

Lecture02: An Introduction to Machine Learning

- Course and textbook objective: Democratizing machine learning with commonly accessible soft- and hardware and regain control over our destiny rather than letting robots or artificial intelligence (AI) agents rule us.
- The job of robots and AI agents should be to make our life easier.
- However, we are running the risk of becoming too dependent on or manipulated by them.

Useful Websites

- <https://www.r-bloggers.com/in-depth-introduction-to-machine-learning-in-15-hours-of-expert-videos/>
- <https://cran.r-project.org/web/views/MachineLearning.html>
- www.kaggle.com

Historical Background

- AI objective: Perform predictions of outcomes y or future actions by using specific AI rules (i.e., functions $f(x)$) based on observable information x .

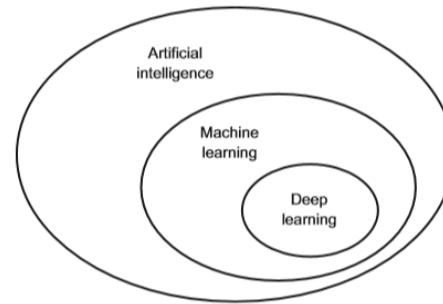


Figure 1.1 Artificial intelligence, machine learning, and deep learning

Artificial Intelligence

- In the 1950 computer pioneers start asking the question “Can computers think?”.
- Ultimately its definition was reduced to automating intellectual tasks normally performed by humans.
- The tasks were hardcoded as “rules” and lead to symbolic AI and expert systems in the 1980.
- AI can solve well-defined logical problems, but it cannot hand complex and fuzzy problems.

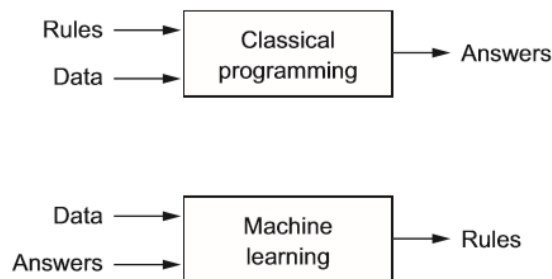


Figure 1.2 Machine learning: a new programming paradigm

Machine Learning (ML)

- ML started in the 1990s asking the question “can a computer discover hidden rules by processing data?”
- In **supervised ML** humans input data and answers, whereas the computer finds rules connecting data with answers.
- This is achieved by training the “machine”.
- Once rules are established with some confidence, they can be used to perform predictions of answers based on new input data.
- In **unsupervised ML** the algorithm search for a structure (i.e., a pattern) in the data to establish homogeneous groups of observations.
- In general, ML is tightly related to statistics but is applied on large and complex data sets for which standard statistical models would be too rigid.
- Progress in ML is mainly driven by statisticians and computer scientist. Their different backgrounds can lead to inconsistencies.
- ML is applied in big data analysis which are characterized by [a] volume, [b] variety, [c] velocity, and [d] veracity (that is, certainty or the lack thereof).
- Criticism of ML: it is data drive and ideas are “proven” empirically rather than being rooted in theory.

Deep Learning

- Deep Learning (DL) is a specialized subfield of ML based neural networks. Small computational neural networks emerged in the early 1960ies. They were able to handle only a small set of hidden layers and thus provided a coarse approximation of the underlying data structure.
- With the increase in affordable computing power and an increasing availability of data, DL with many hidden layers became prominent in the early 2010.
- “Deep” does not refer to a deeper understanding of the task.
- It is based on successive layers of rules with increasingly more meaningful representations of the outcome y .
- DL is based on the mathematical framework of neural networks but does not have any relationship to neurobiology (i.e., neurobiology studies how our brains work).
- One can think of neural network as a series of filters (the layers) that distill the input information x into a purified version that is closely related to the output y .
- Building blocks of understanding neural networks are logistic regression and gradient descent optimizers.
- A sequence of layers weight feature input vectors \mathbf{x}_i that will generate an output label y_i :
$$y_i = f_{NN}(\mathbf{x}_i) = f_3(f_2(f_1(\mathbf{x}_i))).$$

- DL involves more than two hidden layers. In $y_i = f_{NN}(\mathbf{x}_i) = f_3(f_2(f_1(\mathbf{x}_i)))$ the function f_3 is the output layer.

