

Assignment Project Exam Help

Predictive Analytics

Week 4: Model Selection and Estimation I

<https://powcoder.com>

Semester 2, 2018

Discipline of Business Analytics, The University of Sydney Business School

Add WeChat powcoder

Week 4: Model Selection

1. Model Selection and Evaluation

2. The bias-variance decomposition

3. KNN bias-variance decomposition

4. Cross validation

5. The bias-variance derivation (details optional)

Reading: Chapter 2.1, 2.2 of ISL.

Exercise questions: Chapter 2.4 of ISL, Q1, Q3, Q5, Q6 and Q7.(a). Only focus on “regression” if the questions contain “classification” content.

Assignment Project Exam Help

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

~~Model Selection and Evaluation~~
<https://powcoder.com>

Add WeChat powcoder

Model Selection and Evaluation

Model Selection: estimate the performance of different models in order to choose the (approximate) best one. We select the model that is estimated to have the best predictive ability.

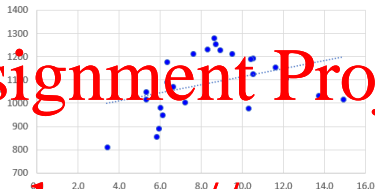
Model Evaluation: after chosen the “different” model, estimate its test error (generalisation error) on new data.

- **Training set:** for exploratory data analysis, model building, model estimation, etc.

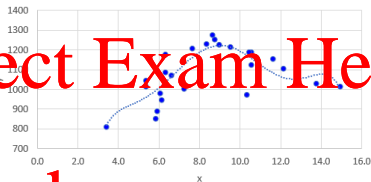
- **Validation set:** for appropriate model selection.

- **Test set:** for model evaluation.

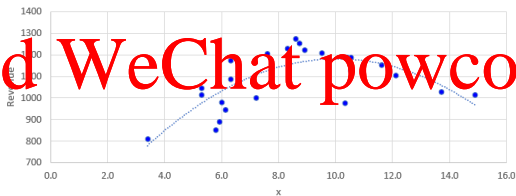
Underfitting and Overfitting



Underfitting: High bias, low variance



Overfitting: High variance, low bias



Best model

Why is overfitting bad?

- Low training error, high generalization error
- Poor predictive performance
- Overreacts to minor fluctuations in training data

How to address overfitting:

- Drop some features
 - Model selection algorithm
 - Manually select features to keep
- Regularization
 - Keep all features, reduce the magnitude/values of parameters.
More details in later sections.

The bias-variance trade-off (key concept)

Assignment Project Exam Help

- Increasing model complexity brings higher flexibility and therefore lower bias. However, this comes at a cost of higher variance: there is higher sample variability when estimating complex models. Overfitting can be a problem.

<https://powcoder.com>

- Decreasing model complexity leads to lower variance.

However, simpler models may not be sufficiently flexible to capture the underlying patterns in the data, leading to higher bias. Underfitting can be a problem.

Add WeChat powcoder

The bias-variance trade-off

Assignment Project Exam Help

To review, we use the training data to estimate the additive error model

<https://powcoder.com>

$$Y = f(X) + \varepsilon$$

leading to an estimator $\hat{f}(x_0)$ for the regression function at given input point x_0 , $f(x_0)$.

Add WeChat powcoder

Assignment Project Exam Help

Linear regression. Adding predictors increases model complexity. Least squares estimates have high variability when the number of inputs is large, but excluding relevant predictors leads to bias.

Review: Lecture 2

<https://powcoder.com>

Add WeChat powcoder

Assignment Project Exam Help

KNN regression. Reducing the number of neighbours increases model complexity. Closer neighbours means lower bias. However, averaging fewer observations increases variance.

Review: Lecture 3
Add WeChat powcoder

Approaches to model selection

Assignment Project Exam Help

Train, validation and test sets. Validation set is used for selecting the optimal model complexity.

Resampling methods. Estimate generalisation performance by generating multiple splits of the training data.

