

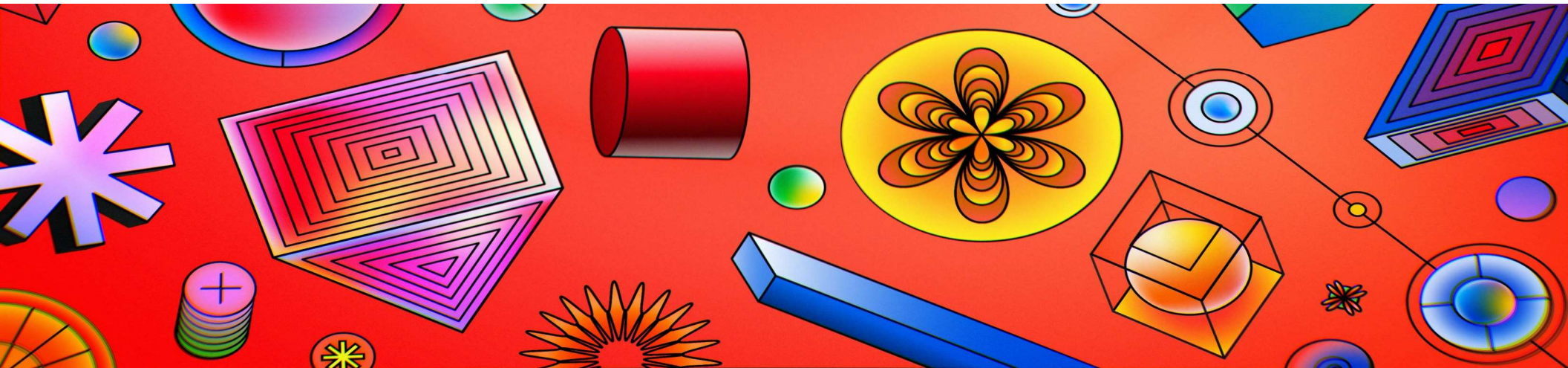
Fall 2023

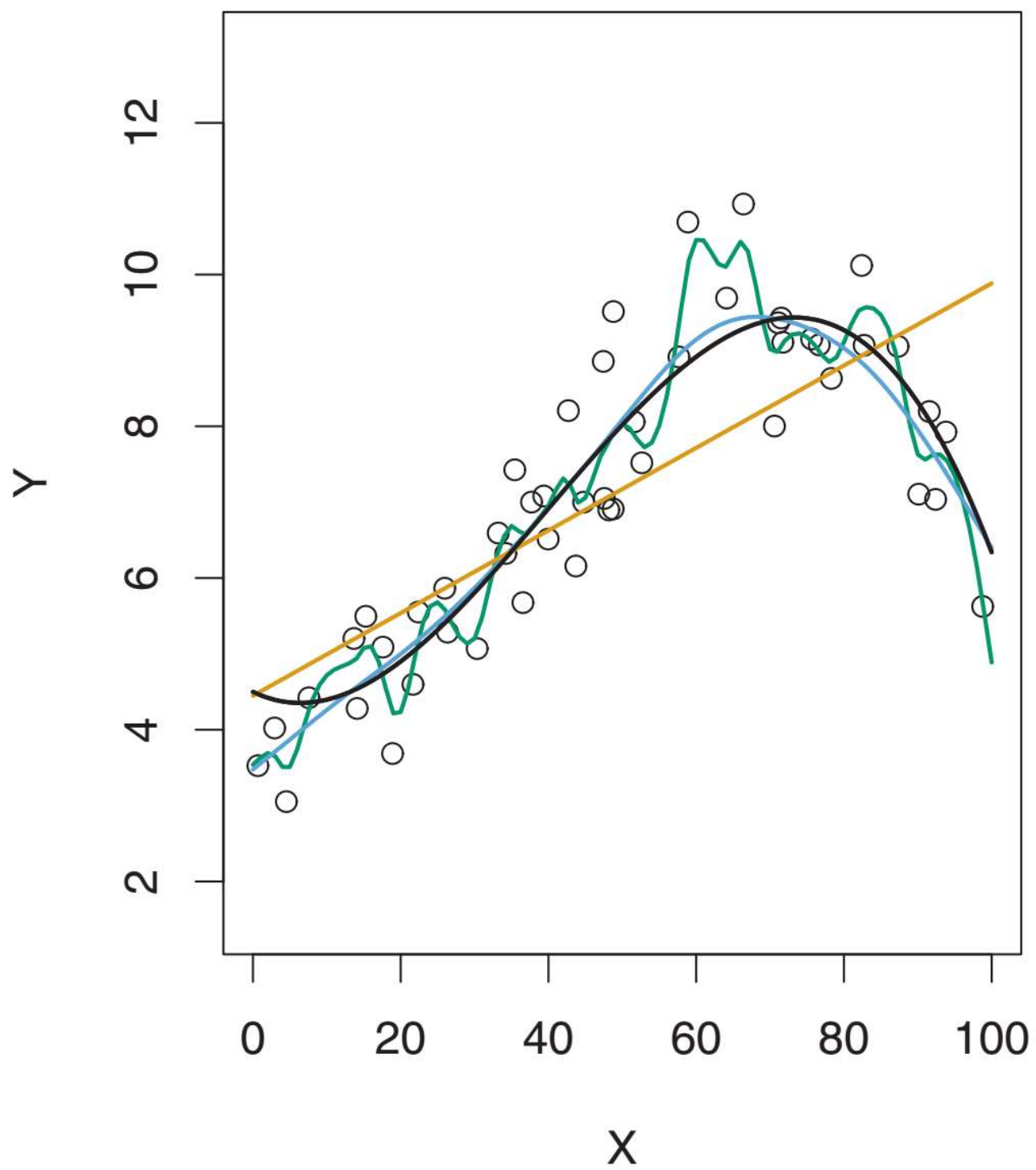
BIF524/CSC463 Data Mining

Statistical Learning