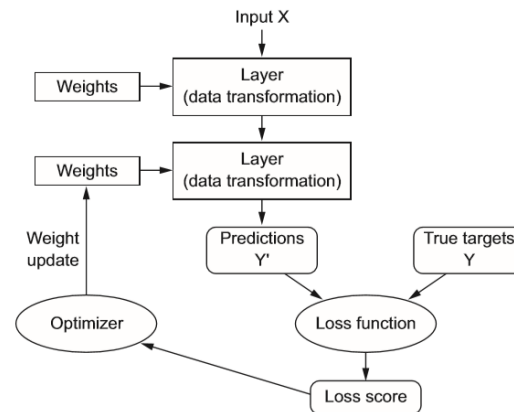


Figure 1.9 The loss score is used as a feedback signal to adjust the weights.

- The objective of specifying the l -th layer functions $f_l(\)$ is to minimize the prediction error, which is measured by a loss score functions.

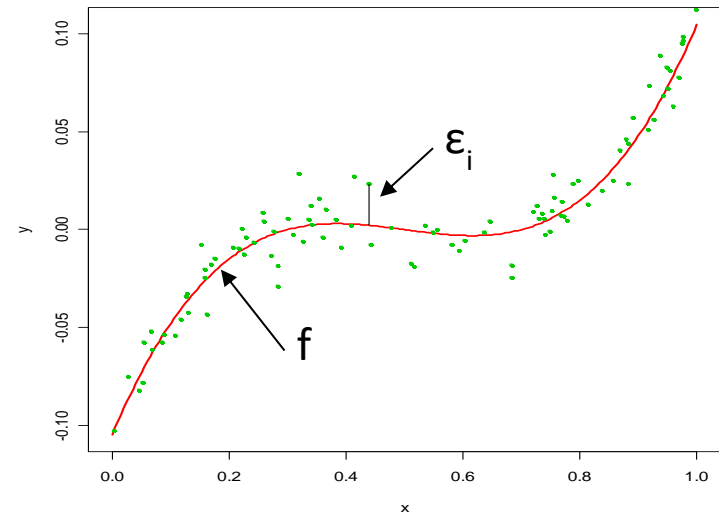
Fundamental Setting of Supervised ML

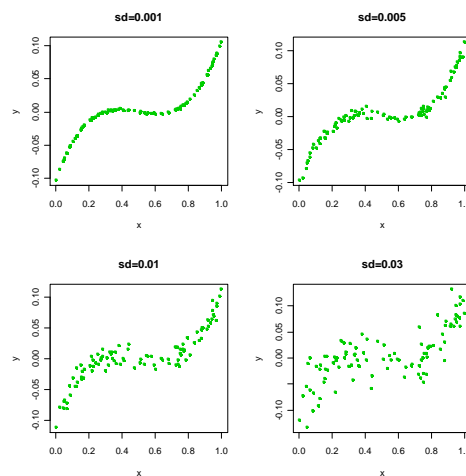
Some of the figures in this presentation are taken from Chapter 2 of "An Introduction to Statistical Learning, with applications in R (Springer, 2013) with permission from the authors: G. James, D. Witten, T. Hastie and R. Tibshirani "

- In supervised learning both the outcome y_i and the features \mathbf{x}_i of object i are observable.
- The outcome and the features are connected by a *true* function $f(\cdot)$ plus an unknown random error ε_i :

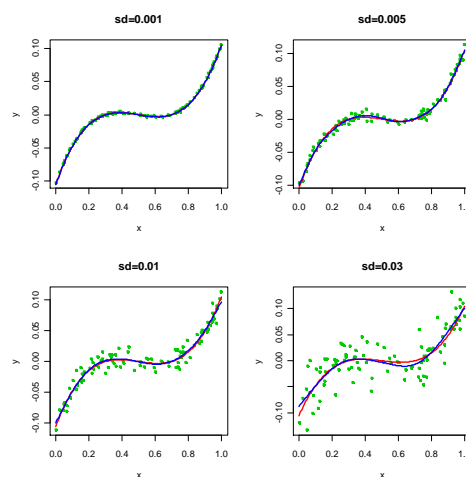
$$y_i = f(\mathbf{x}_i) + \varepsilon_i$$

- As the variance of the error increases the exact function form of $f(\cdot)$ becomes more and more uncertain.





- There are different estimates $\hat{f}(\cdot)$ for true but unknown $f(\cdot)$



- In supervised learning a **training sample** with (y_i, \mathbf{x}_i) is used to obtain an “accurate” estimate $\hat{f}(\cdot)$ of $f(\cdot)$.

