

Sec 6: Main concepts of statistical learning

STA 35C – Statistical Data Science III

Instructor: Akira Horiguchi

Fall Quarter 2025 (Sep 24 – Dec 12)
MWF, 12:10 PM – 1:00 PM, Olson 158
University of California, Davis

Outline

Based on Chapters 1 and 2 of ISL book James et al. (2021).

1 Statistical learning

2 Unsupervised learning

3 Supervised learning
■ Assessing model accuracy

Statistical learning

How does advertising affect sales? (supervised learning)

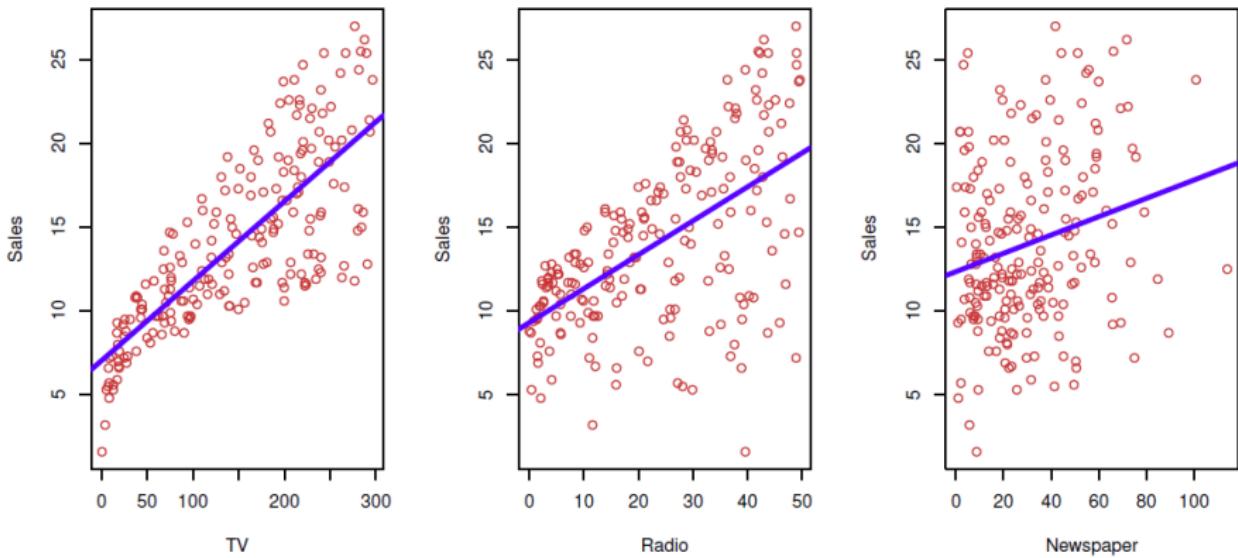


Figure: Image by James et al. (2021), based on the Advertising data set in R. The plot displays sales in thousands of units depending on the input TV, radio and newspaper (advertising) budgets, in thousand dollars, for 200 different markets.

Who will default on credit card payment? (supervised learning)

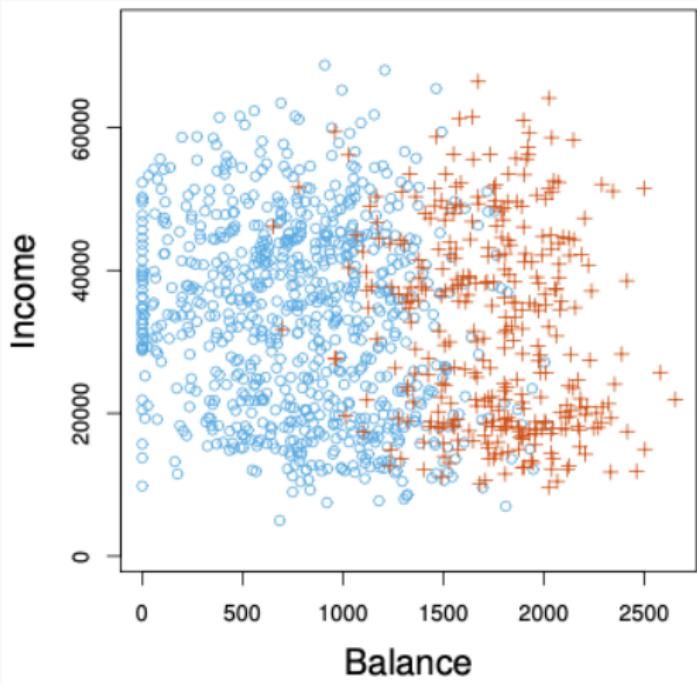


Figure: Image by James et al. (2021). The Default data set. The annual incomes and monthly credit card balances of a number of individuals. Orange +s indicate individuals who defaulted on their credit card payments; blue circles indicate individuals who did not default.

Flow cytometry (unsupervised learning)

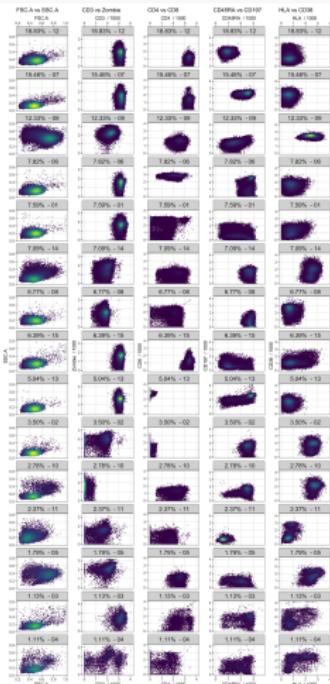


Figure: Image by Horiguchi et al. (2024) –

<https://projecteuclid.org/journals/bayesian-analysis/advance-publication/A-Tree-Perspective-on-Stick-Breaking-Models-in-Covariate-Dependent/10.1214/24-BA1462.full>

Statistical learning: supervised vs unsupervised

Statistical learning refers to a vast set of tools for understanding data.