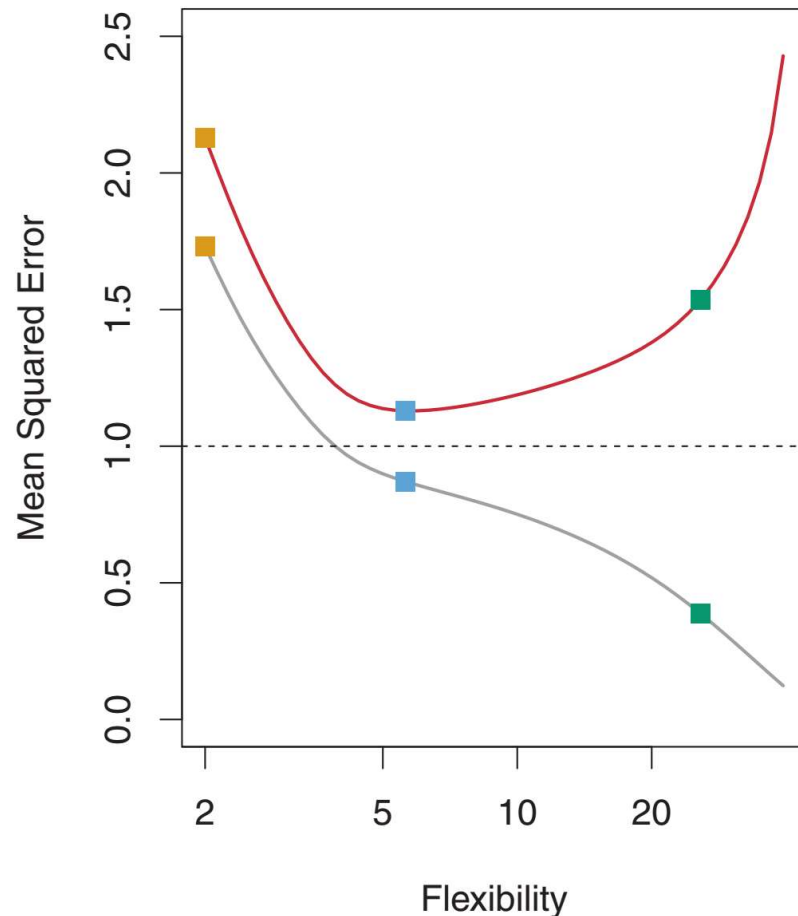
Eileen Marie Hanna, *PhD*

12/09/2023



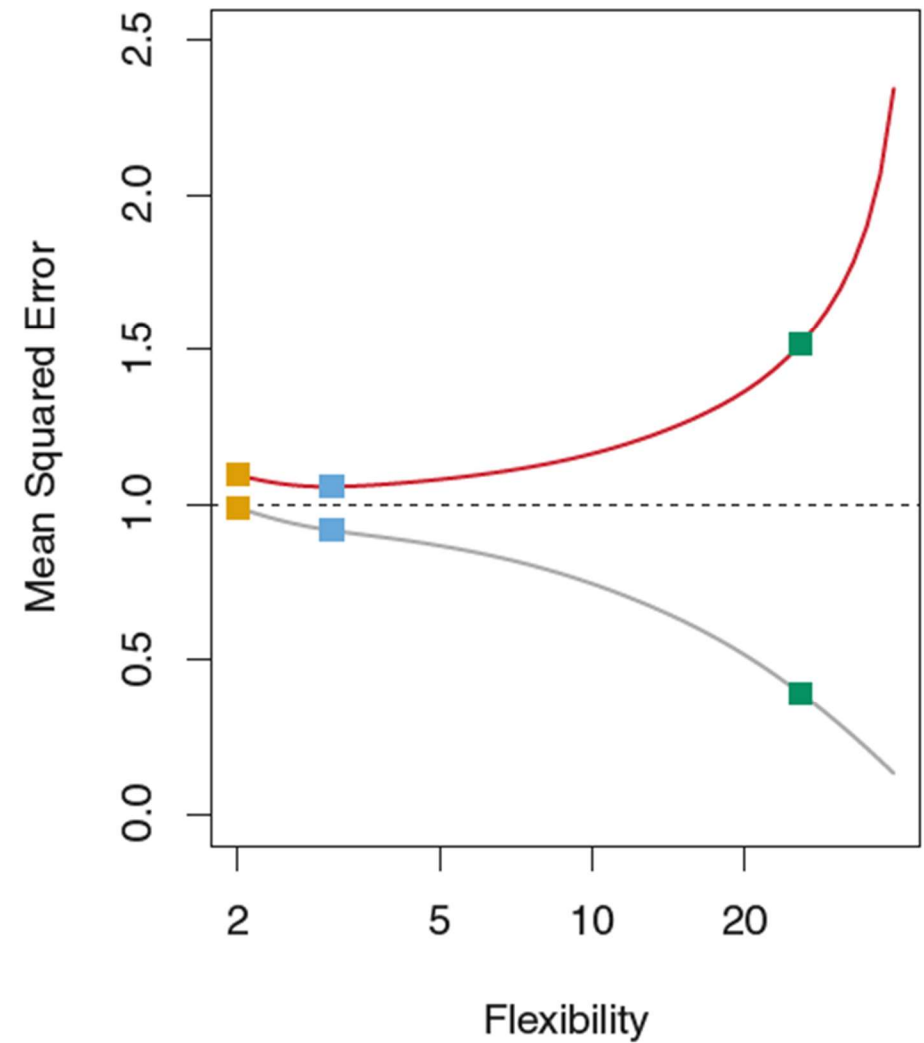
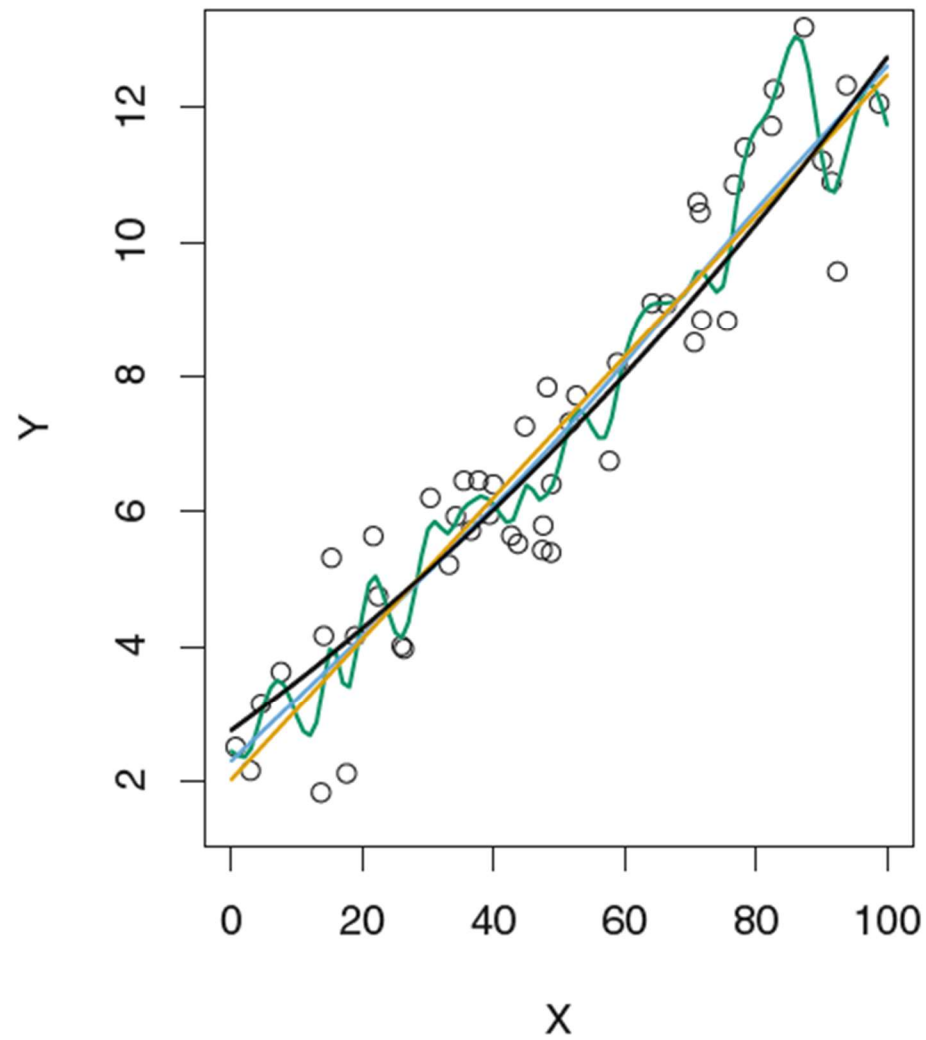


