_1.$$

For simplicity in deriving an interpretable expression, we assume that the Hessian matrix is diagonal [6]. For coordinate descent, we look at the *j*-th coordinate (dimension):

where  $\mathbf{x} = [x_1, \dots, x_d]^\top$ ,  $\mathbf{x}^* = [x_1^*, \dots, x_d^*]^\top$ ,  $h_j$  is the (j, j)-th element of the diagonal Hessian matrix, and c is a constant term with respect to  $x_j$  (not dependent to  $x_j$ ). Taking derivative with respect to  $x_j$  gives us:

$$\frac{\partial \widetilde{J}(x_j;\theta)}{\partial x_j} = 0 + (x_j - x_j^*) h_j + \alpha \operatorname{sign}(x_j) \stackrel{\text{set}}{=} 0 \implies x_j^{\dagger} = x_j^* - \frac{\alpha}{h_j} \operatorname{sign}(x_j) = \begin{cases} x_j^* - \frac{\alpha}{h_j} & \text{if } x_j > 0, \\ x_j^* + \frac{\alpha}{h_j} & \text{if } x_j < 0, \end{cases}$$
which is a soft thresholding function.

#### Theory for $\ell_1$ Norm Regularization



- If  $|x_i^*| < (\alpha/h_j)$ , the solution to the regularized problem, i.e.,  $x_j^{\dagger}$ , is zero. Recall that in  $\ell_2$  norm regularization, we shrank the weak solutions close to zero; however, here in  $\ell_1$  norm regularization, we are setting the weak solutions exactly to zero. That is why the solutions are relatively sparse in  $\ell_1$  norm regularization.
- In l<sub>1</sub> norm regularization, as shown in this figure, even the strong solutions are a little shrunk (from the x<sub>i</sub><sup>†</sup> = x<sub>i</sub><sup>\*</sup> line), the fact that we also had in l<sub>2</sub> norm regularization.

#### Comparison of $\ell_2$ and $\ell_1$ Regularization

- Another intuition for why the <u>l</u> norm regularization is sparse is illustrated below (credit if Tibshirani - <u>1996</u> [8]).
- The objective J(x; θ) has some contour levels like a bowl (if it is convex). The regularization term is also a norm ball, which is a sphere bowl (cone) for ℓ<sub>2</sub> norm and a diamond bowl (cone) for ℓ<sub>1</sub> norm [15].
- For l<sub>2</sub> norm regularization, the objective and the penalty term contact at a point where some of the coordinates might be small; however, for l<sub>1</sub> norm, the contact point can be at some point where some variables are exactly zero. This again shows the reason of sparsity in l<sub>1</sub> norm regularization.



Overfitting, Cross Validation, and Regularizat

#### Acknowledgement

- Some slides are inspired by the textbook: <u>Trevor Hastie, Robert Tibshirani, Jerome</u> Friedman, "The elements of statistical learning: Data Mining, Inference, and Prediction", Springer, 2009 [7].
- Another textbook suitable for sparsity in machine learning is: Robert <u>Tibshirani</u>, Martin Wainwright, Trevor <u>Hastie</u>, "<u>Statistical learning with sparsity: the lasso and generalizations</u>", Chapman and Hall/CRC, 2015 [13].
- Some slides are inspired by the lectures of <u>Prof. Ali Ghodsi</u> at the <u>University</u> of Waterloo, <u>Department of Statistics</u> and Actuarial Science, see his YouTube channel: <u>https://www.youtube.com/@DataScienceCoursesUW</u>
- Some slides are also inspired by the lectures of Prof. Hoda Mohammadzade at the <u>Sharif</u> University of Technology, Department of Electrical Engineering.
- Some slides are based on our <u>tutorial paper</u>: "<u>The Theory Behind Overfitting</u>, Cross Validation, <u>Regularization</u>, <u>Bagging</u>, and Boosting: Tutorial" [2]

#### References

- C. M. Stein, "Estimation of the mean of a multivariate normal distribution," The annals of Statistics, pp. 1135–1151, 1981.
- B. Ghojogh and M. Crowley, "<u>The theory behind overfitting, cross validation</u>, regularization, bagging, and boosting: tutorial," arXiv preprint arXiv:1905.12787, 2019.
- [3] S. Arlot and A. Celisse, "A <u>survey</u> of <u>cross-validation</u> procedures for model selection," *Statistics surveys*, vol. 4, pp. 40–79, 2010.
- [4] V. Barnett, *Elements of sampling theory*. English Universities Press, London, 1974.
- [5] B. Ghojogh, H. Nekoei, A. Ghojogh, F. Karray, and M. Crowley, "Sampling algorithms, from survey sampling to Monte Carlo methods: Tutorial and literature review," arXiv preprint arXiv:2011.00901, 2020.
- [6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*. MIT press, 2016.
- J. Friedman, T. Hastie, and R. Tibshirani, *The elements of statistical learning: Data Mining, Interence, and Prediction*, vol. 2.
   Springer series in statistics, New York, NY, USA, 2009.
- [8] R. Tibshirani, "Regression shrinkage and selection via the lasso," Journal of the Royal Statistical Society: Series B (Methodological), vol. 58, no. 1, pp. 267–288, 1996.

## References (cont.)

- M. Schmidt, "<u>Least squares optimization</u> with <u>l1-norm regularization</u>," CS542B Project Report, vol. 504, pp. 195–221, 2005.
- [10] Y. Chang, "L<sub>2,1</sub> norm and its applications," Technical Report, University of Central Florida.
- [11] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, "<u>Convex optimization with sparsity-inducing norms</u>," *Optimization for Machine Learning*, vol. 5, pp. 19–53, 2011.
- [12] F. Bach, R. Jenatton, J. Mairal, and G. Obozinski, "Optimization with sparsity-inducing penalties," Foundations and Trends (R) in Machine Learning, vol. 4, no. 1, pp. 1–106, 2012.
- [13] R. Tibshirani, M. Wainwright, and T. Hastie, *Statistical learning with sparsity: the lasso and generalizations*. Chapman and Hall/CRC, <u>2015.</u>
- [14] P. Domingos, "The role of Occam's razor in knowledge discovery," Data mining and knowledge discovery, vol. 3, no. 4, pp. 409–425, 1999.
- [15] S. Boyd and L. Vandenberghe, *Convex optimization*. Cambridge university press, 2004.
- [16] D. L. Donoho, "For most large underdetermined systems of linear equations the minimal *l*<sub>1</sub>-norm solution is also the sparsest solution," Communications on Pure and Applied Mathematics: A Journal Issued by the Courant Institute of Mathematical Sciences, vol. 59, no. 6, pp. 797–829, 2006.

## References (cont.)

- [17] N. Parikh and S. Boyd, "Proximal algorithms," Foundations and Trends R in Optimization, vol. 1, no. 3, pp. 127–239, 2014.
- [18] S. J. Wright, "Coordinate descent algorithms," Mathematical Programming, vol. 151, no. 1, pp. 3–34, 2015.
- [19] T. T. Wu and K. Lange, "Coordinate descent algorithms for lasso penalized regression," The Annals of Applied Statistics, vol. 2, no. 1, pp. 224–244, 2008.
- [20] G. Casella and E. I. George, "Explaining the Gibbs sampler," The American Statistician, vol. 46, no. 3, pp. 167–174, <u>1992.</u>