- **Supervised** statistical learning: predict or estimate an output based on one or more inputs. (STA 142A)
- **Unsupervised** statistical learning: learn relationship or structure among observations. (STA 142B)
- (Are there outputs to “supervise” the learning task?)

Unsupervised learning

Overview

Recall: learn relationship or structure among observations. Example tasks:

- **Dimension reduction**: derive a low-dimensional set of features from higher-dimensional observations X_1, \dots, X_n .
 - ▶ Uses: plotting 2-d representations of higher-dimensional data, regression.
 - ▶ **Principal components analysis** is a popular approach.
- **Cluster analysis**: partition observations X_1, \dots, X_n into distinct groups.

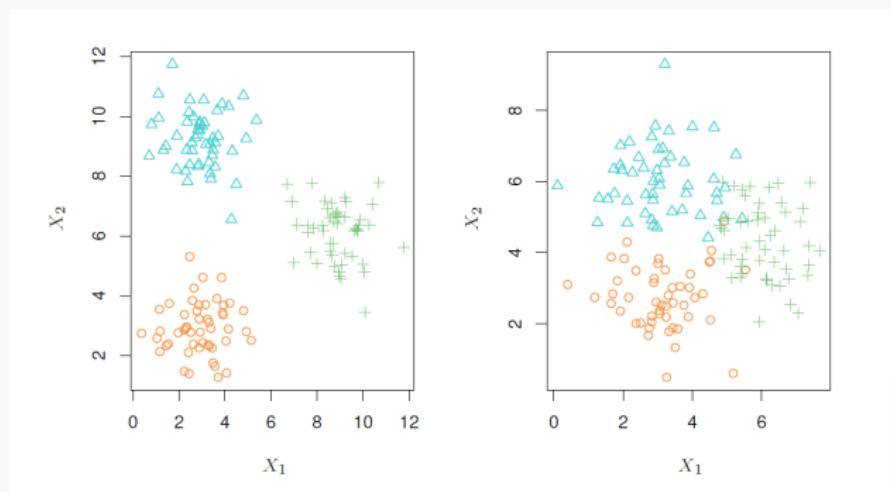


Figure: Image by James et al. (2021). Clustering in a data set involving three groups.

Supervised learning

Non-linear regression

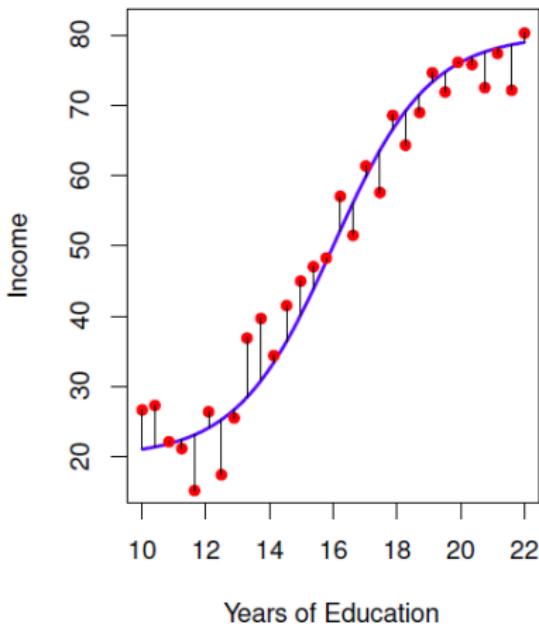
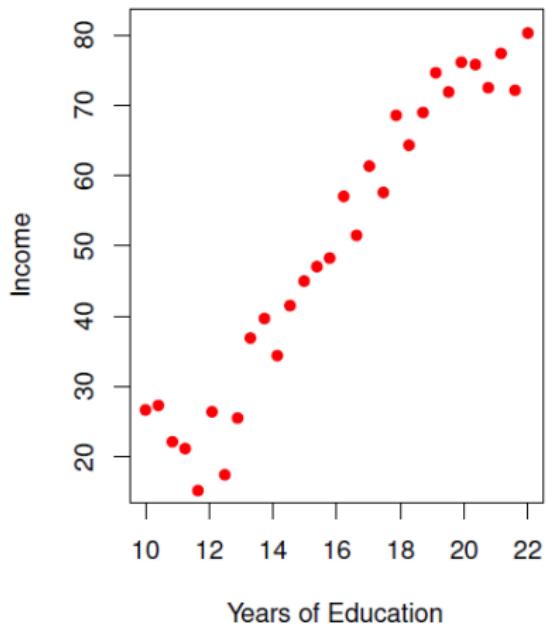


Figure: Image by James et al. (2021), based on the Income data set in R. The red dots are the observed values of income in tens of thousand dollars and years of education for 30 individuals.