Mathematically

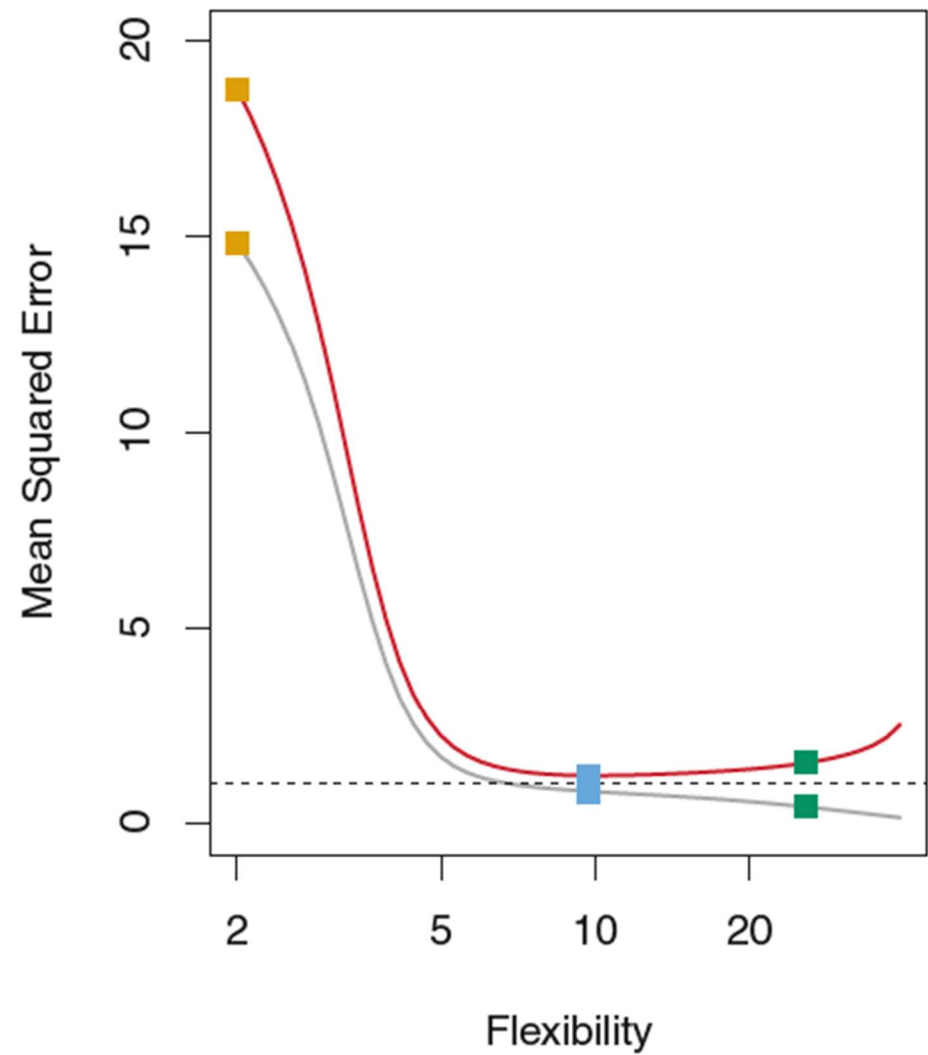
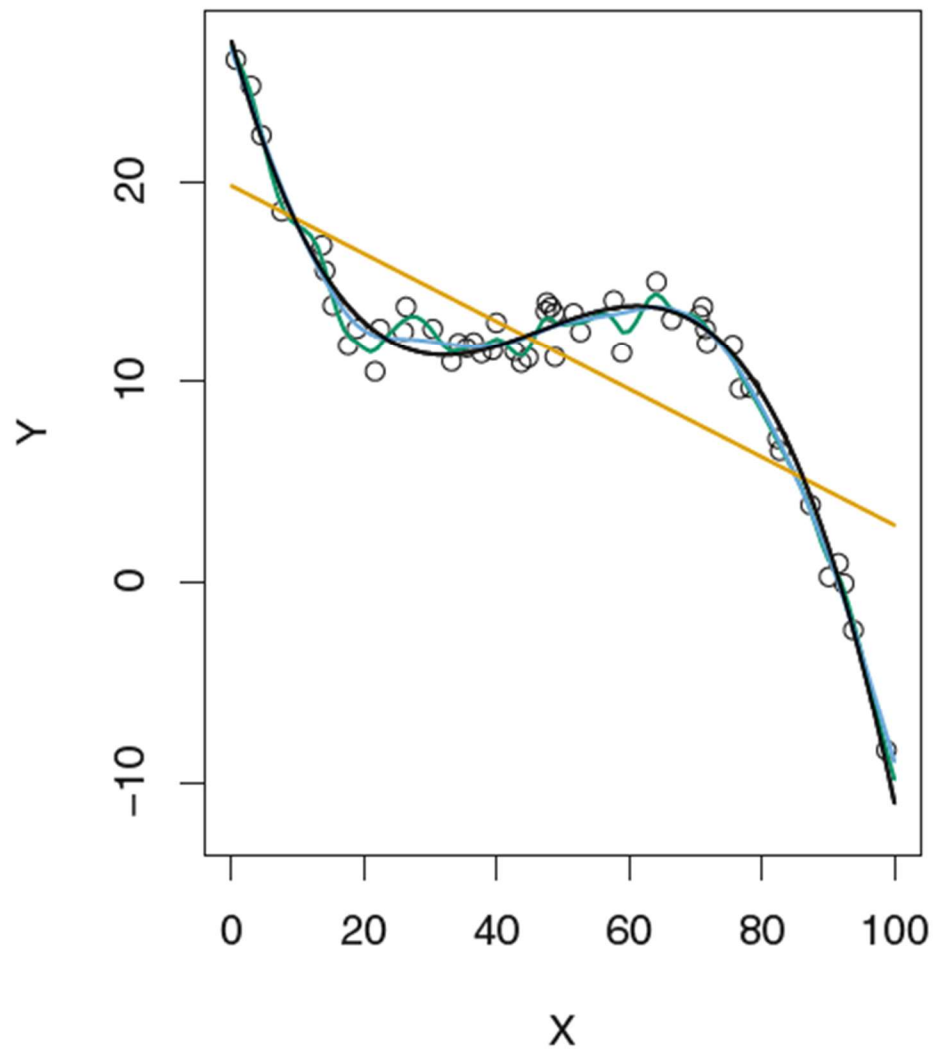


- In general, this MSE behavior (decrease and U-shape) holds regardless of the dataset and of the statistical method.
- As the flexibility increases, the training MSE will decrease, but the test MSE may not.
- When we have a small training MSE and a large test MSE -> **overfitting!**
 - The method may be picking patterns from the training data that are caused by random chance and those patterns don't actually exist in the test data -> increased test MSE .

Other examples



Other examples



The bias-variance trade-off

- Two **competing properties** of statistical learning approaches
- The test *MSE* for a given value x_0 can be decomposed into the sum of three fundamental measures:

$$\underbrace{E \left(y_0 - \hat{f}(x_0) \right)^2}_{\text{expected test } MSE} = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

The bias-variance trade-off

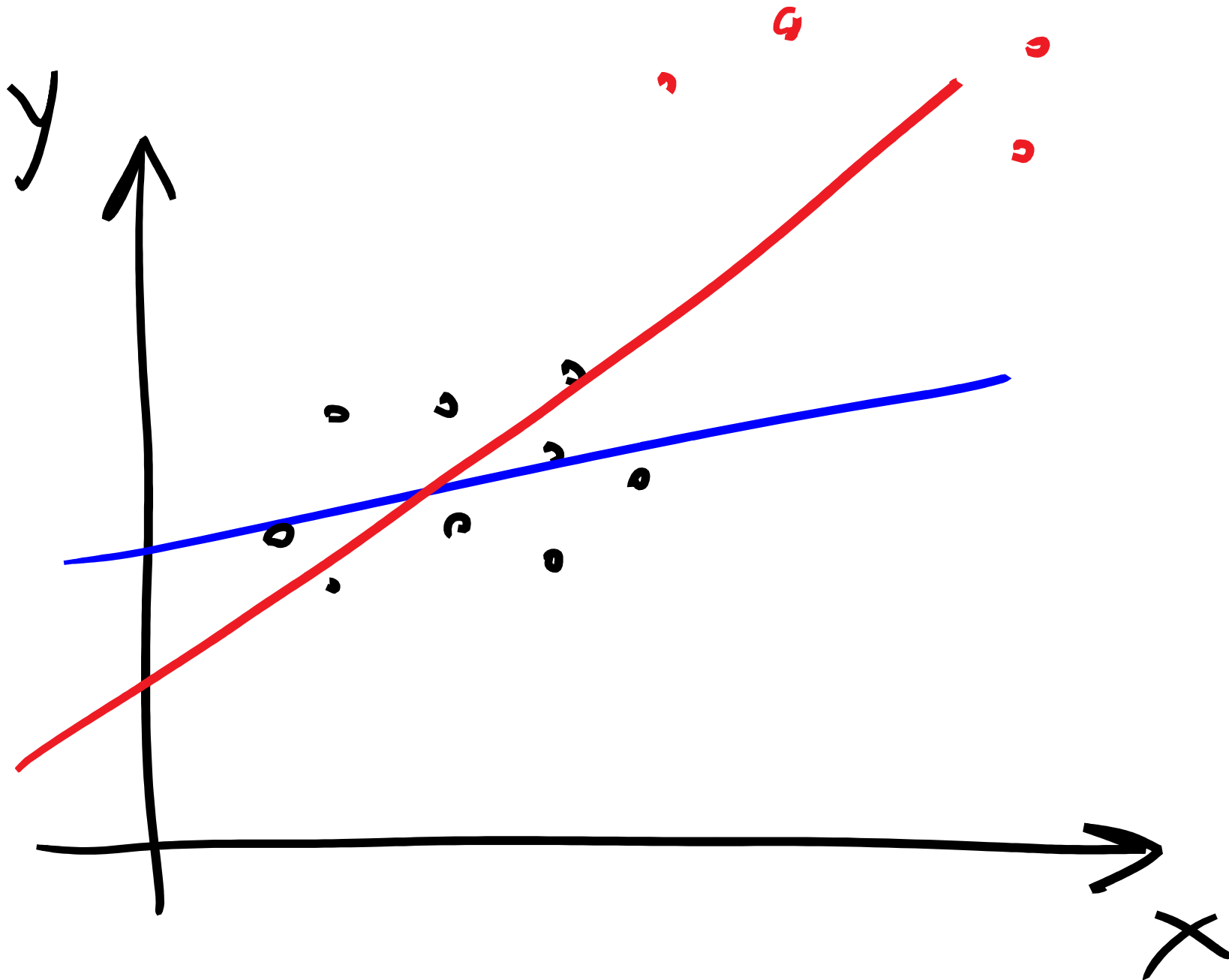
$$\underbrace{E \left(y_0 - \hat{f}(x_0) \right)^2}_{\text{expected test } \mathbf{MSE}} = \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

- The overall expected test **MSE** can be computed by averaging this measure for all possible values of x_0 in the test data.
- **What does this equation infer?**
 - In order to **minimize the expected test error**, we need to choose a method that would satisfy **low variance and low bias**.
 - Note that the **expected test MSE** cannot be less than the **irreducible error** given by $\text{Var}(\epsilon)$.