Analytical criteria. Use analytical results to penalise training performance to account for overfitting.

Training, validation, and test split

Assignment Project Exam Help

In the validation set approach, we randomly split the training data into a training set, a validation set and a test set. We select the model with the best predictive performance in the validation set.

<https://powcoder.com>



Add WeChat powcoder

Typically, we use 50-80% of the data for the training set.

Assignment Project Exam Help

1. Estimate different models on the training data.

2. Predict the observations in the validation set.

3. <https://powcoder.com>
Select the model with best validation set performance.

4. **Re-estimate** the selected model by combining the training and validation sets.

5. Predict the test data with the selected model.

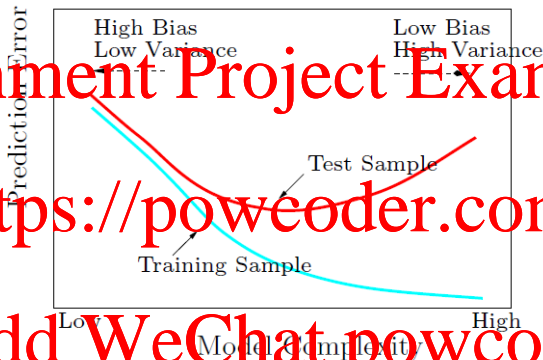
Add WeChat powcoder

Assignment Project Exam Help

The validation set approach has serious limitations when the size of the training data is not large. The model may not have enough data to train on, and there may not be enough cases in the validation set to reliably estimate generalisation performance.

We turn instead to resampling methods.

Learning curve



Assignment Project Exam Help

<https://powecoder.com>

Add WeChat powecoder

We cannot use training set to select the best model. The model minimize the loss of this training data set, not necessarily minimize the loss of the new date sets.

Diagnosing learning curve

Assignment Project Exam Help

Suppose your training loss is low, while validation/test loss is high.

- Underfitting or overfitting problem?

<https://powcoder.com>

Suppose your training loss is high, while validation/test loss is also high.

- Underfitting or overfitting problem?

Add WeChat powcoder

Assignment Project Exam Help

Underfitting: training error is high, validation error is **slightly** > training error.

<https://powcoder.com>

Overfitting: training error is low, validation error is **significantly** > training error.

Add WeChat powcoder

Assignment Project Exam Help

~~The bias-variance decomposition~~
<https://powcoder.com>

Add WeChat powcoder

Expected prediction error

- Consider again the additive error model

$$Y = f(X) + \varepsilon,$$

where we assume that $E(\varepsilon) = 0$ and $\text{Var}(\varepsilon) = \sigma^2$

- In the previous section, we treated $\hat{f}(\cdot)$ as given since our objective was to estimate the test error. Now, we discuss the fundamental problem of choosing a method to learn a predictive function $\hat{f}(\cdot)$.

Expected prediction error

We define the **expected prediction error** for a new input point

$X \doteq x_0$ as

$$\text{Err}(x_0) = E \left[\left(Y_0 - \hat{f}(x_0) \right)^2 \mid X = x_0 \right],$$

where $\hat{f}(x_0) = \hat{f}(x_0) + \epsilon_0$. The expectation is over ϵ_0 and the training sample, i.e. over the sampling distribution of $\hat{f}(\cdot)$. The EPE is a expected loss.

Add WeChat powcoder

Note that this is different from the generalisation error, where $\hat{f}(x_0)$ is an estimate (not an estimator), and the expectation is over the population $P(X, Y)$.

Expected prediction error decomposition

We can write the expected prediction error as:

Assignment Project Exam Help

<https://powcoder.com>

$$= E \left[\left(f(\mathbf{x}_0) + \varepsilon - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right]$$

Add WeChat powcoder

$$= \sigma^2 + E \left[\left(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right]$$

= Irreducible error + Reducible error

Expected prediction error decomposition

$$\begin{aligned}\text{Err}(\mathbf{x}_0) &= \sigma^2 + E_{\mathcal{D}} \left[(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2 | X = \mathbf{x}_0 \right] \\ &= \text{Irreducible error} + \text{Reducible error}\end{aligned}$$

- The first term is the variance of the response around its true mean $f(\mathbf{x}_0)$. We cannot avoid this source of error, and it puts an upper bound on the accuracy of the prediction.

