

Statistical Learning and Data Mining

Module 1: Introduction to Statistical Learning

Semester 2, 2017

Discipline of Business Analytics, The University of Sydney Business School

Module 1: Introduction to Statistical Learning

1. Introduction
2. Business examples and data
3. Notation
4. Statistical decision theory
5. Evaluating model performance
6. Key concepts and themes

Introduction

Introduction

Machine learning is a set of methods for automatically detecting patterns in data and using them for predicting future data and guiding decision making. In other words, learning from data.

Introduction

Two trends bring machine learning to the forefront of successful business decision making:

- We are in the era of **big data**. The Internet and increasing presence of data capturing devices (such as mobile phones, cameras, sensors, card readers, etc), combined with large reductions in the cost of storage, brought an unprecedented availability of data, and continued dramatic growth in the size of data sets.
- Advancing computing power (realising **Moore's law**) increases the scope for exploring complex patterns in data.

Types of machine learning

There are two main types of machine learning:

- In **predictive** or **supervised learning**, the objective is to learn a function to predict an output variable Y based on observed input variables x_1, \dots, x_p . We develop methods that learn this function based on labelled data $\{(x_i, y_i)\}_{i=1}^N$, which we call the training data.
- In **descriptive** or **unsupervised learning**, we are given only inputs $\{x_i\}_{i=1}^N$, and the goal is to find “interesting” patterns in this data.

There are other tasks such as semi-supervised learning and reinforcement learning, which we do not explore in this unit.

Supervised learning

In supervised learning, the **output** or **response variable** can be of any type. However, most methods address two main classes of supervised learning problems:

- In **regression**, the response is a quantitative scalar (such as the income of a worker).
- In **classification**, the response is a **nominal** or **categorical** variable $Y \in \{1, \dots, C\}$, where C is the number of classes. When $C = 2$, this is called binary classification; if $C > 2$, this is called multiclass classification.

Example: handwritten digit recognition



A view of the MNIST dataset.

Statistical learning

We can think of **statistical learning** as a framework for machine learning that draws on statistics. There are different reasons for following this approach:

- Generating probabilistic predictions and quantifying uncertainty. We need tools for decision making under uncertainty.
- Studying the statistical properties of learning methods, with the objective of understanding why certain prediction methods work well and informing applications.
- Drawing on insights from statistics to develop learning methods.

Data mining

Data mining is the process of extracting interesting and previously unknown patterns and relations from large databases, drawing on the fields of machine learning, statistics, and database technology.

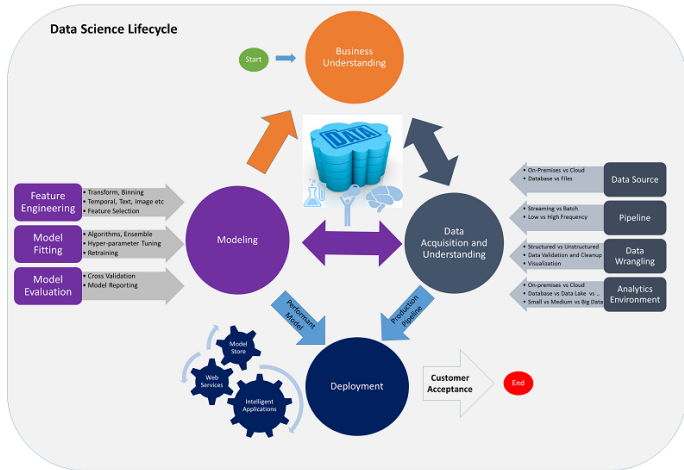
In data mining, the analysis should be both useful and understandable to the data owner.

Data science

Data science is a multidisciplinary field that combines knowledge and skills from statistics, machine learning, software engineering, data visualisation, and domain expertise (in our case, business expertise) to uncover value from large and diverse data sets.

Data scientists often work directly with stakeholders (say, product managers) to translate data analysis results into action. A possible outcome is the creation of a **data product**.

The data science process: a real-world perspective



<https://docs.microsoft.com/en-us/azure/machine-learning/data-science-process-overview>

Data analysis process in this unit

1. Problem formulation.
2. Data collection and preparation.
3. Exploratory data analysis (EDA).
4. Model building, estimation, and selection.
5. Model evaluation.
6. Communicate results.

Unit structure

- Statistical foundations: concepts, statistical thinking, model evaluation, model selection, and model inference.
- Linear methods for regression.
- Linear methods for classification.
- Nonlinear methods for regression and classification.
- Unsupervised learning.