The bias-variance trade-off

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{amount of variation of } \hat{f} \text{ using different training datasets}} + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

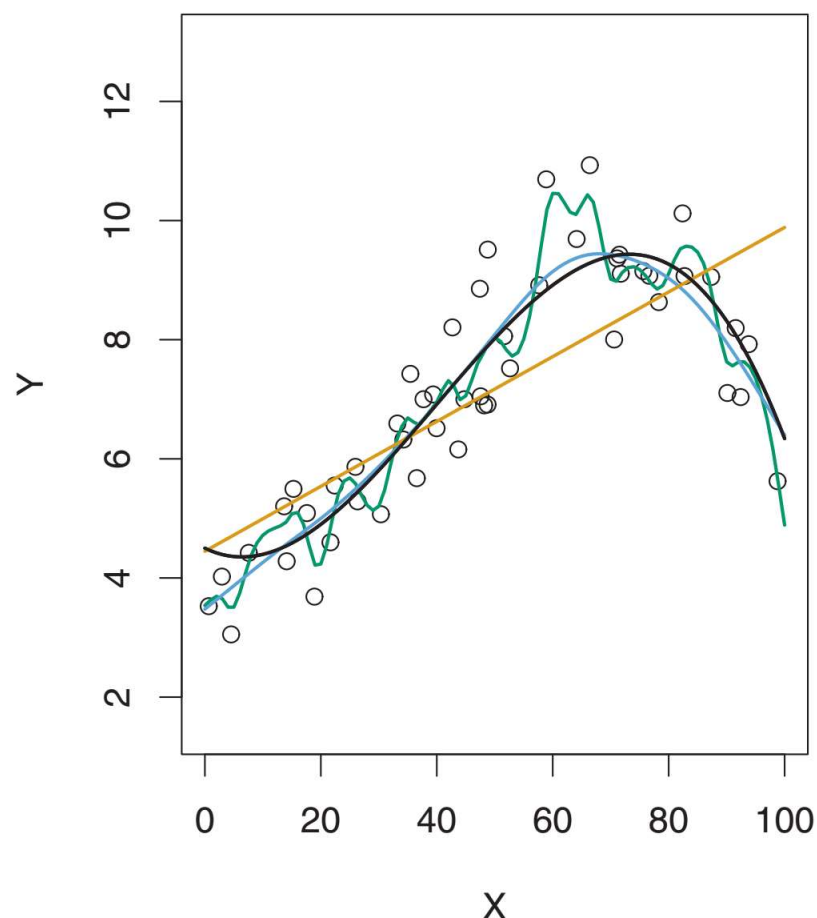
- Ideally the estimate of f **should not vary much** between different training datasets.
- A **high variance** means that **small data changes** -> **large changes in \hat{f}** .
- Generally, **more flexible** methods -> **higher variance**.
Why?



The bias-variance trade-off

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \underbrace{\text{Var}(\hat{f}(x_0))}_{\text{amount of variation of } \hat{f} \text{ using different training datasets}} + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\epsilon)$$

amount of variation
of \hat{f} using different
training datasets



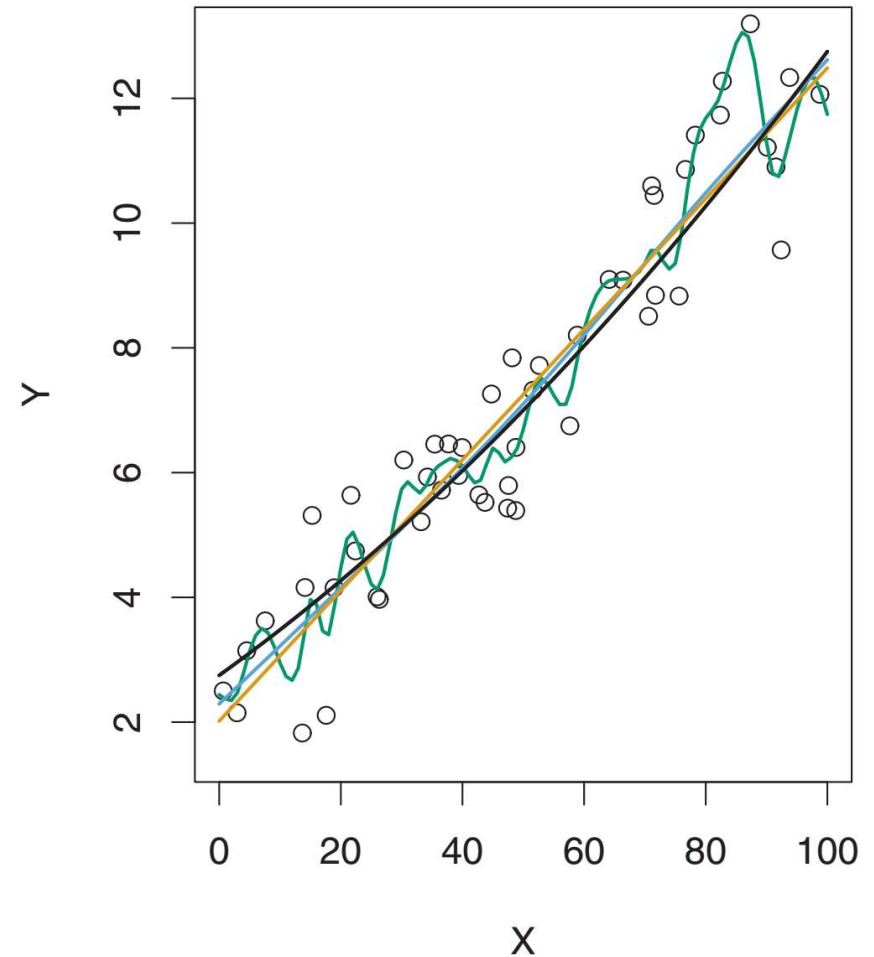
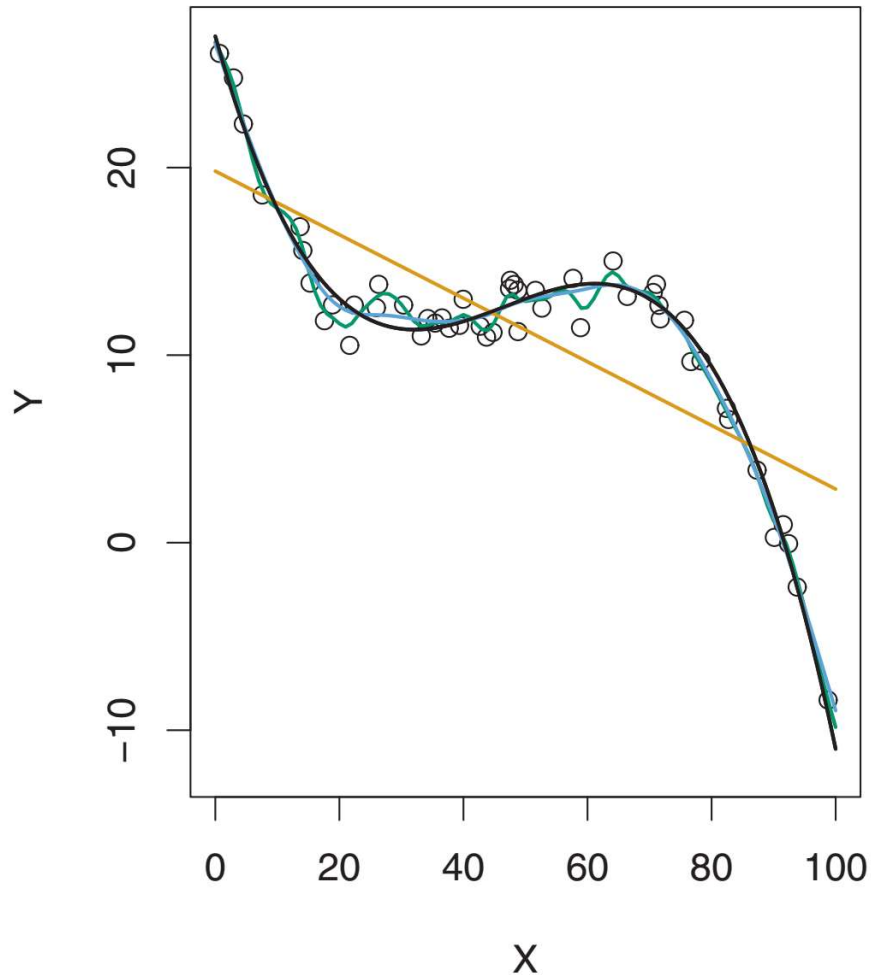
- The **green** curve has a very close trend as the data observations – it is very flexible.
 - Changing data -> higher variance.
- The **orange** curve is less flexible,
 - Moving a single observation would poorly affect the position of the line.

