

Knowledge Check 1.1

1. For each of parts (a) through (d), indicate whether we would generally expect the performance of a flexible statistical learning method to be better or worse than an inflexible method. Justify your answer.

(a) The sample size n is very large, and the number of predictors p is small.

Ans: When the sample size is large, the simple method can better capture the complex relationship between the predictors and the response variable. The larger the sample size, the more data points available to accurately estimate relationships. With fewer predictors, the simpler method can accommodate more complex models that closely follow the data, resulting in better performance in terms of prediction accuracy.

b) The number of predictors p is very large, and the number of observations n is small.

Ans: When the sample size is large, when the number of predicted values is very large, the risk of overfitting increases. Simplified methods

have great potential to capture complex relationships, but with too few observations, they can introduce noise or spurious patterns into the data, resulting in poor generalization performance.

c) The relationship between the predictors and the response is highly non-linear.

Ans: When the relationship between predictors and responses is highly nonlinear, we generally expect the performance of a simple statistical learning method to be better than of a non-standardized simpler method. Such as non-linear regression or tree-based models, can capture non-linear patterns well in tandem. They can be adjusted.

d) The variance of the error terms, i.e. $\sigma^2 = \text{Var}(\epsilon)$, is extremely high.

Ans: Higher error variances indicate greater noise or randomness in the data. Perhaps for simple methods, which have more power to accommodate complex images, this noise may be overqualified, and the generalization may be poor as a consequence. In contrast, unbiased methods that impose simple models can be more robust to higher error margins and provide better performance by focusing on key features or patterns in the data.

2. Explain whether each scenario is a classification or regression problem, and indicate whether we are most interested in inference or prediction. Finally, provide n and p .

a) We collect a set of data on the top 500 firms in the US. For each firm we record profit, number of employees, industry and the CEO salary. We are interested in understanding which factors affect CEO salary.

Ans: This situation is a regression trouble. We are interested in understanding the factors that affect CEO revenue, that's a continuous variable. The intention is to expect the CEO income primarily based on the recorded elements which includes earnings, range of employees, and industry. In this situation, we are typically inquisitive about inference, this is, knowledge the relationship between the predictors (income, range of personnel, enterprise) and the response variable (CEO salary). The values of n and p are 500 and 4.

b) We are considering launching a new product and wish to know whether it will be a success or a failure. We collect data on 20 similar products that were previously launched. For each product we have recorded whether it was a success or failure, price charged for the product, marketing budget, competition price, and ten other variables.

Ans: This scenario is a categorical classification problem. We need to determine whether or not the brand new product might be a achievement or a failure based at the collected statistics on similar previously released products. The reaction variable in this case is binary (fulfillment or failure). The predictors include variables which includes rate charged for the product, marketing price range, competition price and other variables. The intention is to are expecting the fulfillment or failure of the new product. In this state of affairs, we are ordinarily inquisitive about prediction. The quantity of observations, n , is given as 20, and the number of predictors p is 13.

c) We are interested in predicting the % change in the US Dollar in relation to the weekly changes in the world stock markets. Hence we collect weekly ~~as~~ data for all of 2012. For each week we record the % change in the USD/Euro, the % change in the US market, the % change in the British market, and the % change in the German market.

Sols: This Scenario is a ~~as~~ Regression Trouble.
We propose to are expecting the proportion trade with the US Dollar in relation to the weekly changes in the international stock markets. The response variable is the % exchange in the USD/Euro, and the predictors are the % changes in the US marketplace, the British market, and the German market place. The purpose is to expect the proportion change inside the USD/Euro based totally at the weekly changes in the other markets. In this situation, we are normally inquisitive about prediction. The statistics is collected weekly for the whole year of 2012, so n would be 52 weeks of 2012, and p is same to 3 (the quantity of predictors)

3. (a) Describe three real-life applications in which regression might be useful. Describe the response as well as the predictors. Is the goal of each application inference or prediction? Explain your answer.

1. Real Estate Price Prediction:

* Response Variable : Price of a residence or asset.

* Predictors : Features of the assets together with variety of rooms, vicinity, facilities etc.