Learning outcomes

By successfully completing this unit, you are expected to:

1. Understand the conceptual and theoretical foundations of statistical learning.
2. Develop an in-depth knowledge of regression, classification, and unsupervised learning methods for business applications.
3. Be able to conduct a complete data analysis project based on these foundations and methods.
4. Know how to use Python for your practical workflow under realistic data complexity (including tasks such as data manipulation and visualisation).
5. Effectively communicate your results to guide decision making.

Comments

- This unit is designed as training for real-world data science, which requires a range of skills.
- Practical statistical learning involves more than knowing the methods in the lectures: professionals in this area typically spend a substantial amount of time on tasks such as data management, exploratory data analysis, feature engineering, and implementing methods.
- All of this generally done through coding. Therefore, Python is your bridge between knowledge and practice.
- For these reasons, please note that this unit requires independent work and a substantial time commitment (within the university guidelines).

Business examples and data

Some examples

- Credit card fraud detection: collect data from multiple sources to learn typical customer behaviour, then use this model to detect suspicious transactions for further investigation.
- Customer risk analysis: instead of denying sales (say, auto loans, credit cards, and insurance policies) to higher risk customers, it is usually a better strategy to price risk accordingly using available data.
- Retail: providing better product recommendations using collaborative filtering techniques (predicting user preferences based on preferences of similar users).
- Advertising: making online ads more relevant to users by predicting click-through rates.

Zillow Kaggle competition

- Kaggle is a crowdsourcing platform that allows organisations to post data prediction problem to be solved by public competition.
- Zillow's Home Value Prediction is a current competition (with a 1.2 million dollar cash prize) that invites participants to make predictions about the future sale prices of homes (a regression problem).
- In this competition, the goal is to improve on Zillow's home valuation estimates ("ZEstimates"), which are based on 7.5 million statistical and machine learning models that analyze hundreds of data points on each property.

Customer relationship management

- **Customer relationship management (CRM)** is a set of practices that involve collecting and studying customer information with the objective of maximising **customer lifetime value (CLV)**, the net value of a customer to a firm over his/her entire lifetime.
- CRM may be part of a customer-centric (as opposed to brand-centric) business strategy, which focuses on customer satisfaction and loyalty towards the acquisition and retention of profitable customers.
- CRM has four main areas: customer acquisition, retention, churn, and win-back. Statistical models and machine learning algorithms play a central role in each of these areas.

Customer relationship management

	Customer	Acquisition	First_Purchase	CLV	Duration	Censor	Acq_Expense	Acq_Expense_SQ	Industry	Revenue
0	1	1	433.64	0.0000	384	0	760.36	578147.33	1	30.16
1	2	0	0.00	0.0000	0	0	147.70	21815.29	1	39.80
2	3	0	0.00	0.0000	0	0	252.56	63786.55	1	54.93
3	4	1	225.84	5.7316	730	1	609.73	371770.67	1	45.83
4	5	1	363.04	0.0000	579	0	672.36	452067.97	1	69.03
5	6	0	0.00	0.0000	0	0	435.57	189721.22	0	22.54
6	7	0	0.00	0.0000	0	0	362.90	131696.41	0	32.97
7	8	0	0.00	0.0000	0	0	883.54	780642.93	0	22.48
8	9	1	599.30	6.9161	730	1	452.35	204620.52	1	17.98
9	10	1	271.18	6.0839	730	1	786.72	618928.36	1	38.91
10	11	0	0.00	0.0000	0	0	504.03	254046.24	1	28.85
11	12	0	0.00	0.0000	0	0	842.50	709806.25	0	49.41
12	13	0	0.00	0.0000	0	0	150.51	22653.26	1	41.91

The data is from Kumar and Petersen (2012), and refers to corporate clients.

Customer relationship management

Kumar and Petersen (2012) estimate a model to predict the response

$$Y = \begin{cases} 1 & \text{if the customer was acquired,} \\ 0 & \text{if the customer was not acquired,} \end{cases}$$

based on predictors such as the dollar spent on marketing efforts to acquire the prospect, and characteristics of the prospect's firm such as industry, revenue, and number employees.

This is a binary classification problem.

Market basket analysis

Market basket analysis is an unsupervised learning task that is of wide interest in commercial data mining.