The bias-variance trade-off

$$E \left(y_0 - \hat{f}(x_0) \right)^2 = \text{Var}(\hat{f}(x_0)) + \underbrace{[\text{Bias}(\hat{f}(x_0))]^2}_{\text{error introduced by approximating a very complicated problem by a much simpler problem}} + \text{Var}(\epsilon)$$

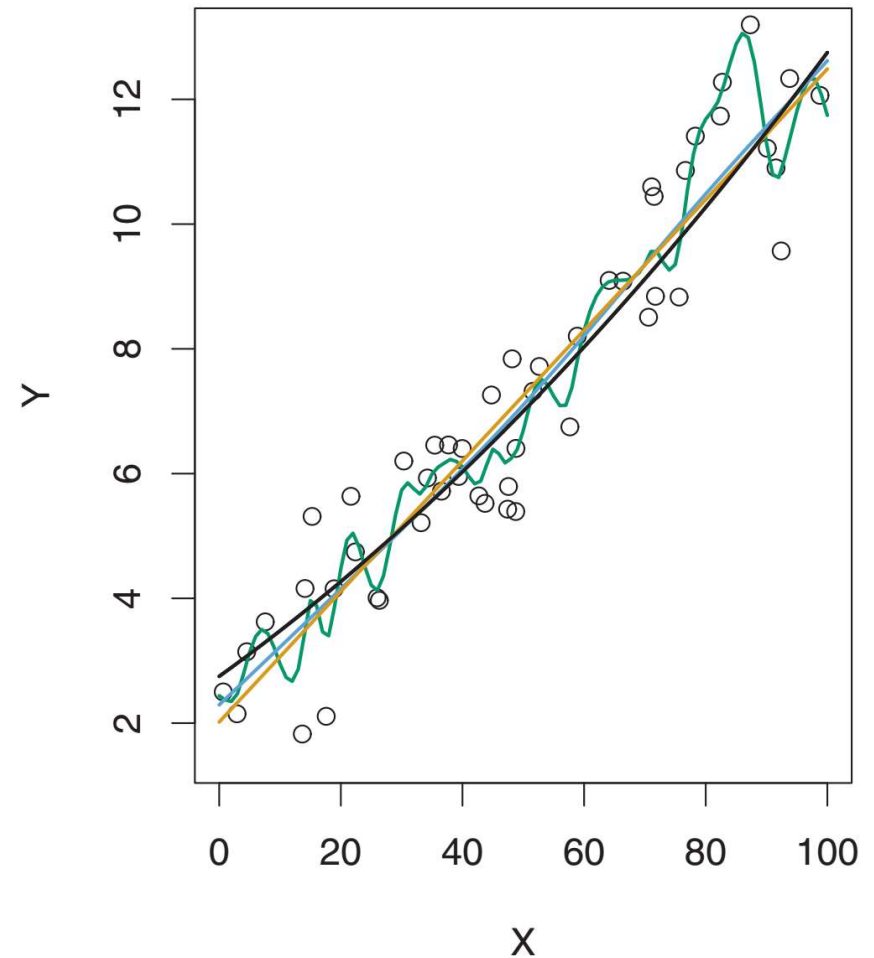
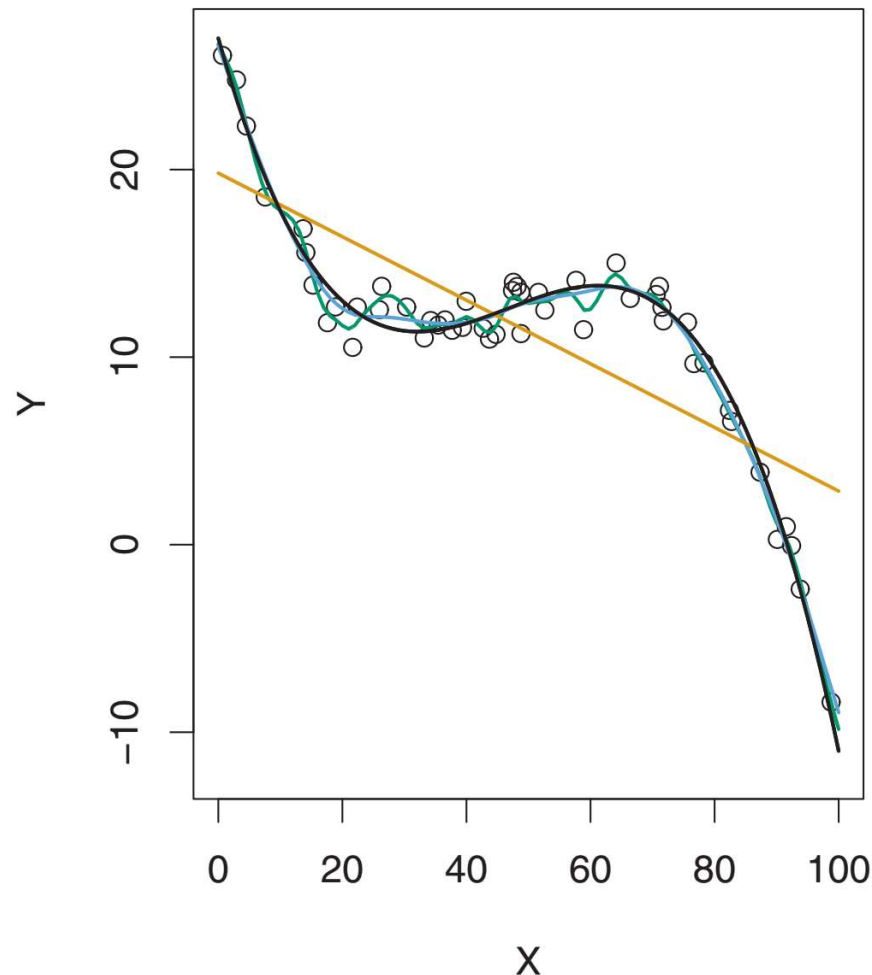
error introduced by approximating a very complicated problem by a much simpler problem

The bias-variance trade-off



- More flexible methods -> _____ (less or more) bias?

The bias-variance trade-off



- For instance, it is less likely that any real-life problem has a truly simple linear relationship between Y and X_1, X_2, \dots, X_p .
- Performing linear regression on such problems would certainly introduce some bias in \hat{f} .
- In general, for **more flexible methods** -> **less bias**.

The bias-variance trade-off – summary

- As a general rule, as flexibility increases, variance increases and bias decreases.
- The corresponding change of rates defines how MSE changes.
- The **bias usually decreases faster** than the variance increase.
 - Expected test MSE decreases.
 - But at some point, the flexibility does not affect the bias anymore, but it starts to significantly increase the variance.

Classification

- In this setting, y_i is not numerical.
- We want to estimate f based on training observations

$$\{(x_1, y_1), \dots, (x_n, y_n)\}$$

where y_1, \dots, y_n are **qualitative**.

Classification

- In this setting, y_i is **not numerical**.
- We want to estimate f based on training observations

$$\{(x_1, y_1), \dots, (x_n, y_n)\}$$

where y_1, \dots, y_n are **qualitative**.

- The most common way to measure the **accuracy of \hat{f}** is the **error rate**.
- For training data, it is the **proportion of mistakes that are made** if we apply \hat{f} to training data.

training error rate: $\frac{1}{n} \sum_{i=1}^n$

$$I(y_i \neq \hat{y}_i)$$

equals 1 if $y_i = \hat{y}_i$,
and 0 otherwise.

the predicted class label of
the i^{th} observation using \hat{f}

Classification

training error rate: $\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i)$

equals 1 if $y_i = \hat{y}_i$,
and 0 otherwise.

the predicted class label of
the i^{th} observation using \hat{f}

- We are interested in the test error rather than the training error.

testing error rate: $\text{Ave} (I(y_0 \neq \hat{y}_0))$

corresponding to the average error
rate on unseen observations (x_0, y_0)

Bayes Theorem

$$\Pr(Y = \text{Default} | x_0) = 0.15$$

$$\Pr(Y = \text{Not Def.} | x_0) = 0.85$$

$$\Pr(Y = c_i | X = x_0) = \frac{\Pr(x_0 | c_i) \Pr(c_i)}{\Pr(x_0)}$$

Posterior prob.