* Goal : Prediction. The goal is to predict the price of a property based totally on its characteristics. This estimate is precious for buyers, dealers and real property marketers.

2. Customer Life time Value Estimation:

* Response Variable : Lifetime fee of a consumer. Overall sales generated by means of a purchaser over their lifetime with a corporation.

* Predictors : Customer demographics, purchasing behavior, historical facts on patron transactions, and many others.

* Goal : Prediction . The objective is to estimate the potential value of a customer to a enterprise, helping with patron segmentation, personalized marketing techniques, and resource allocation.

3. Medical Treatment Outcome Analysis:

* Response Variable : Health final results or remedy effectiveness (e.g. Survival time, improvement in signs, discontent in ache).

* Predictors : Patient traits (age, gender, medical history), remedy variables (dosage, period), and different relevant elements.

* Goal : Inference . The purpose is to apprehend the relationship among predictors and remedy effects, identifying factors that appreciably have an impact on the effectiveness of medical treatments.

(b) Describe three real-life applications in which classification might be useful. Describe the response, as well as the predictors. Is the goal of each application inference or prediction? Explain your answer.

Ans:

1. Email Spam Detection:

* Response Variable: Email class (Junk mail or non-spam)

* Predictors: Email content, sender facts, metadata, and so on.

* Goal: Prediction. The objective is to as it should be classify incoming emails as either spam or legitimate (non-junk mail) to filter undesirable or malicious messages.

2. Credit Risk Assessment:

* Response Variable: Creditworthiness

* Predictors: Financial facts (income, credit score records, debt-to-income ratio, etc), non-public statistics (age, employment reputation), and other relevant facts.

* Goal : Prediction. The purpose is to evaluate the credit danger associated with an character or a agency to to make informed lending choices and control capacity economic dangers.

3. Disease Diagnosis:

* Response Variable : Disease prognosis (presence or absence of a selected sickness).

* Predictors : Medical check results, affected person signs, scientific records, and so on.

* Goal : Prediction. The objective is to classify sufferers into exclusive sickness categories based totally on their symptoms and test consequences, helping in early detection, treatment planning, and patient

Control -

source V

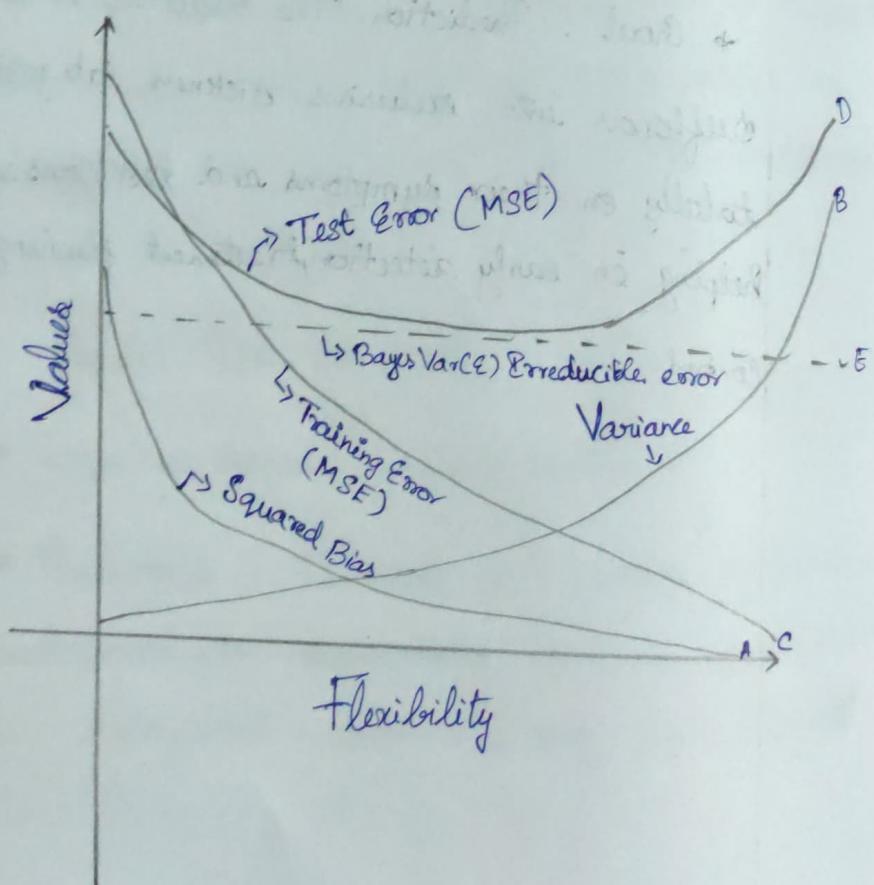
and, before
(DEM)

