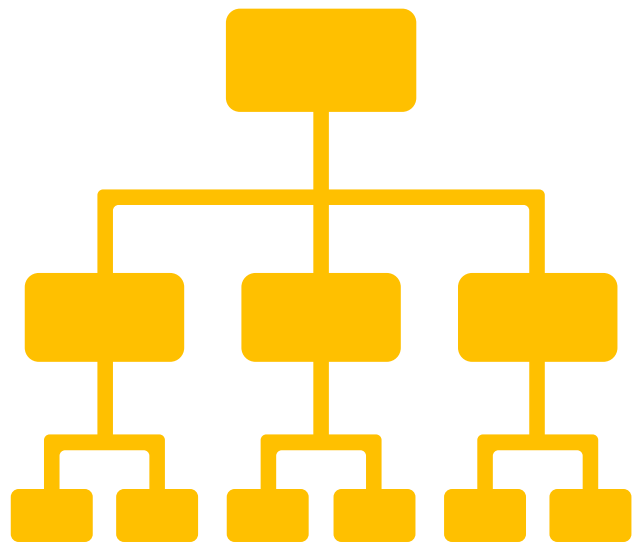# Data
# Divers

## Data Science – Yellow Belt

## Lecture 2

by Sebastian Sauer

# Statistical Learning: What's that?

# Big Picture

Lecture 1

Lecture 2
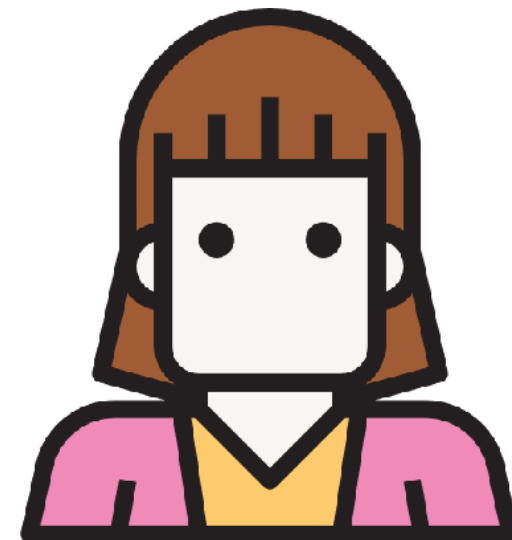
Lecture 3

# Learning goals
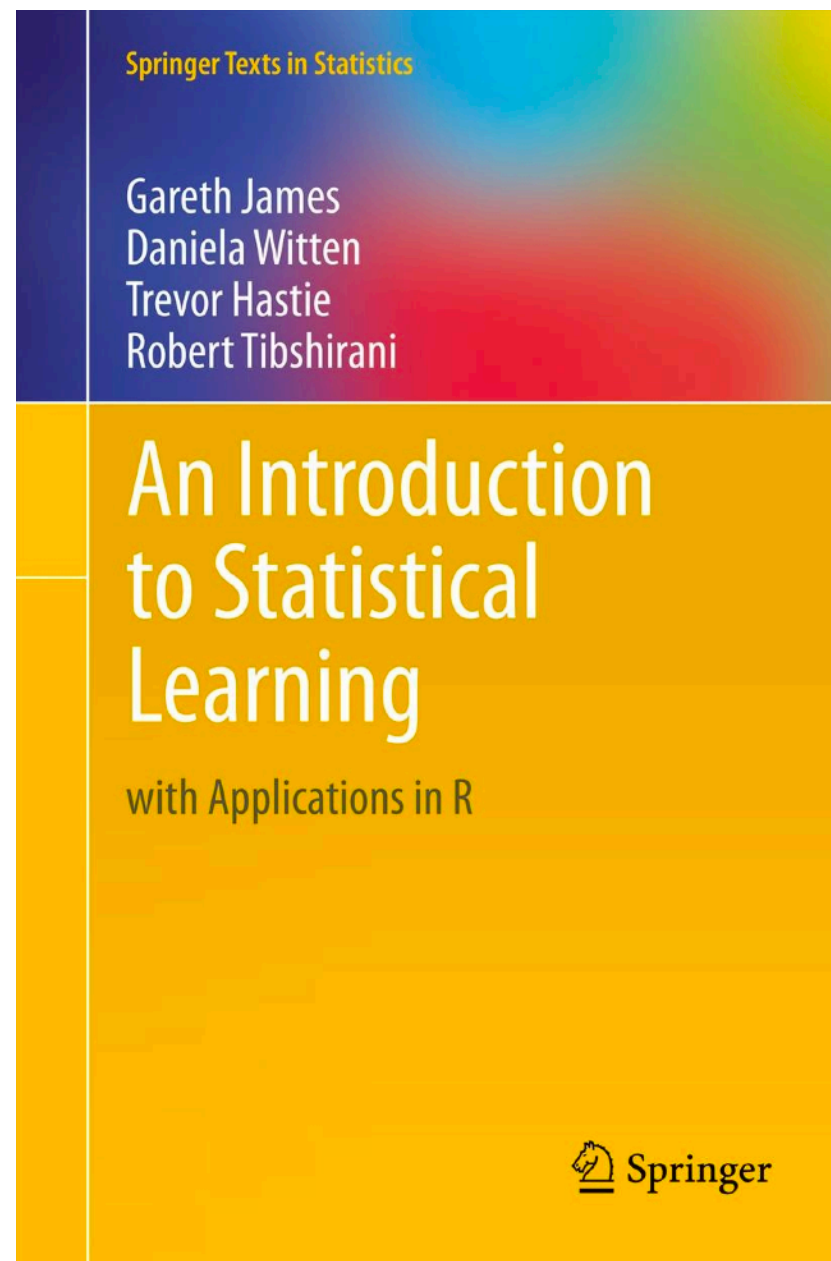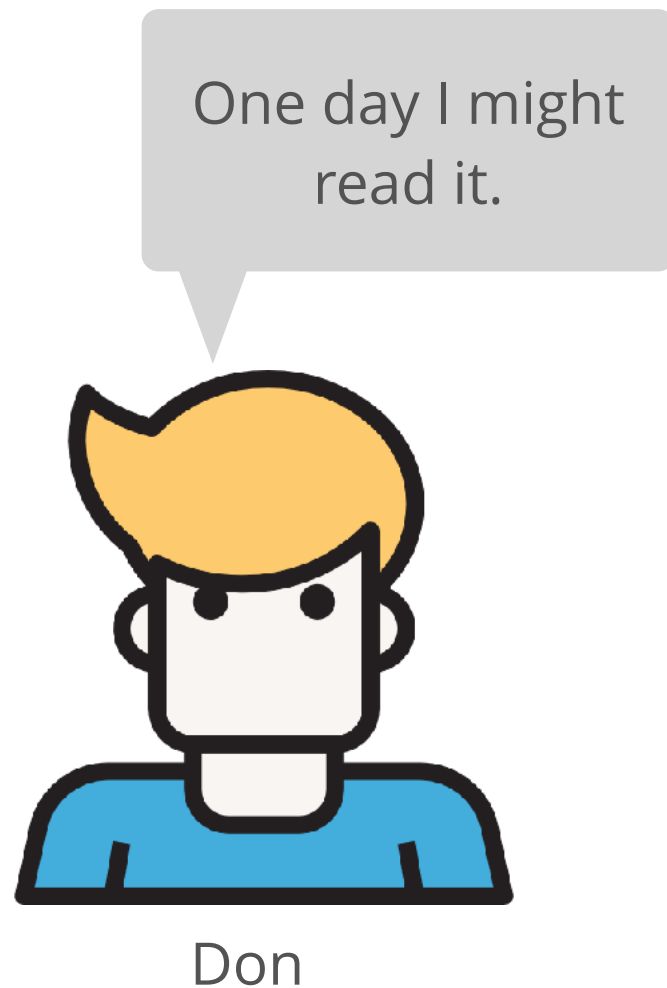
Let's explore the mathematical intricacies!

No, this time we'll get our hands dirty with state-of-the-art case-studies.

OK, but I'll throw in some background where and when needed.

Wolfi

Angi

# The standard source of knowledge

One day I might read it.