A three-dimensional plot

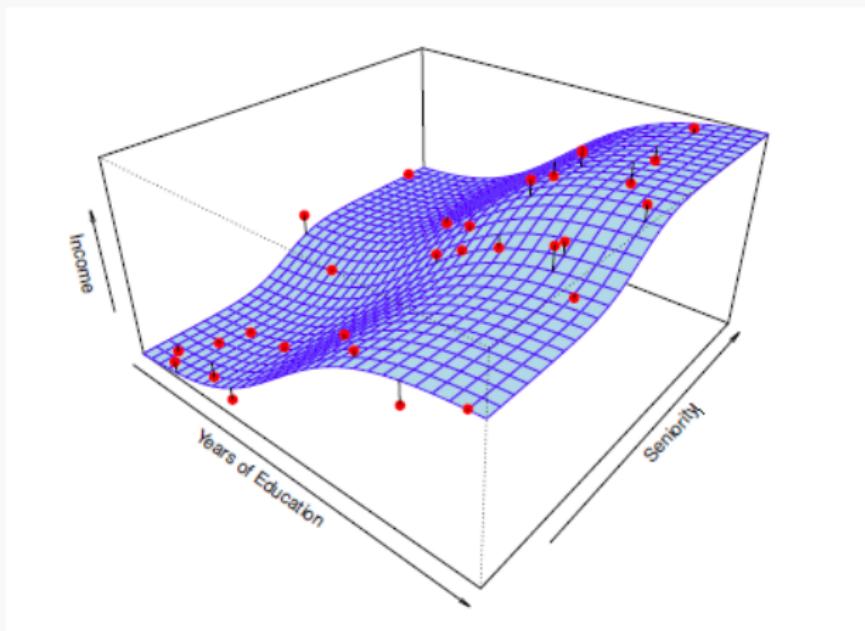


Figure: Image by James et al. (2021), based on the Income data set in R. The income is displayed as a function of years of education and seniority.

Combined plots

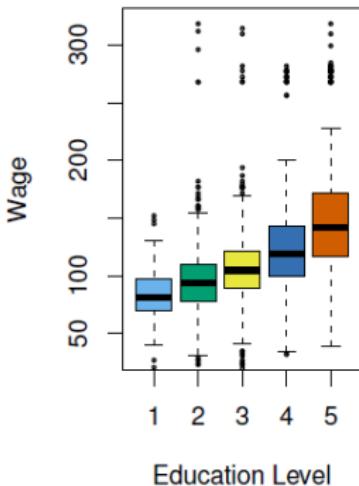
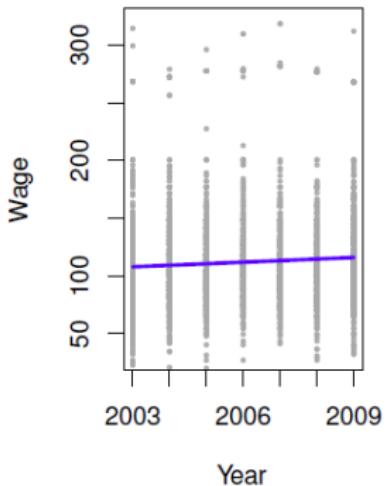
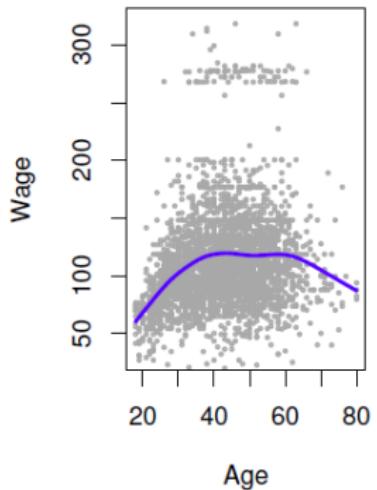


Figure: Image by James et al. (2021), based on the Wage data set in R. The wage is displayed as a function of age, year and education.

Some notation

Recall: predict or estimate an output based on one or more inputs.

- *Input variables* are called *predictors*, *independent variables*, or *features*; denoted by X , and X_1, X_2, X_3 etc. if there is more than one.
- *Output variable* is called *response* or *dependent variable*; denoted by Y .

Example: In Figure 1, the predictors are TV, radio, newspaper, denoted by X_1, X_2, X_3 , respectively, and the response is sales, denoted by Y .

Suppose we observe a numeric response Y and p predictors X_1, \dots, X_p .

- We assume that there is some relationship between Y and $X = (X_1, \dots, X_p)$:

$$Y = f(X) + \varepsilon. \tag{1}$$

- ▶ f denotes a fixed but unknown function of X_1, \dots, X_p .
 - ▶ ε is a random **error term**, which is independent of X , with $E(\varepsilon) = 0$.
- Two main reasons to estimate f : ***prediction*** and ***inference***.

Prediction – i) Idea

Goal: predict response Y at a set of inputs X .

- If X is available, because the error term averages to zero, we can predict Y using

$$\hat{Y} = \hat{f}(X), \quad (2)$$

where \hat{f} denotes the estimate for f , and \hat{Y} the resulting prediction for Y .

- For prediction tasks, \hat{f} is often treated as a **black box** – does it accurately predict Y ?

Example: The blue surface in Figure 6 is an estimate \hat{f} for the unknown function f describing the relationship of the predictors years of education and seniority to the response income:

$$\widehat{\text{income}} = \hat{f}(\text{years of education}, \text{seniority}).$$

Prediction – ii) Accuracy

The accuracy of \hat{Y} for predicting Y depends on *reducible error* and *irreducible error*.