- At **test sample** (y_0, \mathbf{x}_0) is usually withheld to **evaluate the accuracy** of the estimator $\hat{f}(\)$
- Since many different estimators $\hat{f}(\)$ of varying degrees of flexibility can be specified there is usually a trade-off between prediction accuracy and model generality (see `R`-script **ThinPlateOverfitReg.r**).
- Overfitting leads to estimators $\hat{f}(\)$ well adapted a particular given sample and by also modeling its specific random error ε , but they do not generalize to another sample which will have a different random error ε .
→ There is no *free lunch* in statistics: not one estimation model dominates all others over all possible data sets.

Evaluating the Model Fit

- Mean-Square-Error of the training data set is: $MSE = \frac{1}{n} \cdot \sum_{i=1}^n \left(y_i - \hat{f}(\mathbf{x}_i) \right)^2$.
- Other goodness of fit measures, such as the “confusion” matrix and ROC (receiver operating characteristics) are used for classification problems.

- Example: Medium use of estimated parameters and high error variance in the DGP:

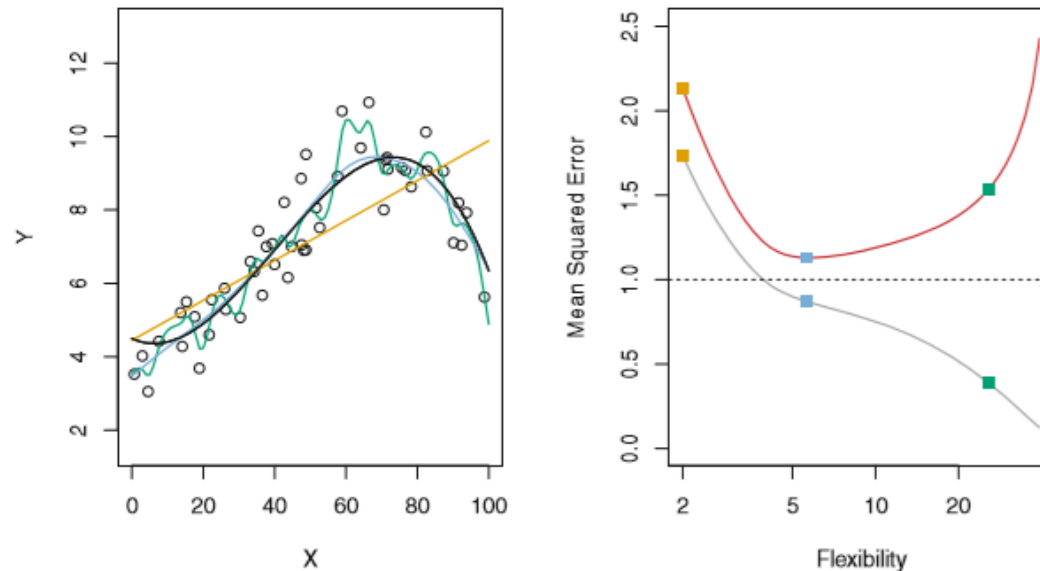


FIGURE 2.9. Left: Data simulated from f , shown in black. Three estimates of f are shown: the linear regression line (orange curve), and two smoothing spline fits (blue and green curves). Right: Training MSE (grey curve), test MSE (red curve), and minimum possible test MSE over all methods (dashed line). Squares represent the training and test MSEs for the three fits shown in the left-hand panel.

- Black line is the true underlying data generating process (DGP) with a random error. The horizontal dashed line is the MSE solely associated with the disturbances.
- Trade-off between Variance and Bias for the prediction dataset y_0 and x_0
- Decomposition of expected squared prediction error for the testing dataset becomes

$$E(y_0 - \hat{f}(x_0))^2 = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon)$$

$E(y_0 - \hat{f}(x_0))^2$	Expected <i>MSE</i> over all possible test samples (y_0, x_0)
$\text{Var}(\hat{f}(x_0))$	Variation of the <i>estimated</i> function $\hat{f}()$ based on the use of different training samples. Simpler models are more stable with less variable in $\hat{f}()$. Simpler models do not run the risk of modeling the sample specific random variability ε .
$[\text{Bias}(\hat{f}(x_0))]^2$	This variation error is introduced by mis-specifying the true model $f()$ of the DGP. Under-fitting a simple model to more complex data variations is a certain cause. A proper selected model $\hat{f}()$ will have no bias.
$\text{Var}(\varepsilon)$	Irreducible error variation. Overfitted models try to capture this random variation. This is the smallest value that the expected MSE can achieve for a perfectly fitting $\hat{f}()$.