Don

**Springer Texts in Statistics**

Gareth James
Daniela Witten
Trevor Hastie
Robert Tibshirani

An Introduction
to Statistical
Learning

with Applications in R

Springer

ebook freely available

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

▶ In break-out groups, work your way through this demonstration of a machine learning pipeline.

▶ Answer the following questions:

   1. What are model parameters?

   2. Which models tend to exhibit a strong bias?

   3. Why do overly-complex exhibit high variance?

   4. What's a way to balance the bias-variance tradeoff?

▶ Feel free to double check part 1 of this demonstration.

https://rstudio.cloud/

https://www.tidymodels.org/

# Here's our tutorial/case study

Stuck? Confused? Ask for help.

# 1 Build a model

**TIDYMODELS PACKAGES:** broom, parsnip

- INTRODUCTION
- THE SEA URCHINS DATA
- BUILD AND FIT A MODEL
- USE A MODEL TO PREDICT
- MODEL WITH A DIFFERENT ENGINE
- WHY DOES IT WORK THAT WAY?
- SESSION INFORMATION

## INTRODUCTION 🔗

How do you create a statistical model using tidymodels? In this article, we will walk you through the steps. We start with data for modeling, learn how to specify and train models with different engines using the parsnip package, and understand why these functions are designed this way.
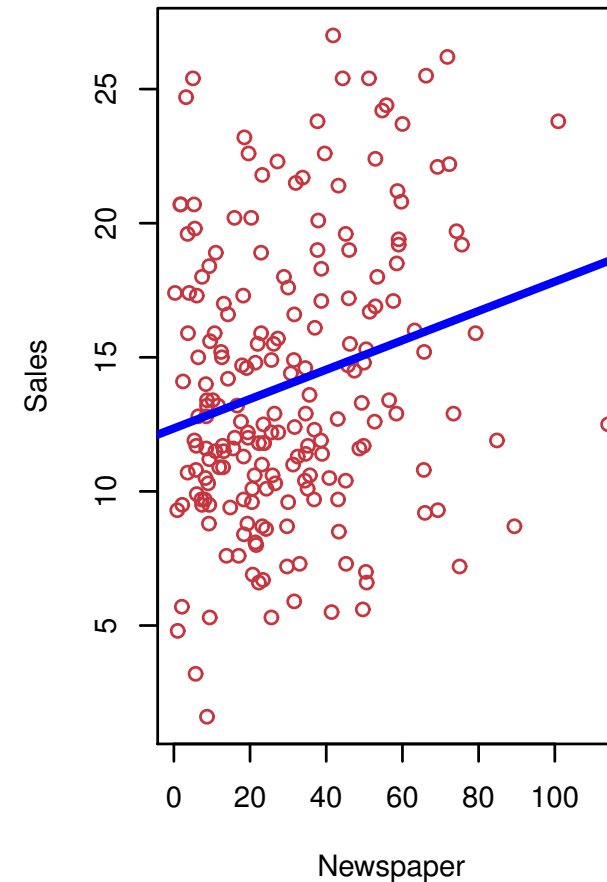
To use code in this article, you will need to install the following packages: broom.mixed, readr, rstanarm, and tidymodels.

https://www.tidymodels.org/start/models/

# Statistical Learning: Finding the pattern of Y and X

$$Y = f(X) + e$$

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Finding patterns is our game

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Why estimating $f$ ?

**Prediction**

**Explanation**

Who cares about „why"
as long you get accurate
predictions!