- In choosing a method, our concern is the reducible error: we want to minimise the estimation error

$$E_{\mathcal{D}} \left[(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2 | X = \mathbf{x}_0 \right].$$

- Here, $E_{\mathcal{D}}(\cdot)$ is used to emphasise that the expectation is over the training data. This notation might be omitted later for simplicity purpose.

The bias-variance trade-off

Assignment Project Exam Help

We can show that (see last optional section of this slides):

<https://powcoder.com>

$$E(f(x_0) - \hat{f}(x_0))^2 = [E(\hat{f}(x_0)) - f(x_0)]^2 + E([\hat{f}(x_0) - E(\hat{f}(x_0))]^2)$$
$$= \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0))$$

Add WeChat powcoder

$$= \text{Bias}^2 + \text{Variance}$$

The bias-variance trade-off

$$E(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2 = \text{Bias}^2(\hat{f}(\mathbf{x}_0)) + \text{Var}(\hat{f}(\mathbf{x}_0))$$

Assignment Project Exam Help

- We would like our model to be flexible enough to be able to approximate (possibly) complex relationships between Y and X .

<https://powcoder.com>

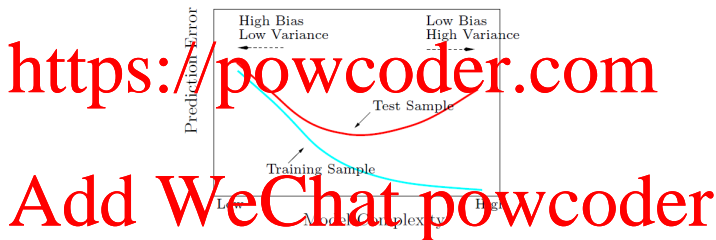
- Typically, the more complex we make the model, the better its approximation capabilities, which translates into lower bias.
- On the other hand, increasing model complexity leads to higher variance. This is due to the larger (effective) number of parameters to estimate.
- Hence, we would like to find the optimal (problem specific) model complexity that minimises our expected loss.

Add WeChat powcoder

The bias-variance trade-off

Increasing model complexity will always reduce the training error, but there is an optimal level of complexity that minimises the test error.

Assignment Project Exam Help



How the plots will be looked like if we fix the model complexity and change x-axis from "model complexity" to "size of the data"?

Assignment Project Exam Help

Just because a model is more “realistic”, it does not mean that it will have higher predictive accuracy. All models are approximations, and our task is to find the most accurate one for our purposes in a data-driven way.

<https://powcoder.com>
Add WeChat powcoder

Assignment Project Exam Help
https://powcoder.com

- Model selection is a set of methods (such as cross validation) that allow us to choose the right model among options of different complexity. It will be a fundamental part of our methodology.
- Similarly to model evaluation, model selection methods are concerned with estimating the generalisation error. However, it is important not to confuse these two steps, which have different goals.

Assignment Project Exam Help

~~KNN bias-variance decomposition~~
<https://powcoder.com>

Add WeChat powcoder

Bias-variance decomposition

Assignment Project Exam Help

Remember the following expression for the expected prediction error:

$$\text{Err}(x_0) = \mathbb{E}[(Y_0 - \hat{f}(x_0))^2 | \mathcal{X} = x_0]$$

$$= \sigma^2 + \text{Bias}^2(\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0))$$

Add WeChat powcoder

The Bias-Variance Decomposition: kNN example

- For kNN regression: Suppose \mathbf{x}_0 is a test data point, its k nearest neighbours are $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_k\}$ with responses $\{y_1, y_2, \dots, y_k\}$. Then the model prediction is

$$\hat{f}(\mathbf{x}_0) = \frac{1}{k} \sum_{\mathbf{x}_i \in N_k(\mathbf{x}_0, D)} y_i,$$