- For a fixed estimator \hat{f} and X , we get

$$\begin{aligned} E[Y - \hat{Y}]^2 &= E[f(X) - \hat{f}(X) + \varepsilon]^2 \\ &= E[f(X) - \hat{f}(X)]^2 + E[\varepsilon^2] + 2E[(f(X) - \hat{f}(X))\varepsilon] \\ &= \underbrace{(f(X) - \hat{f}(X))^2}_{\text{reducible error}} + \underbrace{\text{Var}(\varepsilon)}_{\text{irreducible error}}. \end{aligned}$$

- A more appropriate learning technique might reduce the *reducible error*.
- *irreducible error* comes entirely from the inherent observation noise ε ; independent of how we estimate f .
 - ▶ Even if we would estimate f perfectly, i.e. $\hat{Y} = f(X)$, there is still some *irreducible* prediction error from ε , since $Y = f(X) + \varepsilon$.

Goal: learn relationship between response Y and inputs X_1, \dots, X_p .

- One may want to answer:
 - ▶ Which predictors are associated with the response?
 - ▶ What is the relationship between the response and each predictor?
 - ▶ Can we use a linear equation to describe the relationship between X_1, \dots, X_p to Y , or is there a more complex relationship?
- Knowing more about f allows us to ask questions about Y , such as:
 - ▶ What value of (X_1, \dots, X_p) maximizes Y ?
 - ▶ How much is Y affected by each predictor X_i ?
E.g., in Figure 1 we might have
 - 60% of $\text{Var}(\text{sales})$ can be explained by TV budget,
 - 30% of $\text{Var}(\text{sales})$ can be explained by Radio budget,
 - 8% of $\text{Var}(\text{sales})$ can be explained by Newspaper budget,
 - remaining 2% can be explained by X_4, X_5, \dots, X_p
 - ▶ These questions can be difficult to answer if f is highly non-linear!
 - ▶ <https://www.climateinteractive.org/en-roads/>

How to estimate f ?

Supervised learning: estimate unknown function f using n observed data points

$$(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$$

where $x_i = (x_{i1}, x_{i2}, \dots, x_{ip})^\top \in \mathbb{R}^p$.

- We index the observations by $i = 1, \dots, n$.
- We index the inputs/predictors by $j = 1, \dots, p$.
- Let $x_{ij} \in \mathbb{R}$ represent the value of the j th predictor for the i th observation.
- Let $y_i \in \mathbb{R}$ represent the response variable for the i th observation.

A statistical learning method will use the data to estimate the unknown f .

- I.e., compute a function \hat{f} such that $Y \approx \hat{f}(X)$ for any observed (X, Y) .
- What kind of function is our function estimate \hat{f} allowed to be?
Dictated by our chosen learning method (*parametric* vs *non-parametric*).

Parametric methods

Parametric methods involve two steps:

1. What kind of function is our estimator \hat{f} of $X = (X_1, \dots, X_p)$ allowed to be? E.g.,

$$\hat{f}(X) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p. \quad (3)$$

2. Use the **training data** to **fit/train** the model, i.e., to find the function \hat{f} that best predicts the training responses y_1, \dots, y_n .

Called **parametric** because \hat{f} is determined by a fixed and finite number of parameters.

- Advantage: Reduces the problem of estimating an arbitrary p -variate function f to the easier problem of estimating only $p + 1$ parameters.
- Disadvantage: The best \hat{f} likely might not adequately capture the true form of f .
- The left plot in Figure 1 estimates the function f by

$$\hat{f}(\text{TV, Radio, Newspaper}) = \hat{f}(\text{TV}) \approx 6 + \frac{1}{20} X_1. \quad (4)$$

(Language: “make an assumption about the functional form of f ”)

Non-parametric (or nonparametric) methods

Non-parametric methods: \hat{f} is determined by infinitely many parameters (or # of parameters increases with n).

- Misnomer: non-parametric \neq no parameters!
- Advantage: larger class of possible $\hat{f} \Rightarrow$ more likely that \hat{f} better estimates f .
- Disadvantage: requires a large number of observations to accurately estimate f .
- Disadvantage: inference is often more difficult.

Example: (next slide)

Example – non-parametric

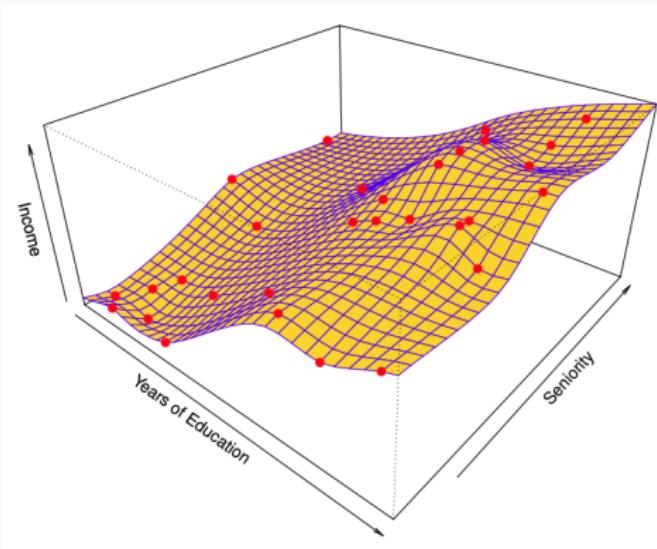


Figure: Image by James et al. (2021), based on the Income data set in R. “A rough thin-plate spline fit to the Income data from Figure 2.3. This fit makes zero errors on the training data.” Here the yellow surface need only be continuous.

- Interpolation theorem: Any n bivariate data points (x_i, y_i) (where x_i s are distinct) can be interpolated by a polynomial of order at most $n - 1$.
https://en.wikipedia.org/wiki/Polynomial_interpolation

Choosing a model: interpretability vs flexibility

Are we more interested in inference or prediction?