- General rule: As the model becomes more flexible the bias will decrease and the variance will increase.

- Example: Low use of parameters and high error variance of DGP:

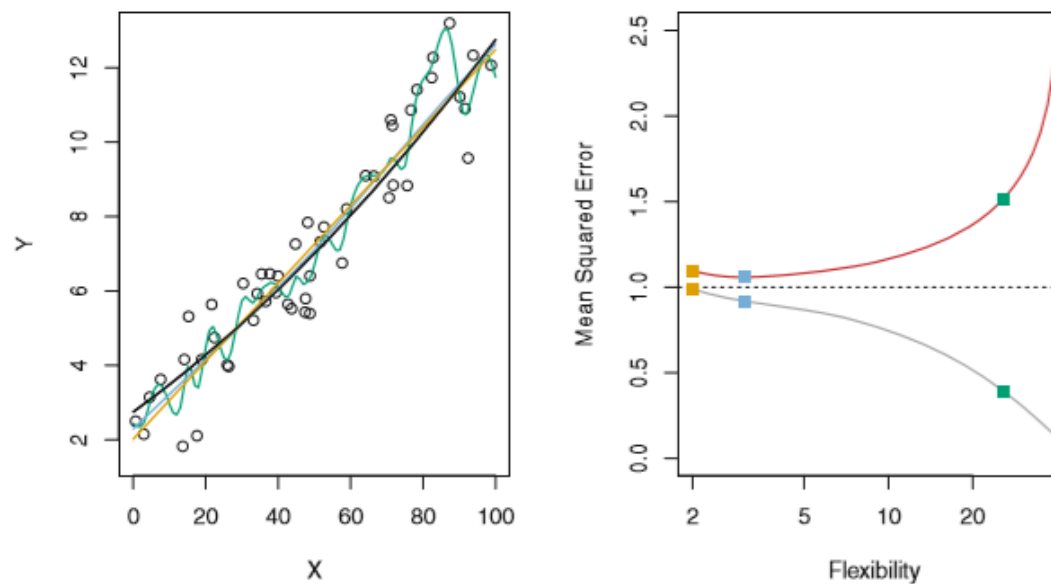


FIGURE 2.10. Details are as in Figure 2.9, using a different true f that is much closer to linear. In this setting, linear regression provides a very good fit to the data.

- Example: High use of estimated parameters and low error variance of DGP

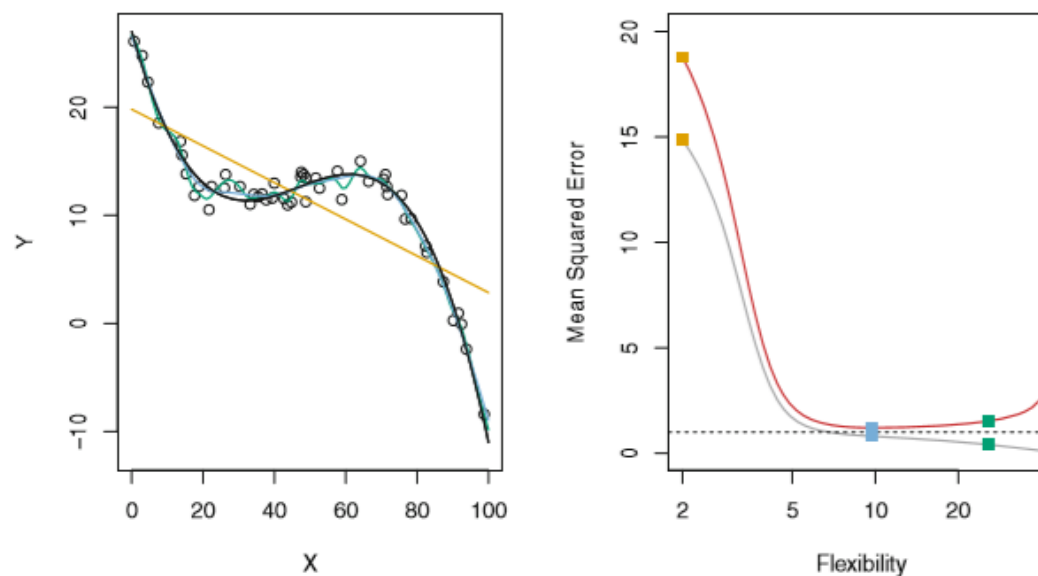


FIGURE 2.11. Details are as in Figure 2.9, using a different f that is far from linear. In this setting, linear regression provides a very poor fit to the data.