- The data for this task consists of a very large and sparse binary matrix where each column represents a product and each row represents a transaction. We set $x_{ij} = 1$ if product j was purchased on transaction i , and $x_{ij} = 0$ otherwise.
- The objective of the analysis is to discover consumer preferences and lifestyles: what goes with what? What products are purchased together (say, Vegemite and bread)?
- The results then inform decisions on product placement in stores, co-marketing promotions, etc.

Market basket analysis

Rule No.	Left-Hand Side (Antecedent)		Right-Hand Side (Consequent)	Support	Confidence	Lift
1	{beef, dairy produce}	=>	{vegetables}	0.030	0.607	2.225
2	{poultry}	=>	{vegetables}	0.029	0.575	2.105
3	{dairy produce, fruit, sausage}	=>	{vegetables}	0.027	0.574	2.103
4	{beef}	=>	{vegetables}	0.046	0.560	2.050
5	{dairy produce, vinegar/oils}	=>	{vegetables}	0.031	0.536	1.962
6	{fruit, sausage}	=>	{vegetables}	0.034	0.529	1.938
7	{bread and baked goods, dairy produce, fruit}	=>	{vegetables}	0.041	0.528	1.933
8	{pork}	=>	{vegetables}	0.030	0.522	1.912
9	{cheese, fruit}	=>	{vegetables}	0.027	0.520	1.904
10	{dairy produce, fruit, non-alc. drinks}	=>	{vegetables}	0.033	0.518	1.899

Table from Miller (2015).

Notation

Study tips

- Always start by making sure that you understand the notation and definitions. Focus first on meaning, then connections.
- If there is a learning challenge, is the root of the problem in understanding notation, concepts, reasoning, or algebra?
- When reading an equation, you should be able to identify parameters and constants, distinguish between random variables and observed values, and distinguish between scalars, vectors, and matrices.
- When there is an expectation or variance operator, what distribution is it over? That is, what random variables do they refer to?

Notation

- We use upper case letters such as Y to denote random variables, regardless of dimension.
- Lower case letters denote observed values. For example, y denotes the realised value of the random variable Y .
- We use i to index the observations, j to index the inputs. For example, y_i is the observed response for sample i , while x_{ij} is the value of predictor j for observation i .
- We use the hat notation (e.g. $\hat{\beta}$) for estimators and estimates. The notation may not distinguish between the two (refer to context).
- Vectors are in lower case bold letters. Matrices are in upper case bold letters.

Vector and matrix notation

Response vector:

$$\mathbf{y} = \begin{pmatrix} y_1 \\ y_2 \\ \vdots \\ y_N \end{pmatrix}$$

Vector and matrix notation

Vector of predictor (features, attributes, covariates, regressors, independent variables) values for observation i :

$$\mathbf{x}_i = \begin{pmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{ip} \end{pmatrix}$$

Vector of observed values for predictor j :

$$\mathbf{x}_j = \begin{pmatrix} x_{1j} \\ x_{2j} \\ \vdots \\ x_{Nj} \end{pmatrix}$$

Vector and matrix notation

Design matrix:

$$\mathbf{X} = \begin{pmatrix} x_{11} & x_{12} & \dots & x_{1p} \\ x_{21} & x_{22} & \dots & x_{2p} \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} & x_{N2} & \dots & x_{Np} \end{pmatrix}$$

Statistical decision theory

Prediction

We define prediction as follows:

1. Train a predictive function $\hat{f}(x)$ using data $\mathcal{D} = \{(y_i, \mathbf{x}_i)\}_{i=1}^N$.
2. Upon observing a new input point \mathbf{x}_0 , make the prediction $\hat{f}(\mathbf{x}_0)$, the predictive function evaluated at \mathbf{x}_0 .

How should we perform this learning task? How do we define our objective? How do we measure success in achieving this objective? To answer these questions, we turn to **decision theory**. We mostly focus on regression problems for simplicity.

Loss function

A **loss function** or **cost function** $L(y, f(\mathbf{x}))$ measures the cost of predicting $f(\mathbf{x})$ when the truth is y . The most common loss function for regression is the **squared loss**:

$$L(y, f(\mathbf{x})) = (y - f(\mathbf{x}))^2$$

For binary classification, a typical loss function is the **0-1 loss**:

$$L(y, \hat{y}) = \begin{cases} 1 & \text{if } y \neq \hat{y} \\ 0 & \text{if } y = \hat{y}, \end{cases}$$

where \hat{y} is the prediction.