$c_1 \equiv \text{Default}$

$c_2 \equiv \text{Not Default}$

$$x_0 = (x_1, x_2)$$
$$= (32, 10)$$

age / years of education

Bayes classifier

- A very simple classifier minimizes this test error on average.
- That is achieved by assigning each observation to the most likely class, given its predictor values.
- A test observation with predictor vector x_0 is assigned to class j for which the following probability is the largest:

$$\Pr(Y = j | X = x_0)$$

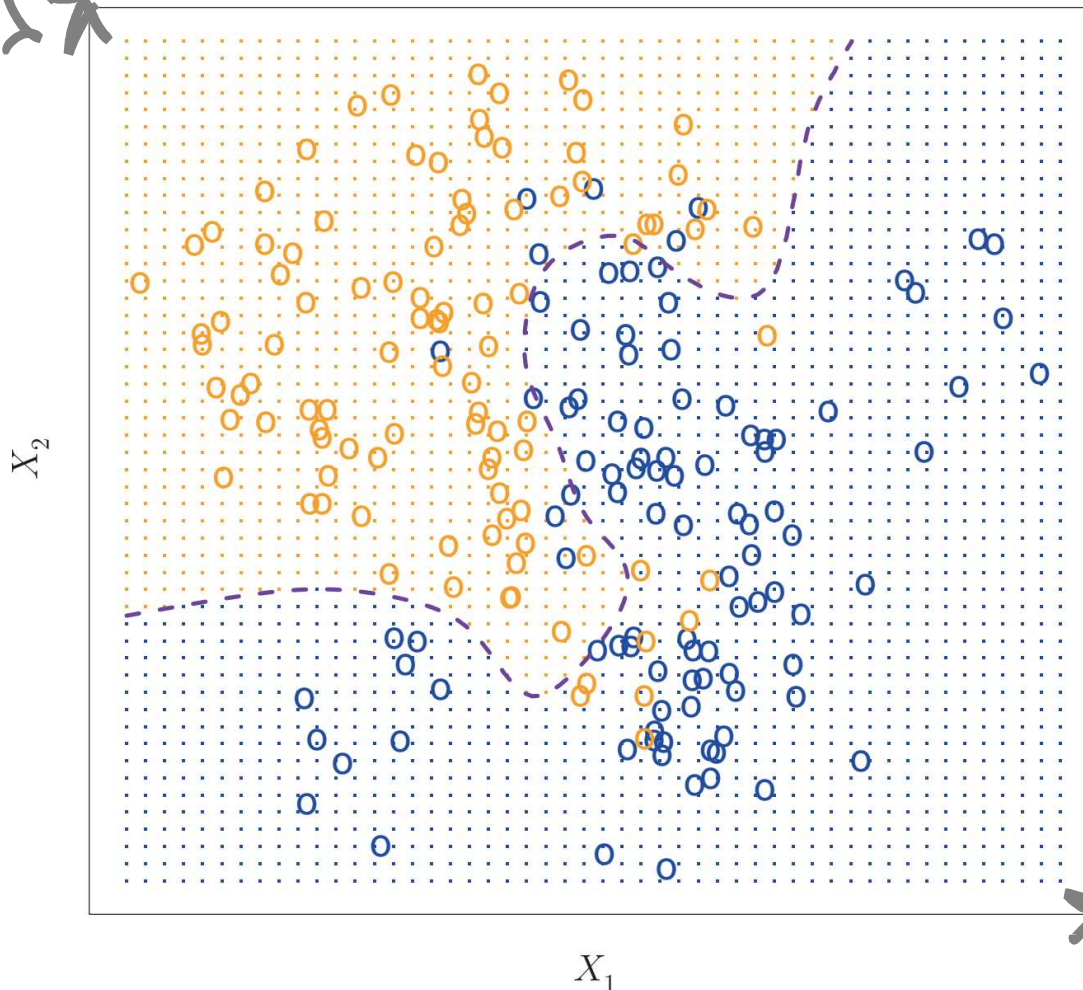
conditional probability
that $Y = j$, given x_0

Bayes classifier – two classes

- Suppose we have a two-class (A or B) problem, the Bayes classifier **assigns x_0 to A if:**

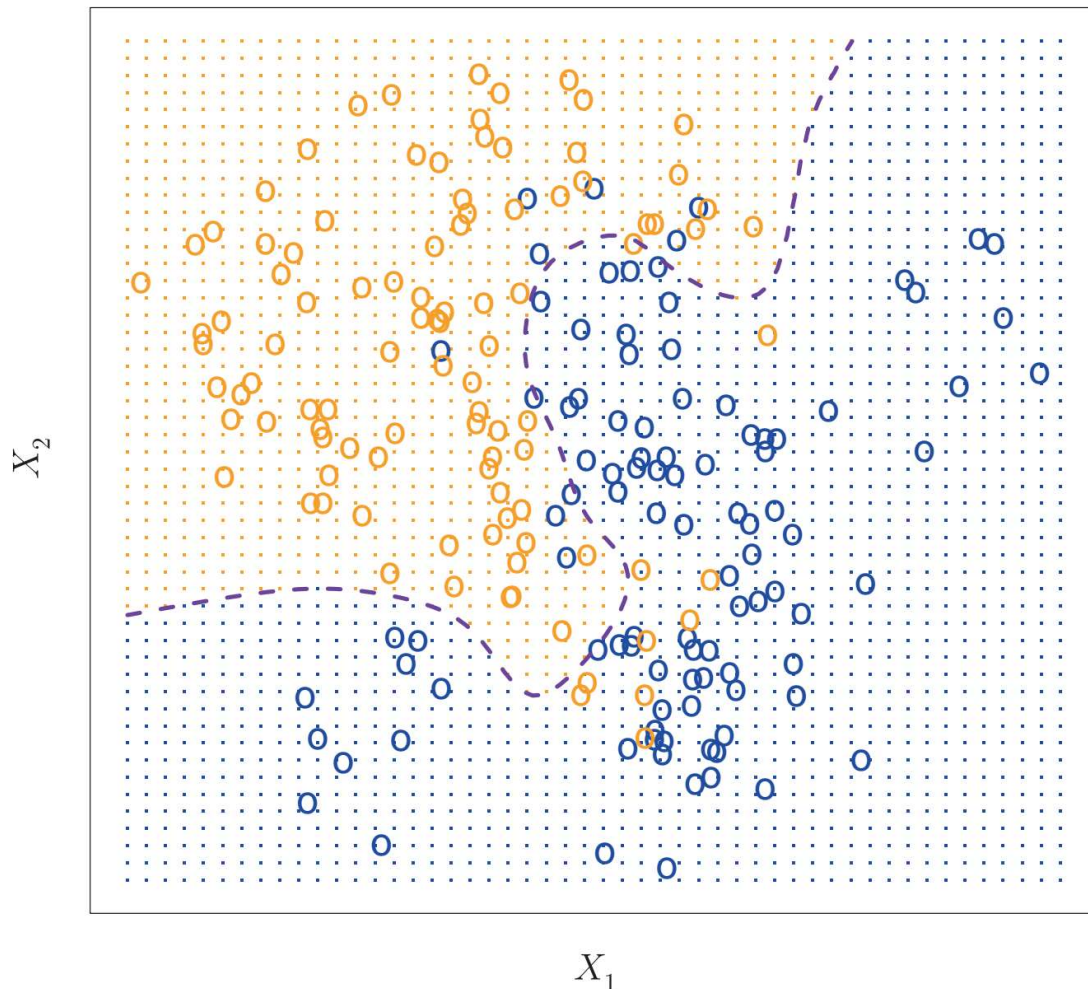
$\Pr(Y = A|X = x_0) > 0.5$ and class B otherwise

years
of education



- Simulated dataset in $2D$ space of predictors X_1 and X_2 .
- Two classes represented by **orange** and **blue** circles.
- We can see that for each value of X_1 and X_2 , there is a different probability of the response (**orange** or **blue**).

Bayes classifier – two classes



- The **orange region** reflects the set of point for which $\Pr(Y = \textit{orange} \mid X) > 0.5$.
- The **blue region** indicates that this probability is less than 0.5.
- The dashed line separating both regions is called the **Bayes decision boundary** and it is where this probability is **exactly 0.5**.

Bayes classifier – two classes

- The Bayes classifier produces the **lowest possible test error** (Bayes error rate).
 - always chooses class for which the probability is largest.

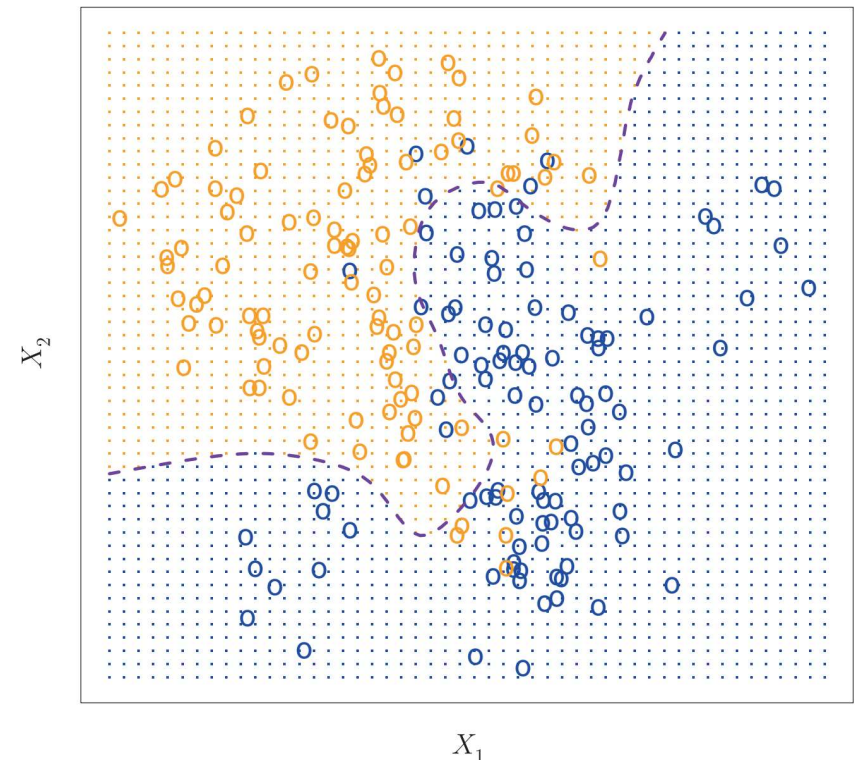
- The error rate at $X = x_0$ is:

$$1 - \max_j \Pr(Y = j | X = x_0)$$

- The overall Bayes error rate is:

$$1 - E(\max_j \Pr(Y = j | X))$$

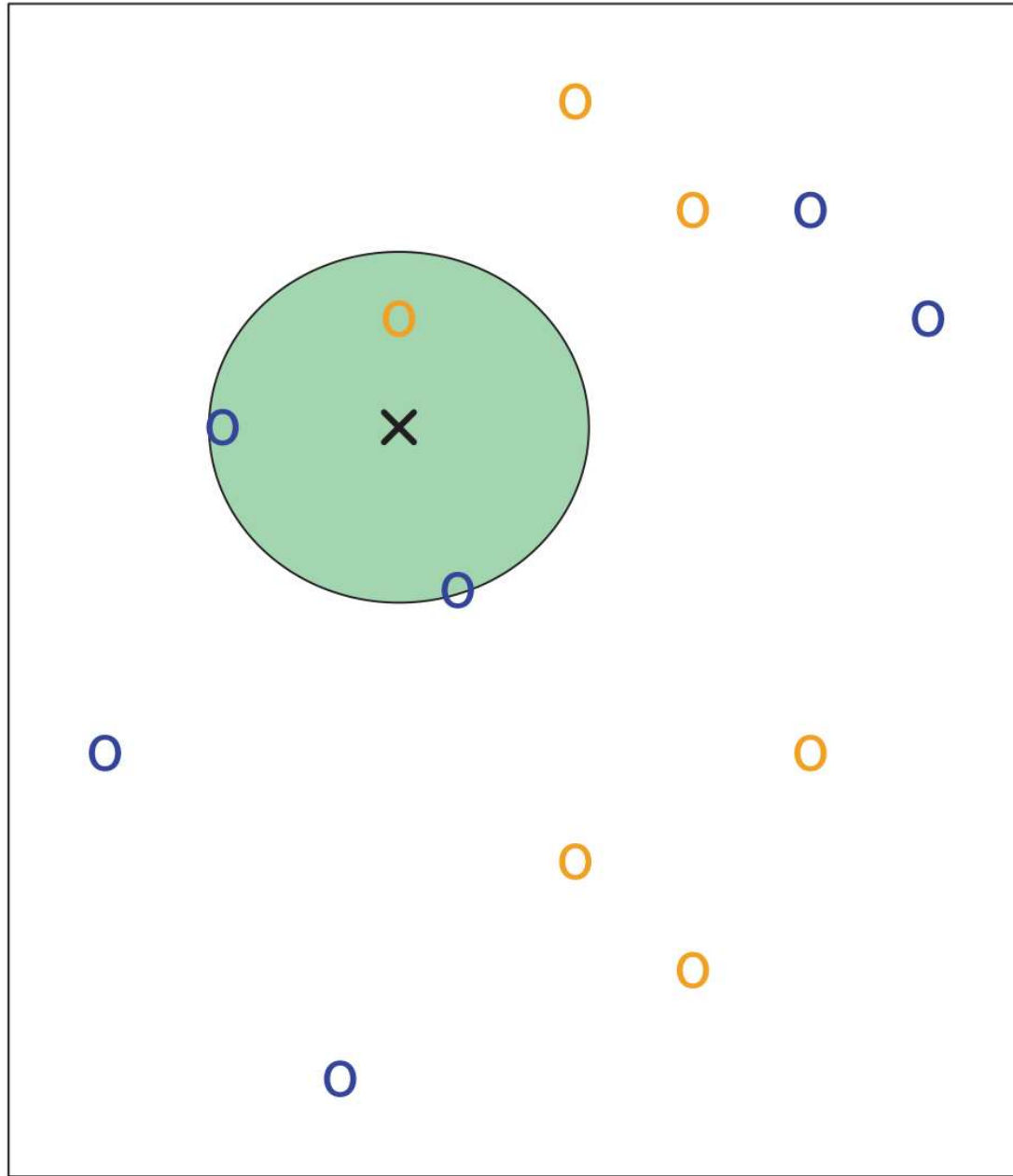
Equivalent to the irreducible error.
It is the lowest possible classification
error that can be reached.



K-nearest neighbors classifier

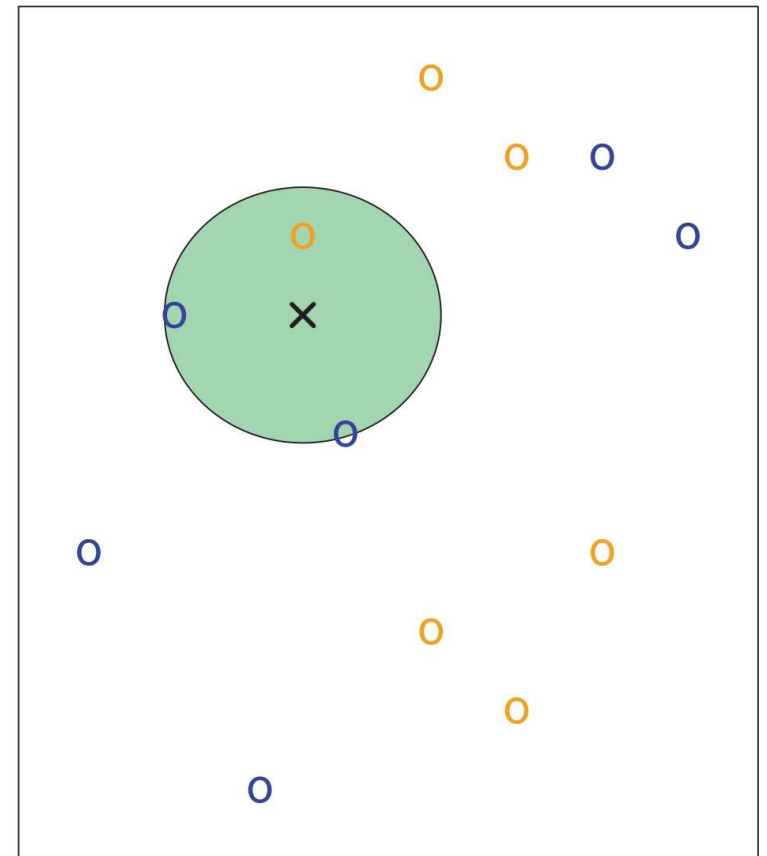
- We do not always know the conditional distribution of Y given X -> using Bayes classifier is not possible.
- In a way, the **Bayes classifier serves as an unattainable gold standard** to which we can compare the performance of other approaches.
- The K -nearest neighbors (**KNN**) classifier thus seeks:
 - to estimate the conditional distribution of Y given X
 - then to classify an observation to the class with the highest estimated probability.

3-nearest neighbors classifier



K-nearest neighbors classifier

- Small data consisting of **6 blue** and **6 orange** observations.
- Goal: make a **prediction** for the point represented by a **black x**.
- Suppose that $K = 3$,
 1. the classifier identifies the **three observations that are the closest** to x .
 2. the nearest neighbors result in **estimated probabilities of $2/3$ for the blue class** and $1/3$ orange class.
 3. KNN will predict that x belongs to the **blue class**.