- Simpler models typically are easier to interpret and make inference on.
- Flexible models might more accurately predict Y but can be less interpretable.

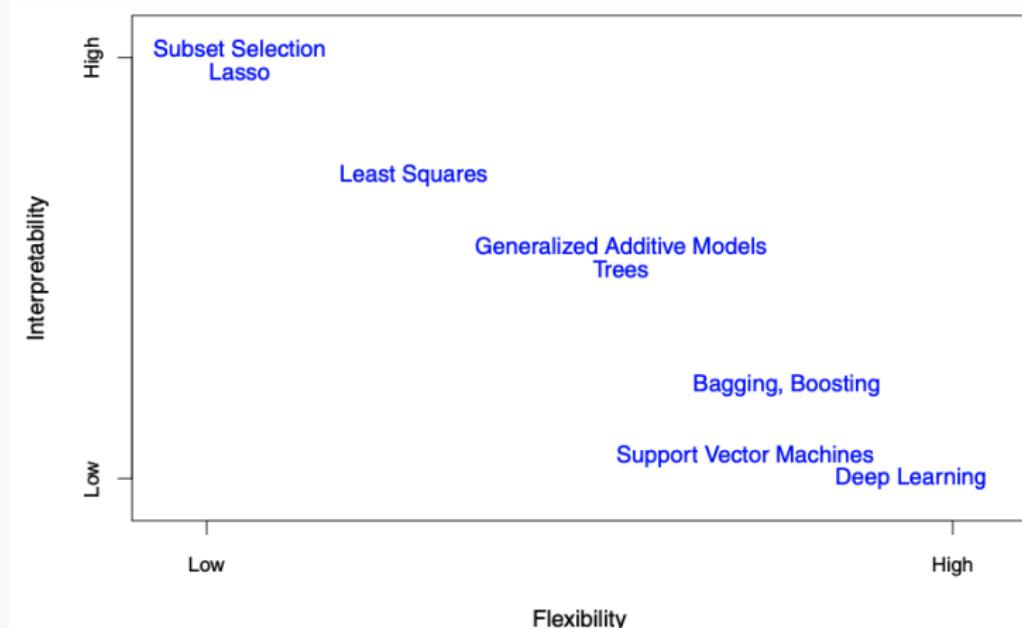


Figure: Image by James et al. (2021). No method dominates *all* other methods.

Regression vs. classification

Describe learning problem based on whether the *response* Y is quantitative or qualitative.

- *quantitative* (numeric). E.g., age, height, income, house value, stock price.
- *qualitative/categorical* (“discrete” – value is one of K classes). E.g., marital status (married or not), brand of product purchased (brand A, B, or C).

Regression vs classification

- Problems w/ a quantitative response: usually referred to as *regression* problems.
- Problems w/ a qualitative response: usually referred to as *classification* problems.

Further notes:

- Distinction is not always crisp; is logistic regression classification or regression?
- Whether the predictors X_1, \dots, X_p are quantitative or qualitative is less important.

Example of classification – K-nearest neighbors

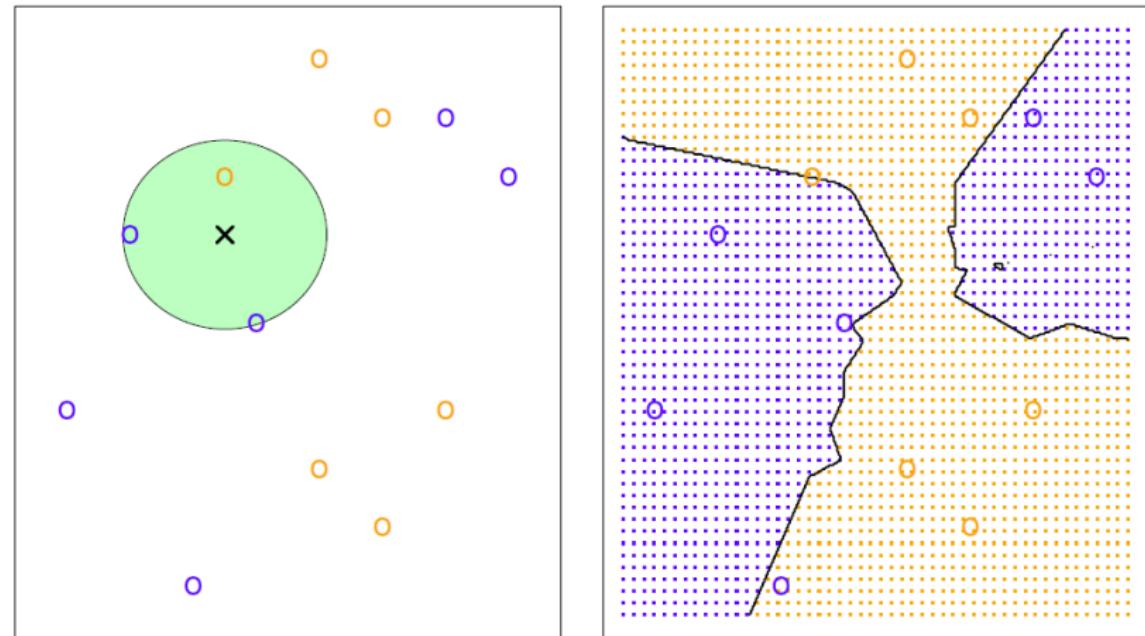


Figure: Image by James et al. (2021). Classification using the K -nearest neighbor approach with $K = 3$. **Left:** We assign a new observation (this is the "x") to the class for which most of three neighbors of "x" belong to. **Right:** A decision line/region for how we would assign a new element for $K = 3$ if they would fall in a certain indicated region.