- Effect of the model specification on variance and bias in dependence of the model complexity:

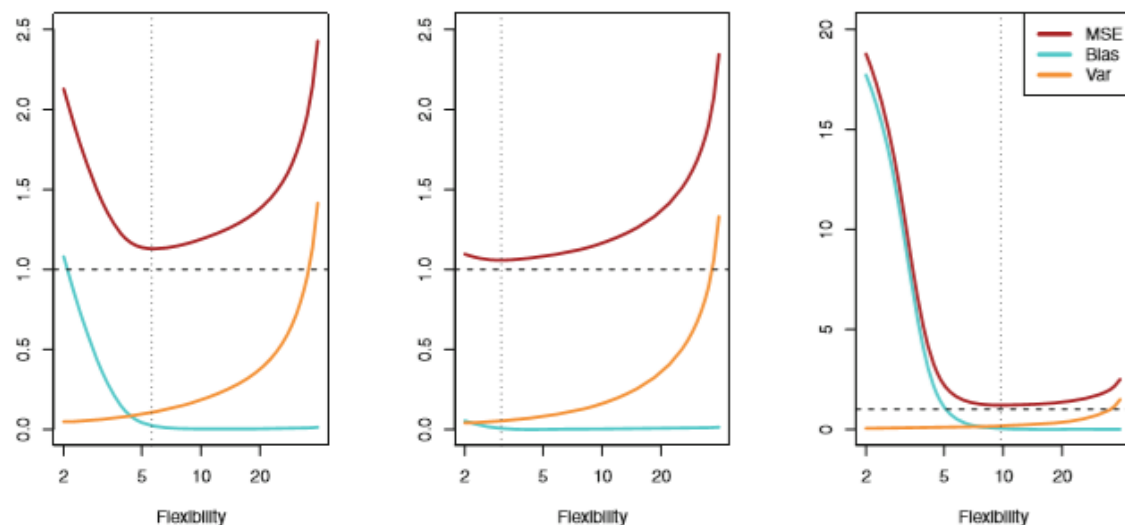


FIGURE 2.12. Squared bias (blue curve), variance (orange curve), $\text{Var}(\epsilon)$ (dashed line), and test MSE (red curve) for the three data sets in Figures 2.9–2.11. The vertical dotted line indicates the flexibility level corresponding to the smallest test MSE.

Classification: Bayesian k-Nearest Neighbor Class Assignment

- In the classification setting the outcome variable is discrete, that is,

$$y_i = k \text{ if object } i \text{ belongs to the } k^{\text{th}} \text{ class}$$

- The error rate, which we aim to minimize, becomes $\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$, where the indicator function $I(\cdot) = 1$ if the condition is true and otherwise it is 0.
- A Bayesian classifier assigns an observation i to the class k if its conditional probability $\Pr(Y = k|X = \mathbf{x}_0)$ is larger than that for any other class.

For just two classes the probability must be $\Pr(Y = k|X = \mathbf{x}_0) > 0.5$.

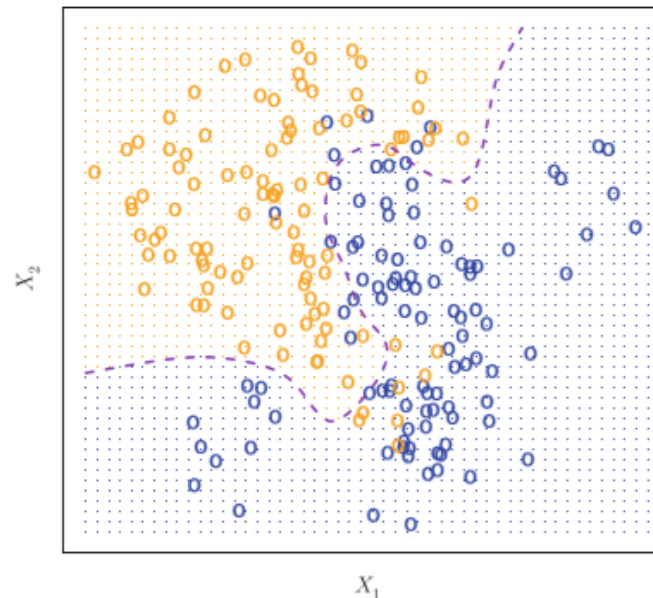


FIGURE 2.13. A simulated data set consisting of 100 observations in each of two groups, indicated in blue and in orange. The purple dashed line represents the Bayes decision boundary. The orange background grid indicates the region in which a test observation will be assigned to the orange class, and the blue background grid indicates the region in which a test observation will be assigned to the blue class.