- Then the test error

$$\begin{aligned} E_{\mathbf{x}_0}(\sigma^2) &= \text{Bias}^2(\hat{f}(\mathbf{x}_0)) + \text{Var}(\hat{f}(\mathbf{x}_0)) \\ &= \sigma^2 + \left[E(\hat{f}(\mathbf{x}_0)) - f(\mathbf{x}_0) \right]^2 + E([\hat{f}(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0))]^2) \end{aligned}$$

The Bias-Variance Decomposition: kNN example

- First we check, noting $y_l = f(\mathbf{x}_l) + \varepsilon_l$

$$E(\hat{f}(\mathbf{x}_0)) = E\left[\frac{1}{k} \sum_{l=1}^k y_l\right] = \frac{1}{k} \sum_{l=1}^k E[f(\mathbf{x}_l) + \varepsilon_l]$$

$$= \frac{1}{k} \sum_{l=1}^k f(\mathbf{x}_l)$$

- For the third term, we have

$$\hat{f}(\mathbf{x}_0) - E(\hat{f}(\mathbf{x}_0)) = \frac{1}{k} \sum_{l=1}^k y_l - \frac{1}{k} \sum_{l=1}^k f(\mathbf{x}_l) = \frac{1}{k} \sum_{l=1}^k \varepsilon_l$$

hence, noting the independence of ε_l

$$\text{Var}(\hat{f}(\mathbf{x}_0)) = E\left[\left(\frac{1}{k} \sum_{l=1}^k \varepsilon_l\right)^2\right] = \frac{1}{k^2}(\sigma^2 + \dots + \sigma^2) = \frac{1}{k}\sigma^2$$

The Bias-Variance Decomposition: kNN example

Assignment Project Exam Help

- Finally we have

$$\text{Err}(\mathbf{x}_0) = E[(Y - \hat{f}(\mathbf{x}_0))^2 | X = \mathbf{x}_0]$$

$$= \sigma^2 + \left[f(\mathbf{x}_0) - \frac{1}{k} \sum_{\ell=1}^k f(\mathbf{x}_\ell) \right]^2 + \frac{1}{k}$$

Add WeChat powcoder

where we need the true model values $f(\mathbf{x}_0)$ and $f(\mathbf{x}_\ell)$ (on k neighbours of \mathbf{x}_0).

Assignment Project Exam Help

$$\text{Err}(x_0) = \sigma^2 + \left[f(x_0) - \frac{1}{k} \sum_{\ell=0}^k f(x_\ell) \right]^2 + \frac{\sigma^2}{k}$$

- The model complexity decreases with the number of neighbours k .
- With a small k , the bias will be relatively small since the regression function evaluated at the neighbours $f(x_\ell)$ will be close to $f(x_0)$. However, a small k means that we are averaging only a few observations, leading to high variance ($\frac{\sigma^2}{k}$ is high).
- As we increase k we reduce the variance, at the cost of higher bias.

<https://powcoder.com>
Add WeChat powcoder

Assignment Project Exam Help

~~Cross validation~~
<https://powcoder.com>

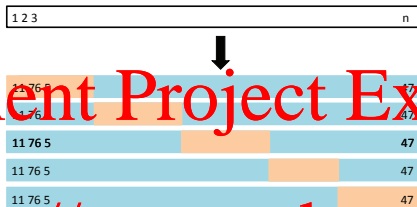
Add WeChat powcoder

Assignment Project Exam Help

Cross validation methods are based on multiple random training/validation set splits. Unlike in the validation set approach, each observation gets a turn at being predicted.