We need to
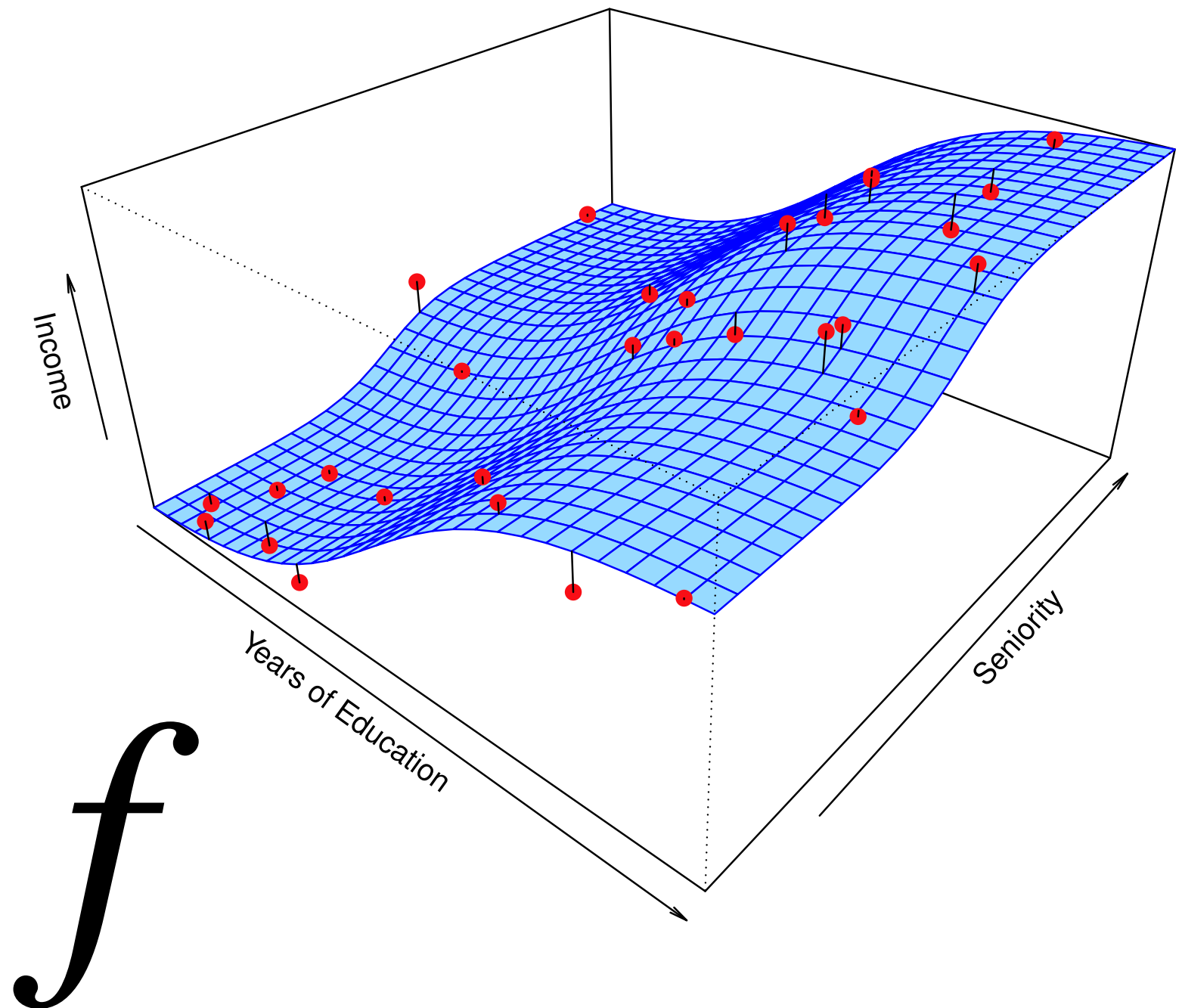