- The Bayesian class membership probabilities in K-nearest neighborhood (kNN) is calculated as:

- Determine those K training points which are the closest to the prediction point \mathbf{x}_0 . Closeness is measured in terms of the Euclidian distance between two objects i and j with their feature vectors $\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{iP})^T$ and $\mathbf{x}_j = (x_{j1}, x_{j2}, \dots, x_{jP})^T$ by

$$d_{ij} = \sqrt{\sum_{p=1}^P (x_{ip} - x_{jp})^2}.$$

- Subsequently estimate the Bayesian group probabilities by $\Pr(Y = k|X = \mathbf{x}_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = k)$, where \mathcal{N}_0 is the neighborhood around \mathbf{x}_0 with K training points.
- Assign the observation to class j with the highest probability.

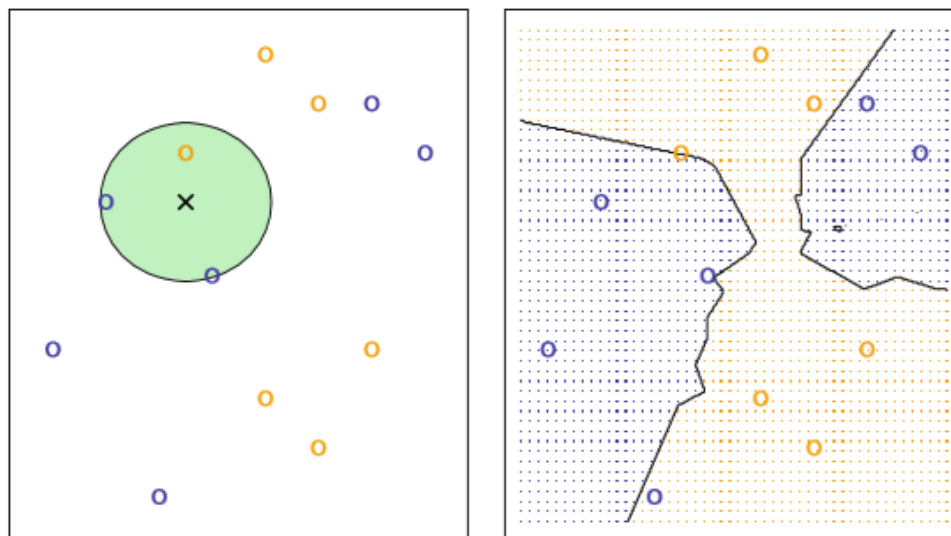


FIGURE 2.14. The KNN approach, using $K = 3$, is illustrated in a simple situation with six blue observations and six orange observations. Left: a test observation at which a predicted class label is desired is shown as a black cross. The three closest points to the test observation are identified, and it is predicted that the test observation belongs to the most commonly-occurring class, in this case blue. Right: The KNN decision boundary for this example is shown in black. The blue grid indicates the region in which a test observation will be assigned to the blue class, and the orange grid indicates the region in which it will be assigned to the orange class.

- The user sets the hyper-parameter K to the desired number of nearest neighbors. This hyper-parameter K determines flexibility of kNN algorithm:

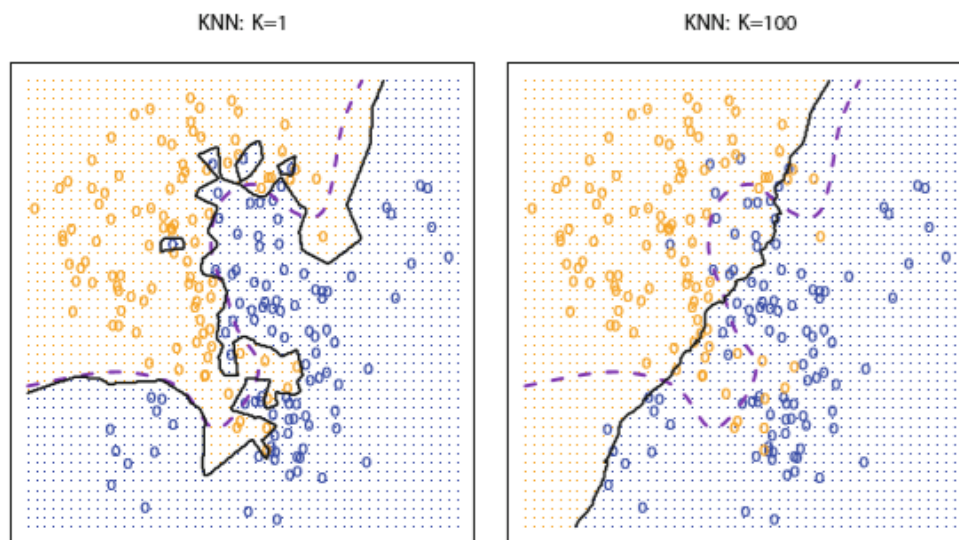


FIGURE 2.16. A comparison of the KNN decision boundaries (solid black curves) obtained using $K = 1$ and $K = 100$ on the data from Figure 2.13. With $K = 1$, the decision boundary is overly flexible, while with $K = 100$ it is not sufficiently flexible. The Bayes decision boundary is shown as a purple dashed line.

- As well as the test and training error rates:

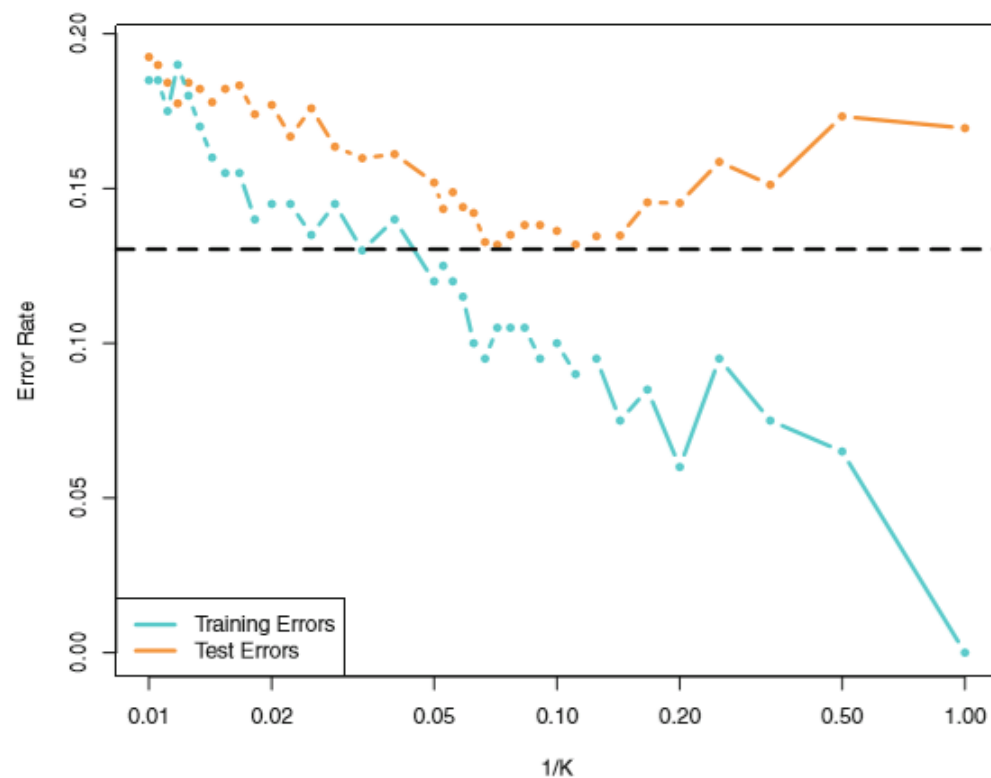



FIGURE 2.17. The KNN training error rate (blue, 200 observations) and test error rate (orange, 5,000 observations) on the data from Figure 2.13, as the level of flexibility (assessed using $1/K$) increases, or equivalently as the number of neighbors K decreases. The black dashed line indicates the Bayes error rate. The jumpiness of the curves is due to the small size of the training data set.