<https://powcoder.com>
Add WeChat powcoder

K-fold cross-validation (key concept)



Assignment Project Exam Help

<https://powcoder.com>

The idea of K-fold cross validation is simple:

1. We randomly split the training sample into K **folds** of roughly equal size.
2. For each fold $k \in \{1, \dots, K\}$, we estimate the model on all other folds combined, and use k as the validation set.
3. The cross validation error is the average error across the K validation sets.

K-fold cross-validation

Assignment Project Exam Help

5-fold and 10-fold CV. $K = 5$ or $K = 10$ folds are common choices for cross validation.

<https://powcoder.com>

Leave one out cross validation (key concept). If we set $K = N$, this is called leave one out cross validation, or **LOOCV**. For each observation i , we train the model on all other observations, and predict i .

Add WeChat powcoder

Leave one out CV

Algorithm Leave one out CV for regression

- 1: **for** $i=1:N$ **do**
- 2: Assign observation i to the validation set
- 3: Assign observations $1, \dots, i-1, i+1, \dots, N$ to the training set \mathcal{D}_{-i} .
- 4: Estimate the model using the training set \mathcal{D}_{-i}
- 5: Compute the prediction $\hat{f}^{-i}(\mathbf{x}_i)$.
- 6: Compute the squared error $(y_i - \hat{f}^{-i}(\mathbf{x}_i))^2$.
- 7: **end for**
- 8: Compute the leave-one-out MSE:

$$\text{MSE}_{\text{cv}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}^{-i}(\mathbf{x}_i))^2$$

LOOCV vs K-fold cross-validation

Assignment Project Exam Help

LOOCV Approximately unbiased estimator of the expected prediction error. However, it can have high variance in some settings (since the training sets are very similar for every prediction) and a high computational cost (except in special cases).

<https://powcoder.com>

K-fold Lower computational cost and may have lower variance. However, it is subject to bias since the training sets are smaller than N .

Add WeChat powcoder

Cross validation: recommendations

One standard deviation rule. Pick the simplest model within one standard deviation of the model with the lowest cross validated errors.

Choice of K . There are no general guidelines for choosing K since the trade-off between variance, bias, and computational cost is highly context specific. The variance of LOOCV tends to be relatively low with stable estimators such as linear regression.

Many predictors. When there are many predictors, pre-screening based on the entire training set may result in misleading CV.

Leave one out CV for linear regression

For a linear regression estimated by OLS, we can use a shortcut to compute the leave-one-out errors without having to reestimate the model. Given the OLS fitted values

$$\hat{\mathbf{y}} = \mathbf{X}(\mathbf{X}^T\mathbf{X})^{-1}\mathbf{X}^T\mathbf{y} = \mathbf{H}\mathbf{y},$$

we can show that the leave-one-out MSE is

$$\text{MSE}_{\text{CV}} = \frac{1}{N} \sum_{i=1}^N (y_i - \hat{f}^-(x_i))^2 = \frac{1}{N} \sum_{i=1}^N \left[\frac{y_i - \hat{f}(x_i)}{1 - H_{ii}} \right]^2,$$

where H_{ii} is the i th diagonal element of the hat matrix \mathbf{H} .

Generalised cross validation

The previous method applies to many situations in which we have fitted values of the type

$$\hat{y} = S\mathbf{y}$$

The **generalised cross validation** method approximates the leave one out MSE as

$$\text{MSE}_{\text{GCV}} = \frac{1}{N} \sum_{i=1}^N \left[\frac{y_i - \hat{f}(\mathbf{x}_i)}{1 - \text{tr}(\mathbf{S})/N} \right]^2,$$

where $\text{tr}(\mathbf{S})$ is the trace of \mathbf{S} (the sum of the elements in its diagonal). GCV can be computationally convenient in some settings.