Supervised learning

Assessing model accuracy

How do we assess model accuracy

Need to establish an error metric. Different for regression vs classification.

Regression: idea

In regression, most commonly use the *mean squared error* (MSE):

$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{f}(x_i))^2, \quad (5)$$

for n data points $(x_1, y_1), \dots, (x_n, y_n)$.

- MSE computed on the *training data* is called the *training MSE*, denoted by MSE_{train} .
- We want to choose an estimator \hat{f} that produces a small MSE on data it wasn't trained on:

$$MSE_{test} = \frac{1}{m} \sum_{i=1}^m (y_{n+i} - \hat{f}(x_{n+i}))^2, \quad (6)$$

for m data points $(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})$. We call (6) the *test MSE*.

Emphasize: set of training data should be disjoint from set of test data.

Regression: a sketch

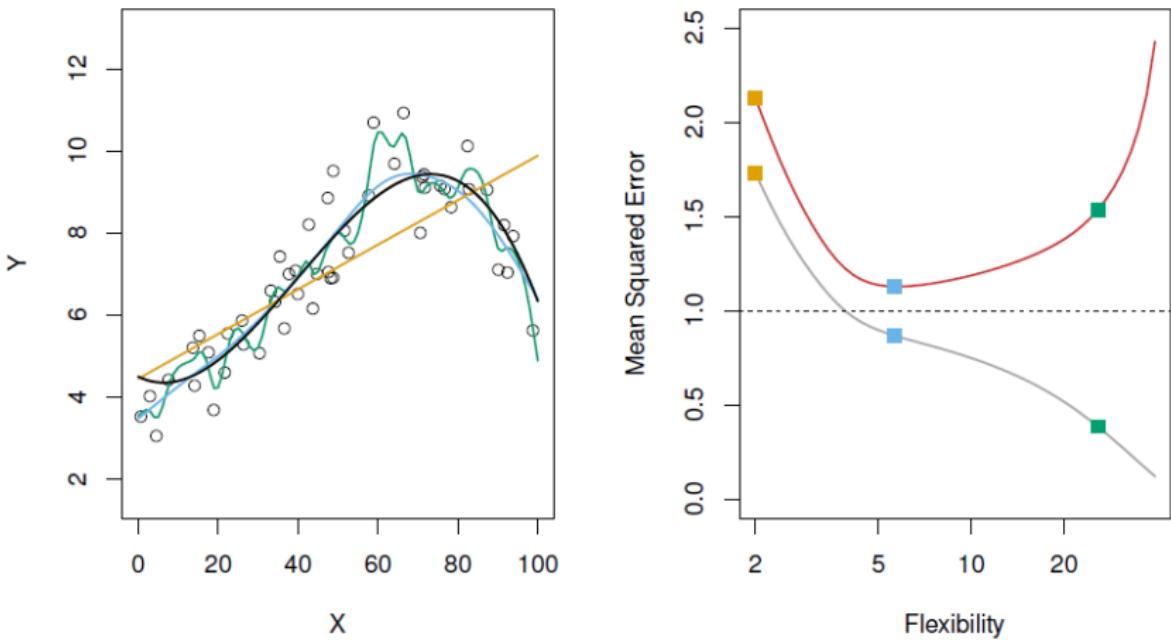


Figure: Image by James et al. (2021). **Left:** Data simulated from f (in black). Three estimates for f : linear regression (in orange), and two smoothing splines (in blue and green). **Right:** Training MSE (in grey), test MSE (in red), and minimum possible test MSE over all methods (dashed line). The squares represent the three fits from the left panel.

Regression: comments on the sketch

- The dashed line in Figure 8 indicates the irreducible error (this is $\text{Var}(\varepsilon)$) which is the lowest achievable test MSE among all possible methods.
- Training MSE *always* decreases as we increase model flexibility.
- The U-shape in the test MSE in Figure 8 indicates that it decreases up to a certain amount of flexibility, but gets worse afterwards.
- When a given method yields a small training MSE, but a large test MSE, we say that the data points are *overfitted*.