Expected loss

Let Y and X have a joint probability distribution $P(X, Y)$. The idea of decision theory is that we take the action that minimises our **expected loss** or **risk**:

$$R(f) = E [L(Y, f(X))] ,$$

where the expectation is over $P(X, Y)$. Here, the risk is for a given function $f(\cdot)$.

We can use the **law of iterated expectations** to rewrite the expected loss as

$$R(f) = E \left[E (Y - f(X))^2 | X \right] .$$

Optimal prediction

The optimal action is to choose the prediction function $\delta(\cdot)$ that minimises the expected loss. This is equivalent to minimising the expected loss at every input point \mathbf{x} :

$$\delta(\mathbf{x}) = \operatorname{argmin}_{f(\cdot)} E(L(Y, f(\mathbf{x})) | X = \mathbf{x})$$

The solution for the squared loss (see module notes) is the conditional expectation:

$$\delta(\mathbf{x}) = E(Y | X = \mathbf{x})$$

Concept: under the squared error loss, the optimal prediction of Y at any point $X = \mathbf{x}$ is the conditional mean $E(Y | X = \mathbf{x})$.

Statistical modelling

- Our regression problem reduces to the estimation of the conditional expectation function $E(Y|X = \boldsymbol{x})$. In order to learn this function, we need to introduce assumptions.
- Assumptions lead to statistical models.
- For example, the linear regression model assumes that $E(Y|X = \boldsymbol{x})$ is linear:

$$E(Y|X = \boldsymbol{x}) = \boldsymbol{x}^T \boldsymbol{\beta}$$

Additive error model

The **additive error model** is our basic general model for regression. It assumes that the relationship between Y and X is described as

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

where $f(\cdot)$ is an unknown **regression function**, and ε is a random error with mean zero ($E(\varepsilon) = 0$).

Under this model,

$$E(Y|X = \mathbf{x}) = E(f(\mathbf{x}) + \varepsilon) = f(\mathbf{x}),$$

since $E(\varepsilon) = 0$.

Example: linear regression

In the special case of the linear regression model, we assume that

$$f(X) = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p,$$

leading to the model

$$Y = \beta_0 + \beta_1 X_1 + \dots + \beta_p X_p + \varepsilon,$$

and predictions

$$\hat{f}(\mathbf{x}) = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \dots + \hat{\beta}_p x_p,$$

where $\hat{\beta} = (\hat{\beta}_0, \hat{\beta}_1, \dots, \hat{\beta}_p)$ is the vector of least squares estimates of the model parameters.

Statistical decision theory

Our discussion of statistical decision theory lays the foundation for the rest of our discussion.

- Evaluating model performance: estimating the expected loss of a trained model.
- Choosing a learning method: finding and estimating an appropriate model such that we minimise our expected loss.

Evaluating model performance

Evaluating model performance

Model evaluation consists of estimating the expected loss of a trained model. To incorporate model assessment into our analysis, we split the dataset into two parts.

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

Training and test data

- Because we are interested on the estimating how well a model will predict future data, the test set should be kept in a “vault” and brought in strictly at the end of the analysis. The test set does not lead to model revisions.
- We generally allocate 50-80% of the data to the training sample.
- A higher proportion of training data leads to more accurate model estimation, but higher variance in estimating the expected loss.
- The split of the data into the training and test sets is often random, but sometimes there are reasons to consider alternative schemes.

Evaluating test performance

Suppose that we have test observations $\{(\tilde{y}_i, \tilde{\mathbf{x}}_i)\}_{i=1}^M$ and corresponding predictions $\hat{f}(\tilde{\mathbf{x}}_i)$ for $i = 1, \dots, M$. We evaluate model performance by computing the **empirical risk** for the test set:

$$\hat{R}_{\text{test}} = \frac{1}{M} \sum_{i=1}^M L(\tilde{y}_i, \hat{f}(\tilde{\mathbf{x}}_i))$$

Below, we drop the specific notation for test observations for simplicity.

Mean squared error

The choice of loss function leads to a measure of predictive accuracy. Suppose that we have observations y_i and predictions $\hat{y}_i = \hat{f}(\mathbf{x}_i)$ for an arbitrary sample, $i = 1, \dots, n$. The **mean squared error** is:

$$\text{MSE} = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$$

The test mean squared error is the MSE evaluated for the test set.