Assignment Project Exam Help

The bias-variance derivation details
(optional) <https://powcoder.com>

Add WeChat powcoder

The Bias-Variance Decomposition

Assignment Project Exam Help

- We first focus on the squared error loss function. Assume the true model $Y = f(X) + \epsilon$ where $E(\epsilon) = 0$ and $\text{Var}(\epsilon) = \sigma^2$.
- For a test point \mathbf{x}_0 , first we look at

$$(Y - \hat{f}(\mathbf{x}_0))^2$$

$$= \left[Y - f(\mathbf{x}_0) + f(\mathbf{x}_0) - E_{\mathcal{D}} \hat{f}(\mathbf{x}_0) + E_{\mathcal{D}} \hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0) \right]^2$$

$$\begin{aligned} &= (Y - f(\mathbf{x}_0))^2 + (f(\mathbf{x}_0) - E_{\mathcal{D}} \hat{f}(\mathbf{x}_0))^2 + (E_{\mathcal{D}} \hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2 \\ &\quad + 2(Y - f(\mathbf{x}_0))(f(\mathbf{x}_0) - E_{\mathcal{D}} \hat{f}(\mathbf{x}_0)) + 2(Y - f(\mathbf{x}_0))(E_{\mathcal{D}} \hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0)) \\ &\quad + 2(f(\mathbf{x}_0) - E_{\mathcal{D}} \hat{f}(\mathbf{x}_0))(E_{\mathcal{D}} \hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0)) \end{aligned}$$

Add WeChat powcoder

The Bias-Variance Decomposition

Assignment Project Exam Help

- We calculate in test error as, noting that $Y - \hat{f}(\mathbf{x}_0) = \epsilon$,

$\text{Err}(\mathbf{x}_0)$

$$= E_{Y, \mathcal{D}}[(Y - \hat{f}(\mathbf{x}_0))^2 | Y = \mathbf{x}_0]$$
$$= E_{Y, \mathcal{D}}[\epsilon^2] + E_{Y, \mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}f(\mathbf{x}_0))^2]$$

$$+ E_{Y, \mathcal{D}}[(E_{\mathcal{D}}\hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2]$$

$$+ 2E_{Y, \mathcal{D}}[(\epsilon)(f(\mathbf{x}_0) - E_{\mathcal{D}}f(\mathbf{x}_0))] + 2E_{Y, \mathcal{D}}[(\epsilon)(E_{\mathcal{D}}\hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))]$$
$$+ 2E_{Y, \mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))(E_{\mathcal{D}}\hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))]$$

The Bias-Variance Decomposition

Assignment Project Exam Help

- Clearly

$$E_{Y,\mathcal{D}}[\epsilon^2] = E_{\epsilon}[\epsilon^2] = \text{Var}(\epsilon) = \sigma^2$$

- There is no randomness in $(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))^2$, so

$$E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))^2] = (f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))^2$$

Add WeChat powcoder
called the **Squared Bias**

The Bias-Variance Decomposition

- And the **Variance**

Assignment Project Exam Help

- There is no randomness in $(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))$, hence

$$E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))\hat{f}(\mathbf{x}_0)] = E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))]E_{Y,\mathcal{D}}[\hat{f}(\mathbf{x}_0)] = 0$$

- With independence, we have

$$E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))\hat{f}(\mathbf{x}_0)] = E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))]E_{Y,\mathcal{D}}[\hat{f}(\mathbf{x}_0)] = 0$$

- Similarly

$$E_{Y,\mathcal{D}}[(f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))(E_{\mathcal{D}}\hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))] = 0$$

The Bias-Variance Decomposition

- Then we have the following decomposition

$$\text{Err}(\mathbf{x}_0) = \sigma_\epsilon^2 + \text{Bias}^2(\hat{f}(\mathbf{x}_0)) + \text{Var}(\hat{f}(\mathbf{x}_0))$$

- The first term is irreducible error. This exists due to the nature. We cannot make it smaller through any modelling
- the second term is the squared bias, which is defined by

$$\text{Bias}^2(\hat{f}(\mathbf{x}_0)) = (f(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}(\mathbf{x}_0))^2$$

- And the last term is the variance (the expected squared deviation of the estimated $\hat{f}(\mathbf{x}_0)$ around its mean, i.e.,

$$\text{Var}(\hat{f}(\mathbf{x}_0)) = E_{\mathcal{D}}[(E_{\mathcal{D}}\hat{f}(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0))^2]$$

where $E_{\mathcal{D}}$ is the “average” over all the training data.