Regression: another sketch

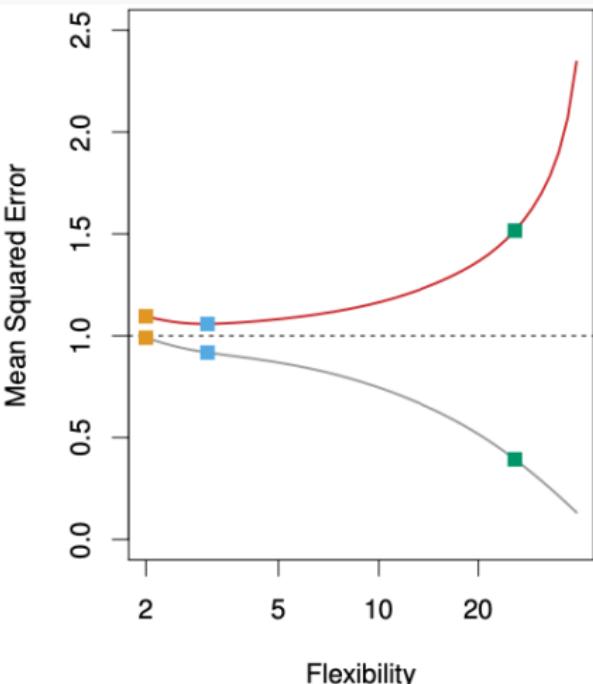
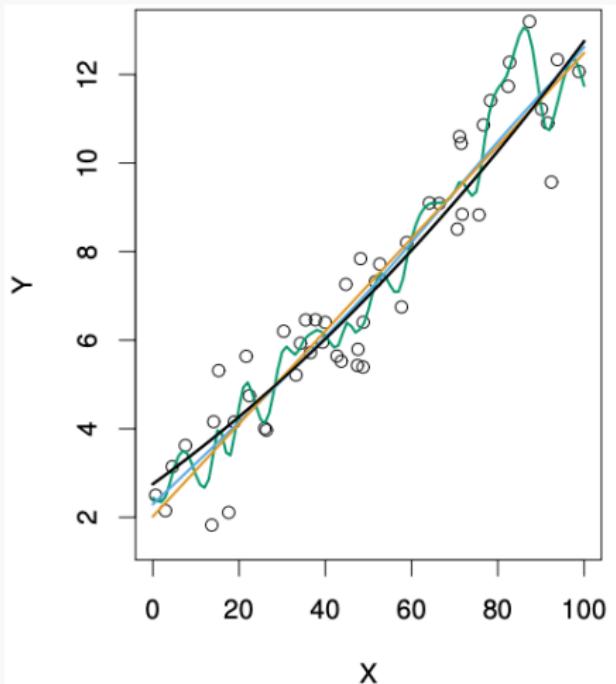


Figure: Image by James et al. (2021). Details are same as in previous figure, but using a different true f that is much closer to linear. In this setting, linear regression provides a very good fit to the data. U-shape is barely noticeable.

Regression: Bias-variance trade-off

U-shape results from two competing properties: *(squared) bias* and *variance*.

- Consider the *expected test MSE* at a new test data point (x, y) .
This is the average test MSE obtained if we repeatedly estimate f over *infinitely many training data sets drawn from some underlying data distribution*.
- The *expected test MSE* at (x, y) can always be decomposed into the sum:

$$E\left[\{y - \hat{f}(x)\}^2\right] = \text{Var}(\hat{f}(x)) + \{f(x) - E[\hat{f}(x)]\}^2 + \text{Var}(\varepsilon). \quad (7)$$

- ▶ $\text{Var}(\hat{f}(x))$ is the amount \hat{f} would change by using different training data sets.
- ▶ $\{f(x) - E[\hat{f}(x)]\}^2$ is the squared bias of $\hat{f}(x)$.
 - Squared bias can be interpreted as the error introduced by approximating f using the given model assumptions (which often do not capture the full complexity of f).
- ▶ $\text{Var}(\varepsilon)$ is the variance of the error term ε (this is a property of the data, not of \hat{f}).
 - The expected test MSE can never be below $\text{Var}(\varepsilon)$. (Dotted horiz line in prev two figures.)

(7) is called a *bias-variance trade-off*. As a general rule, as model flexibility increases, the variance will increase and the (squared) bias will decrease.

Classification: training and test error rate

In classification, we quantify accuracy of \hat{f} using the *training error rate*,

$$Err_{train} = \frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{f}(x_i)), \quad (8)$$

where $(x_1, y_1), \dots, (x_n, y_n)$, with $n \in \mathbb{N}$, are training data.

- $I(y_i \neq \hat{f}(x_i))$ is the *indicator variable* that equals 1 if $y_i \neq \hat{f}(x_i)$ and equals 0 if $y_i = \hat{f}(x_i)$.
- As $\hat{f}(x_i)$ is the predicted class given the observation x_i , the training error rate counts the average number of wrong classifications.
- A good classifier, however, produces a small *test error rate*

$$Err_{test} = \frac{1}{m} \sum_{i=1}^m I(y_{n+i} \neq \hat{f}(x_{n+i})), \quad (9)$$

where $(x_{n+1}, y_{n+1}), \dots, (x_{n+m}, y_{n+m})$, with $m \in \mathbb{N}$, are test data.

Bias-variance trade-off also appears in classification, as we will see.

Classification: Bayes classifier

The test error rate is on average minimized, by a classifier that assigns each observation to the most likely class given its predictor.

- In other words, such a classifier assigns a test observation x to the class

$$\arg \max_j P(Y = j|X = x). \quad (10)$$

- This classifier is called the *Bayes classifier*.
- Special case: if Y must belong to one of only two classes, we predict class 1 if $P(Y = 1|X = x) > 0.5$, and class 2 otherwise.
- The Bayes classifier produces the lowest possible test error rate, called the *Bayes error rate* — analogous to the irreducible error discussed earlier.

Classification: KNN – i) Idea

- In theory, we would always like to predict using the Bayes classifier, but in practice we don't know the conditional distribution of Y given X .
- Many approaches attempt to estimate the conditional probabilities.
- The *K-nearest neighbor* (KNN) classifier estimates the Bayes classifier by counting the $K \in \mathbb{N}$ closest values of x in a neighborhood \mathcal{N}

$$P(\widehat{Y = j | X = x}) = \frac{1}{K} \sum_{i \in \mathcal{N}} I(y_i = j), \quad (11)$$

and assigns x to the class j with the highest probability.