and, during L

without

Knowledge Check 1.2

1. We now revisit the bias-variance decomposition.
- A) Provide a sketch of typical (squared bias, variance, training error, test error, and Bayes for irreducible) error curves, on a single plot, as we go from less flexible statistical learning methods towards more flexible approaches. The x-axis should represent the amount of flexibility in the method, and the y-axis should represent the values of for each curve. There should be five curves. Make sure to label each one.



B) Explain why each of the five curves has the shape.

1) Squared Bias Curve (Curve A):

The Squared Bias Curve initially starts high and decreases as the flexibility of the statistical learning method increases. This is because less flexible models have simpler assumptions and limited capacity to capture complex patterns in the data. As flexibility increases, the model becomes more capable of fitting the underlying relationships, resulting in a decrease in bias. However, if the model becomes too flexible, it may start overfitting the noise in the data, leading to an increase in bias again.

2) Variance Curve (Curve B):

The Variance curve initially starts low and increases as the flexibility of the statistical learning method increases. Less flexible models have lower variance because they have fewer parameters or assumptions, resulting in more stable predictions. As flexibility increases, the models become more complex and have higher variance because they can fit the training data more closely. However, this increased flexibility can also lead to overfitting, causing the variance to increase.

3) Training Error Curve (Curve C):

The training error curve starts high with low flexibility and decreases as flexibility increases. Less flexible models have limited expressive power and may not capture the true underlying patterns in the data. As flexibility increases, the models can better fit the training data, leading to a ~~decrease~~ decrease in training error. In extreme cases, very flexible models can even achieve zero training error by perfectly fitting the training data.

4) Test Error Curve (Curve D):

The test error curve initially starts high with low flexibility, decreases and may reach a minimum point before increasing again. Initially, with low flexibility, the model may underfit the data, resulting in high test error. As flexibility increases, the model captures more complex patterns, leading to a decrease in test error. However, after reaching the minimum point, further increasing flexibility may cause the model to overfit the training data, resulting in an increase in test error.

This occurs because the model starts to capture noise or idiosyncrasies in the training data that do not generalize well to unseen data.

5) Bayes (Irreducible) Error Curve (Curve E):

The Bayes error curve represents the irreducible error curve represents the irreducible error or the inherent noise in the data that cannot be reduced regardless of the flexibility of the model. It is typically a flat horizontal line since it is independent of the model's flexibility. This error arises due to unobservable or unpredictable factors in the data or limitations in the measurement process. The Bayes error provides a theoretical lower limit on the achievable error, and no model can surpass it.

2. What are the advantages and disadvantages of a very flexible (versus a less flexible) approach for regression or classification? Under what circumstances might a very flexible approach be preferred to a less flexible approach? When might a less flexible approach be preferred?

Ans:

Advantages of a Very Flexible Approach:

1. Capture Complex Relationships: Highly flexible models have the ability to capture intricate and non-linear relationships between predictors and the response variable.
2. Improved Predictive Performance:
3. Adaptability to Diverse Data

DisAdvantages of a Very Flexible Approach:

1. Overfitting
2. Interpretability and Transparency
3. Computational demands:

More Flexible approach:

1. Complex Relationships:

When there is a prior belief or evidence that the relationship between predictors and the response variable is complex and non-linear, a more flexible model can better capture these intricate patterns.

2. Large and Diverse Data:

With large datasets in a high-dimensional feature space, a more flexible approach can effectively utilize the abundance of data to model complex relationships and potentially achieve superior performance.