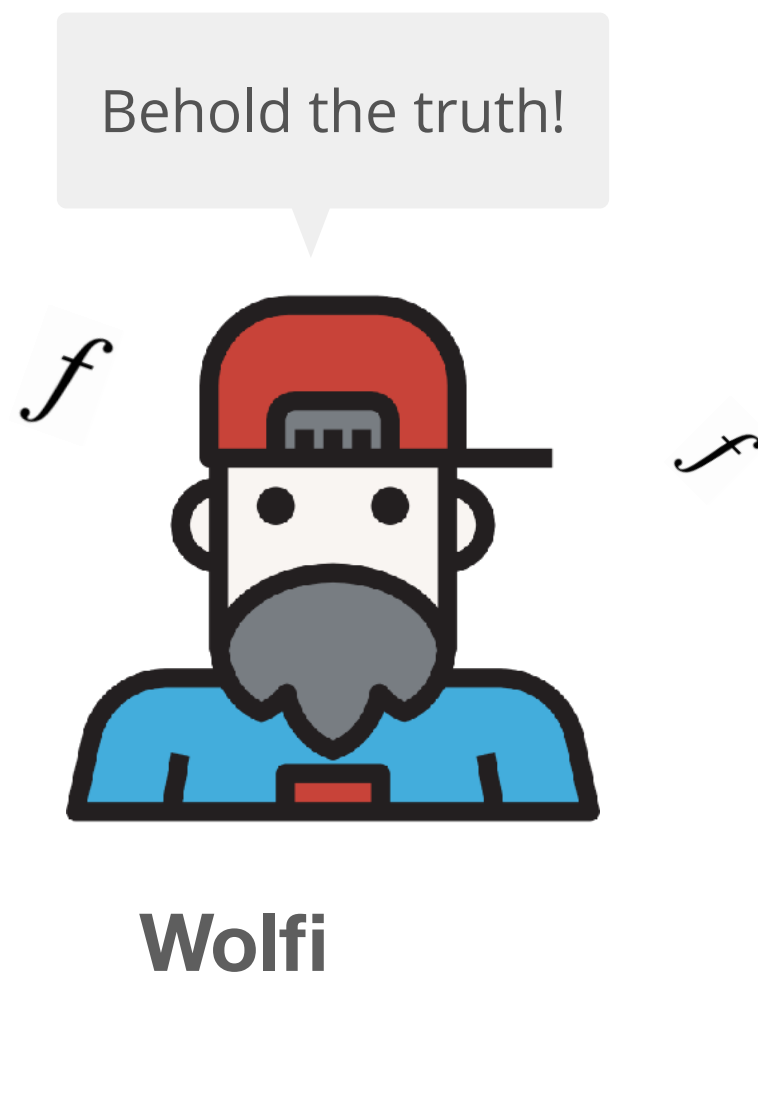understand what's
going on.