K-nearest neighbors classifier

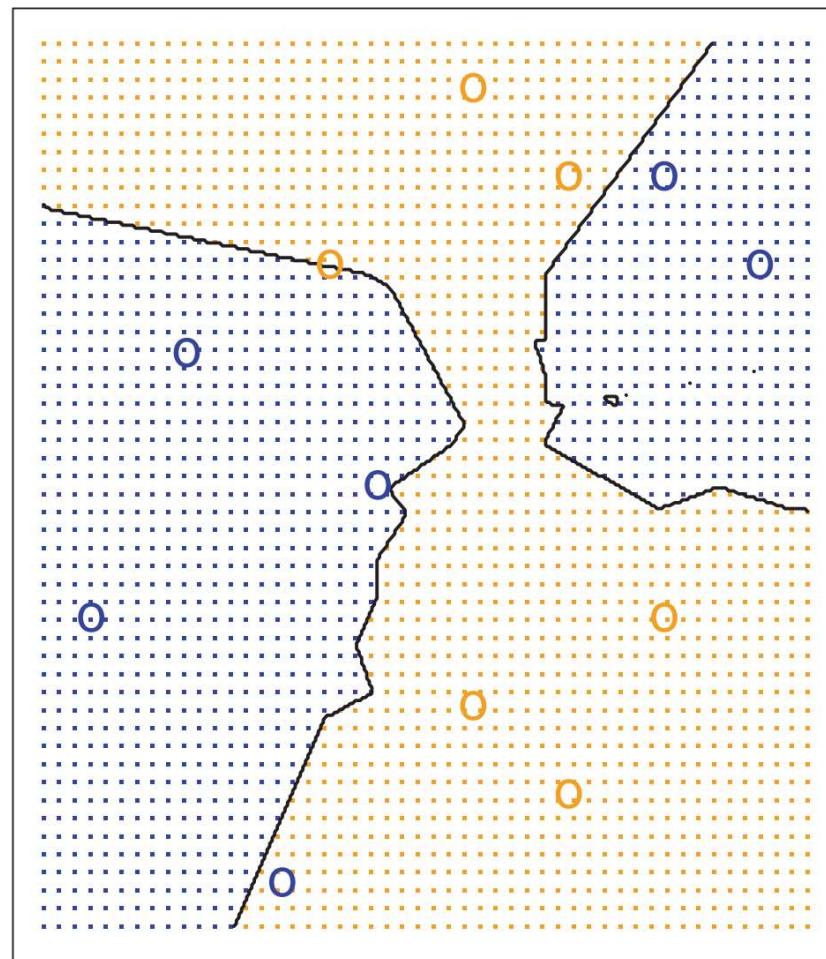
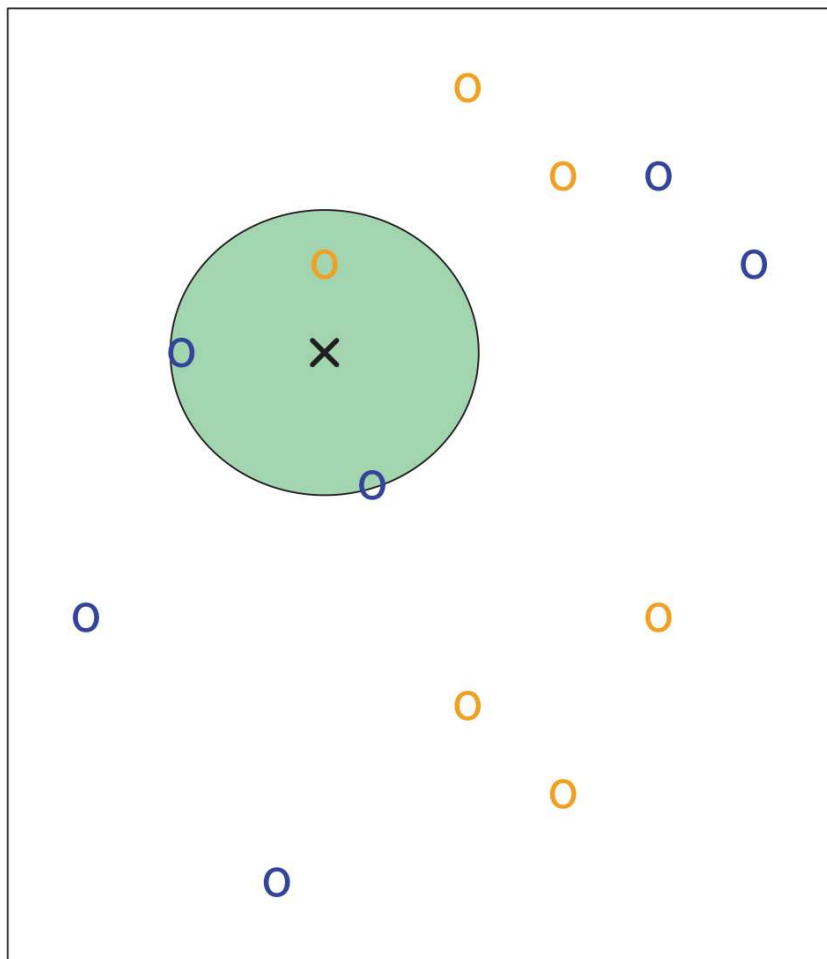
Given a positive integer K and a test observation x_0 :

1. The classifier **identifies** \mathcal{N}_0 consisting of **K points in the training data that are the closest to x_0 .**
2. Then, it **estimates the conditional probability for class j** as the fraction of points in \mathcal{N}_0 with response values equal to j :

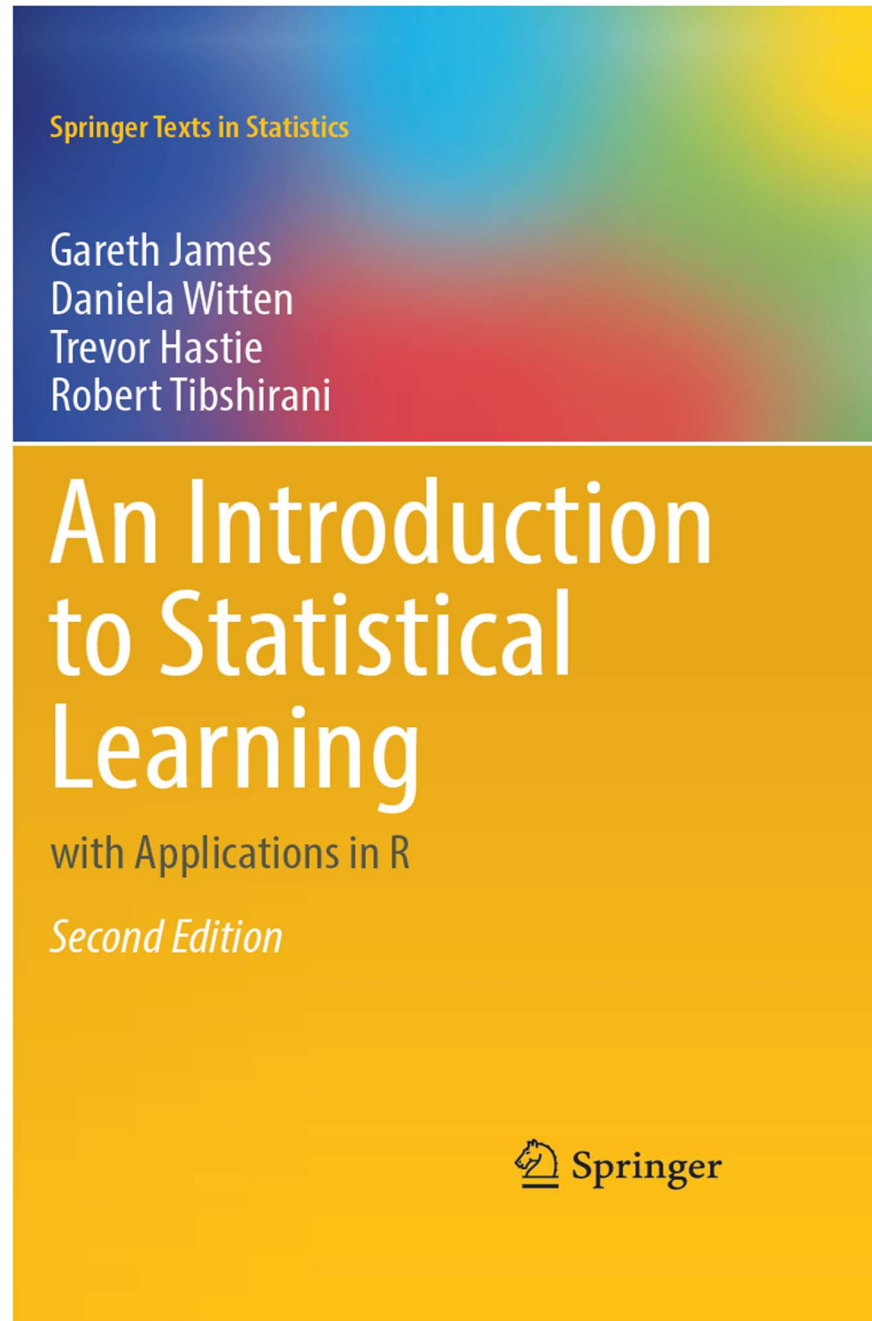
$$\Pr(Y = j | X = x_0) = \frac{1}{K} \sum_{i \in \mathcal{N}_0} I(y_i = j)$$

3. It applies **Bayes rule** and classifies x_0 to the class with the **largest probability**.

K-nearest neighbors classifier



Reference



Springer Texts in Statistics

Gareth James
Daniela Witten
Trevor Hastie
Robert Tibshirani

An Introduction to Statistical Learning

with Applications in R

Second Edition

 Springer