3) Emphasis on Predictive Accuracy:

When the primary goal is to maximize predictive accuracy and there is less concern about model interpretability, a more flexible model can be chosen to optimize performance.

Less Flexible Approach:

1. Limited Data:

When the sample size is small relative to the number of predictors, there is a higher risk of overfitting with more flexible models. In such cases, simpler models with fewer parameters can be more robust and less prone to overfitting.

2. Interpretability and Explainability:

In some scenarios, the interpretability and explainability of the model are crucial. Less flexible models, such as linear regression or decision trees, provide transparent insights into the relationship between predictors and the response variable, making them more suitable when interpretability is a priority.

3. Computational Constraints:

When computational resources are limited, training and deploying a complex, highly flexible model may not be feasible. Less flexible models are computationally less demanding and can be trained and deployed more efficiently.

3. Describe the differences between a parametric and non-parametric statistical learning approach. What are the advantages of a parametric approach to regression or classification (as opposed to a non-parametric approach)? What are its disadvantages?

Parametric Approach:

- * Parametric models make assumptions about the functional form or distribution of the relationship between predictors and the response variable.
- * They have a fixed number of parameters that are estimated from the data.
- * Examples of parametric models include linear regression, logistic regression, and linear discriminant analysis.

Non-Parametric Approach:

- * Non-parametric models do not make explicit assumptions about the functional form or distribution of the relationship between predictors and the response variable.
- * Non-parametric models adapt to the data and have a flexible number of parameters that depend on the complexity of the data.
- * Examples of non-parametric models include decision trees, random forests, support vector machines, and k-nearest neighbors.

4. The table below provides a training data set containing six observations, three predictors, and one qualitative response variable.

	X_1	X_2	X_3	Y
1	0	3	0	R
2	2	0	0	R
3	0	1	3	R
4	0	1	2	G
5	-1	0	1	G
6	1	1	1	R

Suppose we wish to use this data set to make a prediction for Y when $X_1 = X_2 = X_3 = 0$ using k-nearest neighbors.

- a) Compute the Euclidean distance between each observation and the test point, $(X_1, X_2, X_3) = (0, 0, 0)$

$$\text{Distance} = \sqrt{(X_1 - 0)^2 + (X_2 - 0)^2 + (X_3 - 0)^2}$$

$$\text{Observation } 1 : \text{Distance} = \sqrt{(0-0)^2 + (3-0)^2 + (0-0)^2} \\ = \sqrt{9} = 3$$

$$\text{Observation } 2 : \text{Distance} = \sqrt{(2-0)^2 + (0-0)^2 + (0-0)^2} \\ = \sqrt{4} = 2$$

$$\text{Observation } 3 : \text{Distance} = \sqrt{(0-0)^2 + (1-0)^2 + (3-0)^2} \\ = \sqrt{10} = 3.16$$

$$4 \Rightarrow \text{Distance} : \sqrt{(0-0)^2 + (0-0)^2 + (2-0)^2} = \sqrt{4} = 2.23$$

$$5 \Rightarrow \text{Distance} = \sqrt{(-1)^2 + (0)^2 + (1)^2}$$

$$= \sqrt{2} \approx 1.414$$

$$6 \Rightarrow \text{Distance} = \sqrt{(1)^2 + (1)^2 + (1)^2}$$

$$= \sqrt{3} \approx 1.73$$

b) What is our prediction with $k=1$? Why?

Ans: With $k=1$, we consider only the nearest neighbor to the test point. The nearest neighbor is observation 2 with a distance of 2. Therefore, the prediction would be the value of γ associated with observation 2, which is "RED".

c) What is our prediction with $k=3$? Why?

Ans: With $k=3$, we consider the three nearest neighbors to the test point. The three nearest neighbors are 2, 6, 5. Out of these three, there are two "RED" and one "green" label.

Therefore the prediction with $k=3$ would be "RED" since it is the majority label among the three nearest neighbors.

d) If the Bayes decision boundary in this problem is highly nonlinear, then would we expect the best value for K to be large or small?

Ans:

If the Bayes decision boundary is highly nonlinear, a smaller value of K is generally preferred.