- The function `class::knn()` performs kNN model calibrations (see  script `kNN.R`).
- Advantages of the kNN model:

- Simple and effective estimation which gives by default the probability of the predicted group membership.
- Makes no assumptions about the underlying distribution of the features.
- Virtually no training time required.
- Disadvantages of the kNN model:
 - As the number of features available to classify the data increases, the decision space becomes more sparsely populated which increases the overall uncertainty.
 - How are ties broken?
 - Features that are not relevant for the classification task are not eliminated (ML concept of **regularization**) thus there is the potential of overfitting the data.
 - The performance of the algorithm depends on how **comparable** the features are scaled.
 - For each test data point \mathbf{x}_0 the distances need to be calculated to identify the k nearest neighborhood \mathcal{N}_0 .
 - Euclidian distances for nominal scaled features and missing feature information require additional processing.
 - The algorithm's prediction performance depends on the the hyper-parameter K .

Fundamental Setting of Unsupervised ML

- In unsupervised learning, only the object features \mathbf{x}_i but not the outcome y_i are known to the investigator.

- A function $f(\cdot)$ is sought that allows us to guess what y could be.
- An example is by grouping objects together into clusters. With respect to features \mathbf{x}_i the clusters are supposed to be internally as homogenous as possible, but as heterogeneous (distinct) from each other.
- Again, fuzziness between the cluster delimitations determines the success of the classification rule.

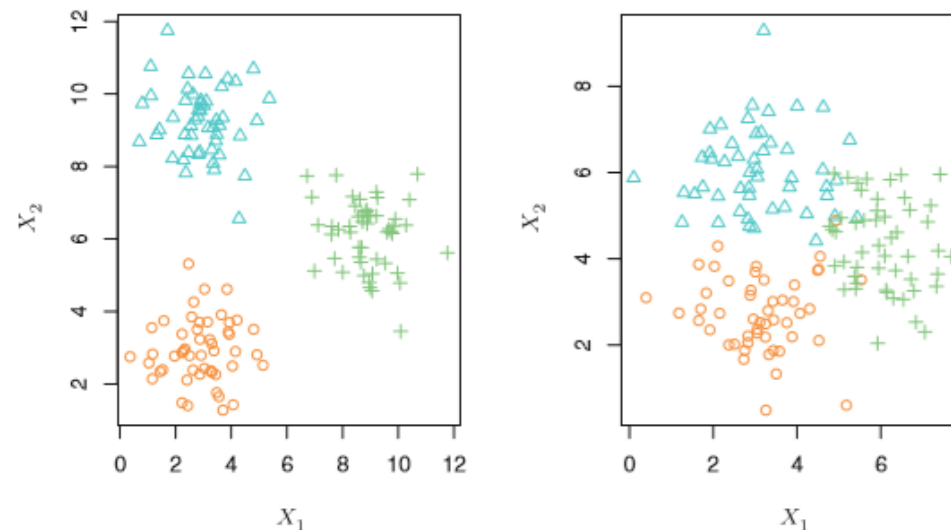



FIGURE 2.8. A clustering data set involving three groups. Each group is shown using a different colored symbol. Left: The three groups are well-separated. In this setting, a clustering approach should successfully identify the three groups. Right: There is some overlap among the groups. Now the clustering task is more challenging.

- Example: Identifying a multivariate mixture distribution
 - The distribution of points with multiple centers and ellipsoidal standard distances can be modeled by a mixture of multi-normal distributions:

$$f \begin{bmatrix} x_i \\ y_i \end{bmatrix} \sim \sum_{p=1}^P \pi_p \cdot \mathcal{N} \left(\begin{bmatrix} \mu_{x,p} \\ \mu_{y,p} \end{bmatrix}, \begin{bmatrix} \sigma_{x,p}^2 & \sigma_{xy,p} \\ \sigma_{xy,p} & \sigma_{y,p}^2 \end{bmatrix} \right) \text{ with } \pi_p > 0 \text{ and } \sum_{p=1}^P \pi_p = 1$$

- Unknowns are:
 - [a] the number of clusters P ,
 - [b] the proportion of points per cluster π_p ,
 - [c] the P cluster centers $\begin{bmatrix} \mu_{x,p} \\ \mu_{y,p} \end{bmatrix}$ and
 - [d] their P ellipsoidal covariance matrices $\begin{bmatrix} \sigma_{x,p}^2 & \sigma_{xy,p} \\ \sigma_{xy,p} & \sigma_{y,p}^2 \end{bmatrix}$.
- The probability $\Pr \left(p \mid \begin{bmatrix} x_i \\ y_i \end{bmatrix} \right)$ of observation $\begin{bmatrix} x_i \\ y_i \end{bmatrix}$ belonging to cluster p is evaluated using the Bayesian theorem.
- See the  script **findClusters.r**