Mean squared error

The root mean-squared error and the prediction R^2 are derived from the MSE and you may be a better way to report the test results:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2}$$

$$\text{Prediction } R^2 = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

Mean absolute error

Another common measure of performance is the **mean absolute error** (MAE):

$$\text{MAE} = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$$

- Implicit in the use of the MAE is the absolute error loss function. The absolute error setting is less mathematically tractable, which is one of the reasons why focus on the squared error loss.
- In this case the optimal prediction is the conditional median, not the mean.

Generalisation error

The **test** or **generalisation** error is the expected loss for the model estimated with the training data \mathcal{D} . We define it as

$$\text{Err} = E \left[L \left(Y, \hat{f}(X) \right) | \mathcal{D} \right],$$

where the expectation is over $P(X, Y)$.

Concept: the test MSE estimates the test error (under the squared error loss).

Standard error

As always, you should report a measure of sample uncertainty for every important estimate in your analysis. The test MSE is a sample average, so obtaining a standard error is straightforward. The formula for a general sample is:

$$\text{SE}(\text{MSE}) = \frac{1}{\sqrt{n}} \sqrt{\sum_{i=1}^n \frac{\left((y_i - \hat{f}(\mathbf{x}_i))^2 - \text{MSE} \right)^2}{n-1}}$$

Inference for the test errors is possible, but we do not pursue this here.

Key concepts and themes

Key concepts and themes

- The bias-variance trade-off and model selection.
- Overfitting.
- Parametric vs non-parametric models.
- No-free lunch theorem.
- Accuracy vs interpretability.

Expected prediction error

- Consider again the additive error model:

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

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

- 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 $X = \mathbf{x}_0$ as

$$\text{Err}(\mathbf{x}_0) = E \left[\left(Y_0 - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right],$$

where $Y_0 = f(\mathbf{x}_0) + \varepsilon_0$. 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.

Note that this is different from the generalisation error, where $\hat{f}(\mathbf{x}_0)$ 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:

$$\begin{aligned}\text{Err}(\mathbf{x}_0) &= E \left[\left(Y_0 - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right] \\ &= E \left[\left(f(\mathbf{x}_0) + \varepsilon - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right] \\ &= \sigma^2 + E \left[\left(f(\mathbf{x}_0) - \hat{f}(\mathbf{x}_0) \right)^2 \mid X = \mathbf{x}_0 \right] \\ &= \text{Irreducible error} + \text{Reducible error}\end{aligned}$$

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, I wrote $E_{\mathcal{D}}(\cdot)$ to emphasise that the expectation is over the training data.

The bias-variance trade-off

We can show that (see module notes):

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

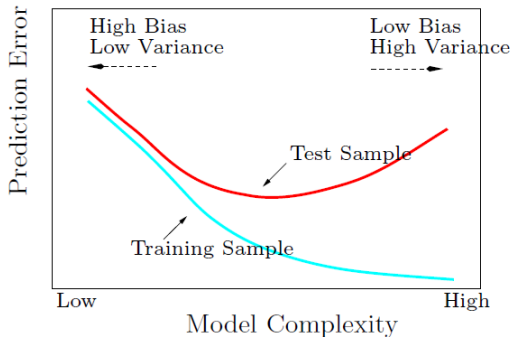
The bias-variance trade-off

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

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

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.



The bias-variance trade-off

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.

Model selection

- **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.
- We conduct model selection on the training data.

Overfitting

- We say that there is **overfitting** when an estimated model is excessively flexible, incorporating minor variations in the training data that are likely to be noise rather than predictive patterns.
- An overfit model has small training errors, but may predict poorly. In essence, it has memorised the training set.
- Not being misled by overfitting is an important reason why we use a test set.

Illustration: predicting fuel economy

- This example uses data extracted from the fuelconomy.gov website run by the US government, which lists different estimates of fuel economy for passenger cars and trucks.
- For each vehicle in the dataset, we have information on various characteristics such as engine displacement and number of cylinders, along with laboratory measurements for the city and highway miles per gallon (MPG) of the car.
- We here consider the unadjusted highway MPG for 2010 cars as the response variable, and a single predictor, engine displacement.

Illustration: predicting fuel economy

A scatter plot reveals a nonlinear association between the two variables. We therefore need a model that is sufficiently flexible to capture this nonlinearity.