Classification: KNN – ii) An example

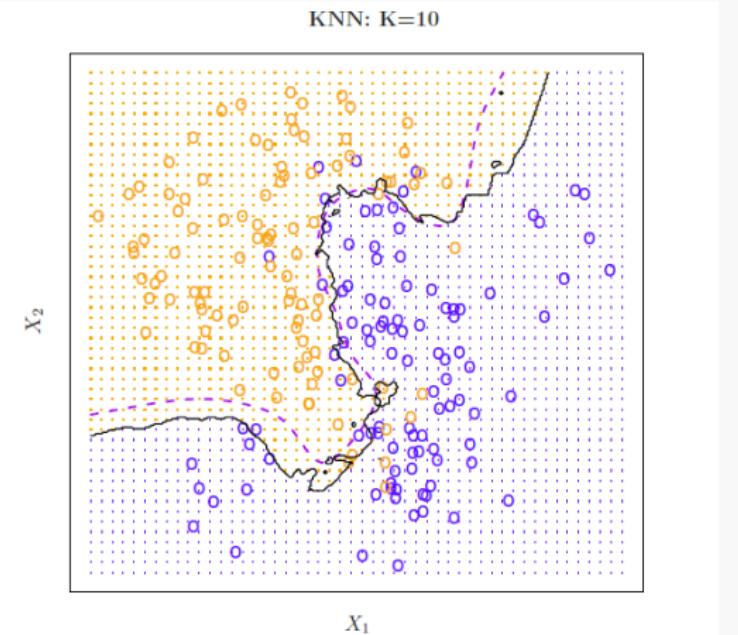


Figure: Image by James et al. (2021). The KNN decision boundary for $K = 10$ (in black), and the Bayes decision boundary (in purple).

Classification: KNN – iii) Different K

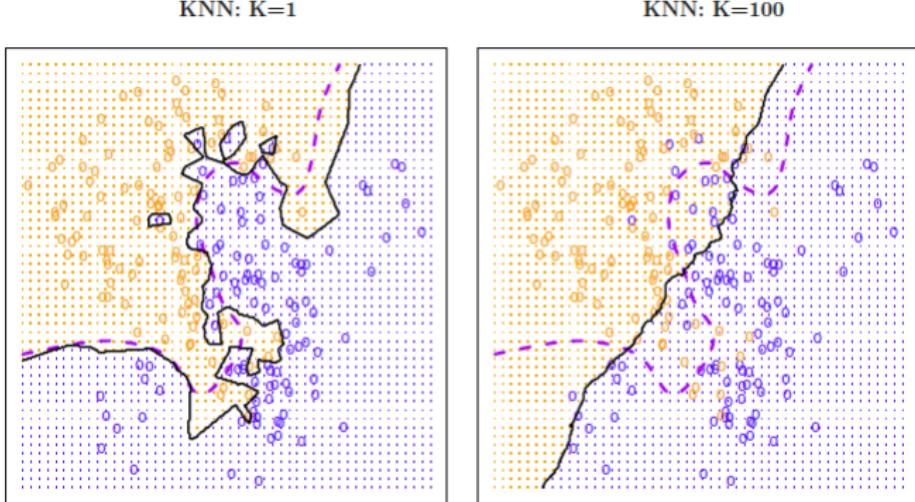


Figure: Image by James et al. (2021). KNN decision boundaries (in black) for $K = 1$ and $K = 100$, and the Bayes decision boundary (in purple). With $K = 1$, the decision boundary is overly flexible, while with $K = 100$ it is not sufficiently flexible.

Classification: KNN – iv) Bias-variance trade-off

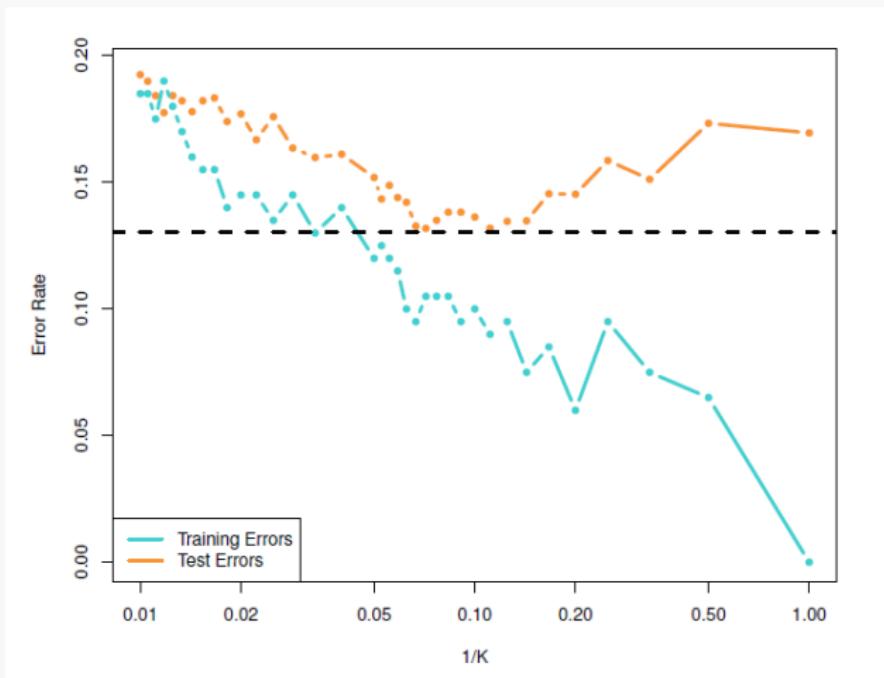


Figure: Image by James et al. (2021). The KNN training error rate (in blue) and test error rate (in orange) and the Bayes error rate (in black). The jumpiness of the curves is due to the small size of the training data set.

General comments

We want to choose the method that produces the smallest test error.

- True regardless of regression or classification.