The Bias-Variance Decomposition: OLS regression example

- For OLS with the estimated model $\hat{f}_{ols}(\mathbf{x}_0) = \hat{\beta}^T \mathbf{x}_0 = \mathbf{x}_0^T \hat{\beta}$,

we have

$$E_{\mathcal{D}} \hat{f}_{ols}(\mathbf{x}_0) = E_{\mathcal{D}} [\mathbf{x}_0^T \hat{\beta}] = \mathbf{x}_0^T E_{\mathcal{D}} [(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}]$$

$$\begin{aligned} &= \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E_{\mathcal{D}} [\mathbf{y}] \\ &= \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T E_{\mathcal{D}} [\mathbf{f} + \epsilon] \end{aligned}$$

$$= \mathbf{x}_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{f}. \quad (\text{Note that } E(\epsilon) = 0)$$

Add WeChat powcoder

where $\mathbf{f} = [f(\mathbf{x}_1), f(\mathbf{x}_2), \dots, f(\mathbf{x}_N)]^T$ the vector of true model values at training data and $\epsilon = [\epsilon_1, \epsilon_2, \dots, \epsilon_N]^T$ given by $y_i = f(\mathbf{x}_i) + \epsilon_i$ ($i = 1, 2, \dots, N$).

Assignment Project Exam Help

- Hence the Bias square is

$$\text{Bias}^2 = (f(x_0) - x_0^T (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{f})^2$$

- Denote

$$\mathbf{h}(x_0) = \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{x}_0 = [h_1(x_0), h_2(x_0), \dots, h_N(x_0)]^T$$

The Bias-Variance Decomposition: OLS regression example

- For the variance term, we have

$$\begin{aligned}\text{Var}(\hat{f}(\mathbf{x}_0)) &= E_{\mathcal{D}}[(\hat{f}_{ols}(\mathbf{x}_0) - E_{\mathcal{D}}\hat{f}_{ols}(\mathbf{x}_0))^2] \\ &= E_{\mathcal{D}}\left[\left(\mathbf{h}(\mathbf{x}_0)^T \mathbf{y} - \mathbf{h}(\mathbf{x}_0)^T \mathbf{f}\right)^2\right]\end{aligned}$$

$$\begin{aligned}&= E_{\mathcal{D}}\left[\left(\mathbf{h}(\mathbf{x}_0)^T \boldsymbol{\epsilon}\right)^2\right] \\ &= E_{\mathcal{D}}\left[h_1(\mathbf{x}_0)\epsilon_1 + h_2(\mathbf{x}_0)\epsilon_2 + \dots + h_N(\mathbf{x}_0)\epsilon_N\right]^2\end{aligned}$$

- According to assumptions, $E_{\mathcal{D}}(\epsilon^2) = \sigma^2$ and all ϵ_i 's are independent of each other. We have

$$\text{Var}(\hat{f}(\mathbf{x}_0)) = h_1(\mathbf{x}_0)\sigma^2 + h_2(\mathbf{x}_0)\sigma^2 + \dots + h_N(\mathbf{x}_0)\sigma^2 = \|\mathbf{h}(\mathbf{x}_0)\|^2 \sigma^2$$

- Finally we have

$$\text{Err}(\mathbf{x}_0) = \sigma^2 + (f(\mathbf{x}_0) - \mathbf{h}(\mathbf{x}_0)^T \mathbf{f})^2 + \|\mathbf{h}(\mathbf{x}_0)\|^2 \sigma^2$$

Assignment Project Exam Help

- How does model selection relate to the bias-variance trade-off?

- <https://powcoder.com>
What is a validation set? How is it different from a test set?

- What is K-Fold cross validation? Describe how it works.

- Add WeChat powcoder
What is the one standard deviation rule?