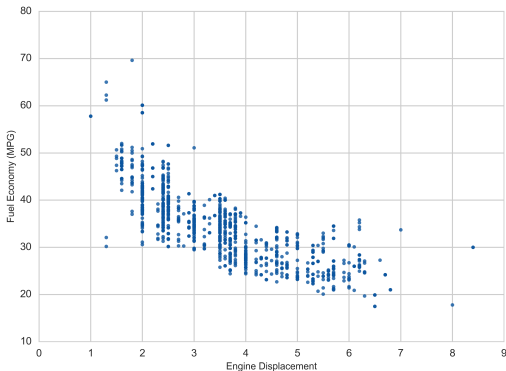
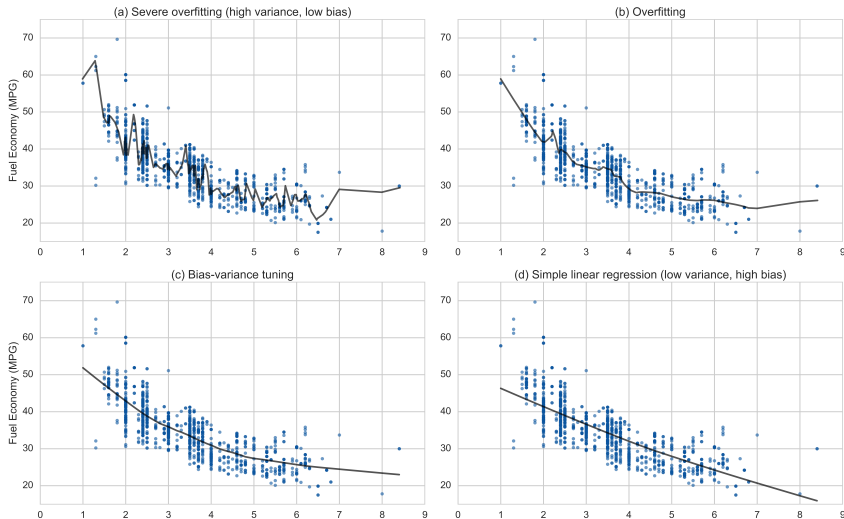


Illustration: predicting fuel economy



Parametric vs nonparametric models

There are many ways to define statistical models, but the most important distinction is the following:

- A **parametric model** has a fixed number of parameters. Parametric models are faster to use, and more interpretable, but have the disadvantage of making stronger assumptions about the data.
- In a **nonparametric model**, the number of parameters grows with the size of the training data. Nonparametric are more flexible, but have larger variance and can be computationally infeasible for large datasets. An example is the K-nearest neighbours method, which we will study in the next module.

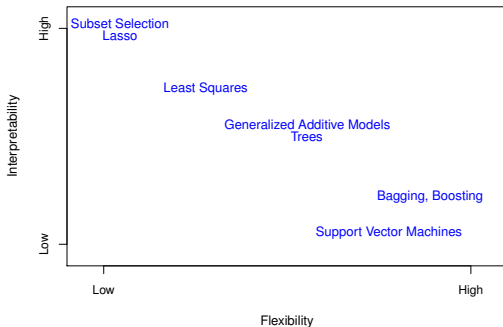
No free lunch theorem

All models are wrong, but some are useful. – George Box

- The field of machine learning proposes a large range of models and algorithms to solve supervised and unsupervised learning problems.
- However, there is no single model or approach that works optimally for all problems. This is sometimes called the **no free lunch theorem**.
- Therefore, applied statistical learning requires awareness of speed-accuracy-complexity trade-offs and data-driven consideration of different approaches for every problem.

Accuracy vs interpretability

Particularly in data mining, interpretability is an important consideration in addition to predictive accuracy. Highly flexible, nonparametric methods, tend to be less interpretable than simpler methods.



Study guide

- Recall three important concepts from these slides, and explain them in your own words.
- Use the review questions in the next slide to self-test on key concepts.
- Study the mathematical details in the module notes.
- Study (or revise) Chapters 1 and 2 of ISL. Reader Chapter 3 before the next module.

Review questions (1/2)

- What is statistical learning?
- What is the difference between supervised and unsupervised learning?
- What is a loss function?
- What do we learn from statistical decision theory for regression problems?
- How do we evaluate model performance with data?

Review questions (2/2)

- What is the difference between the generalisation error and the expected prediction error?
- What is the bias-variance trade-off and why is it important for predictive modelling?
- What is model selection? How is it different from model evaluation?
- What is overfitting?
- What is the difference between parametric and nonparametric models? What are the advantages and disadvantages of each approach?