Don

Wolfi

# Two types of errors: reducible and non-reducible
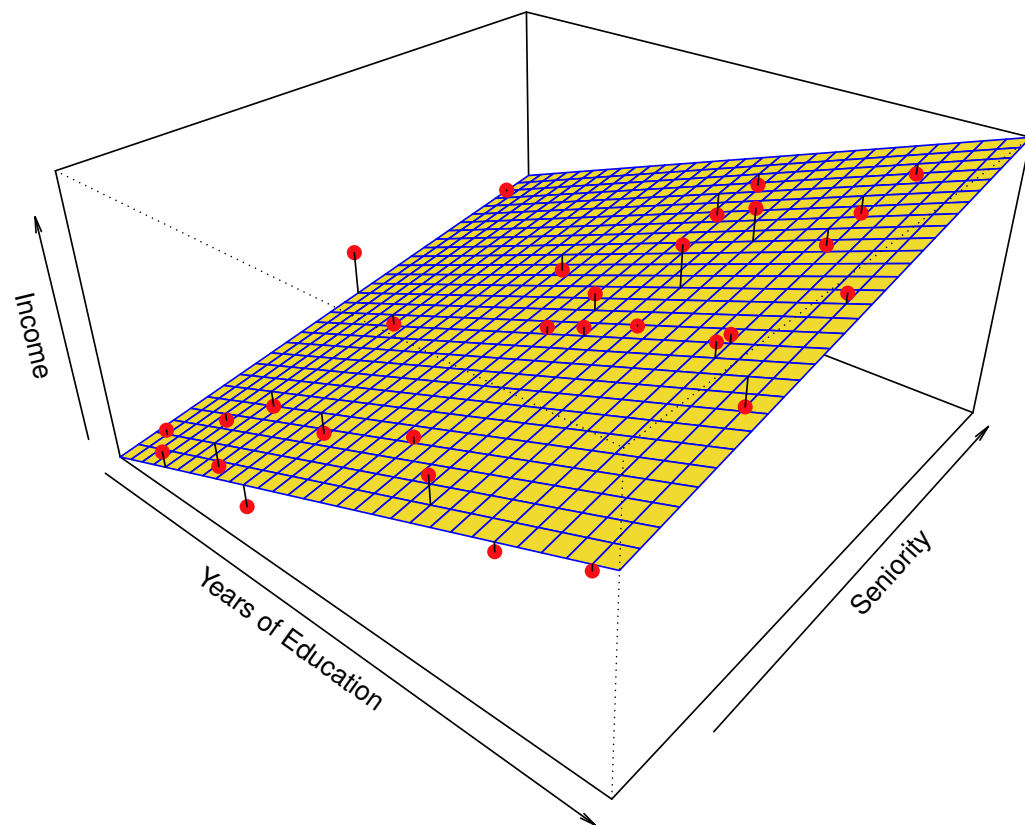
$$e = e_r + e_{nr}$$

$$E(Y - \hat{Y})^2 = E[f(X) + \epsilon - \hat{f}(X)]^2 = \underbrace{E[f(X) - \hat{f}(X)]^2}_{\text{reducible}} + \underbrace{Var(\epsilon)}_{\text{non-reducible}}$$

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Here's the non-reducible error

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf
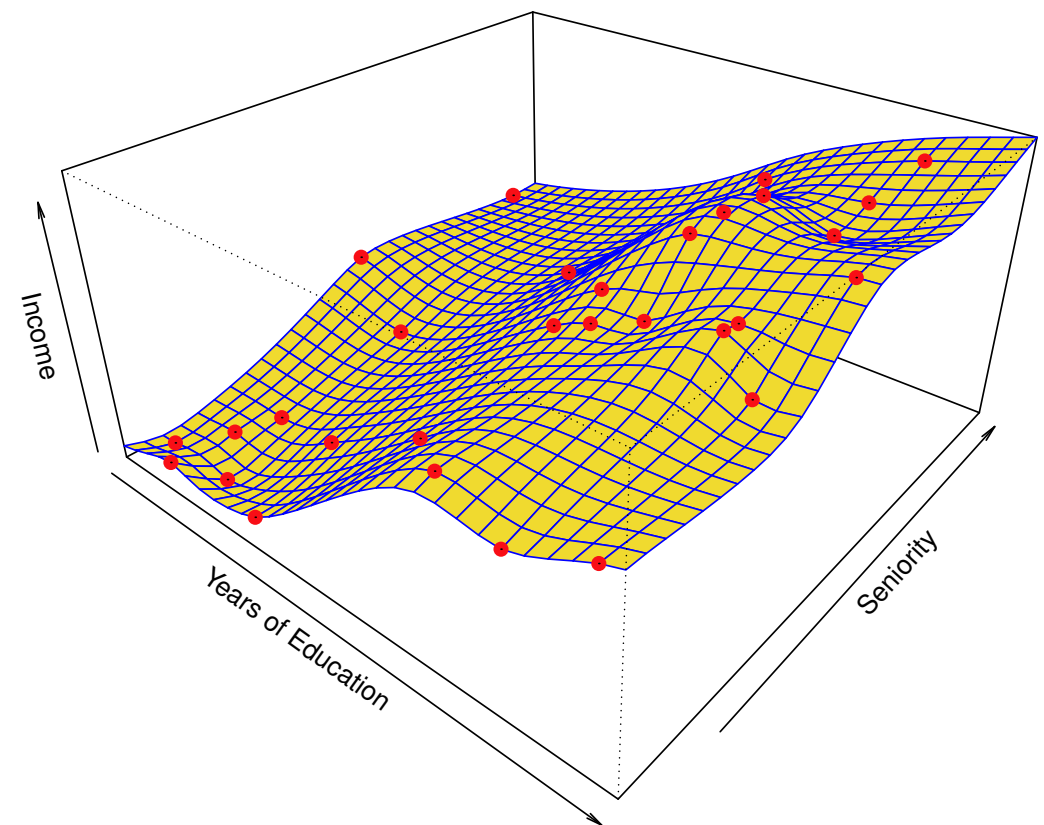
# Parametric and non-parametric models



Parametric model

Non-parametric Model

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf
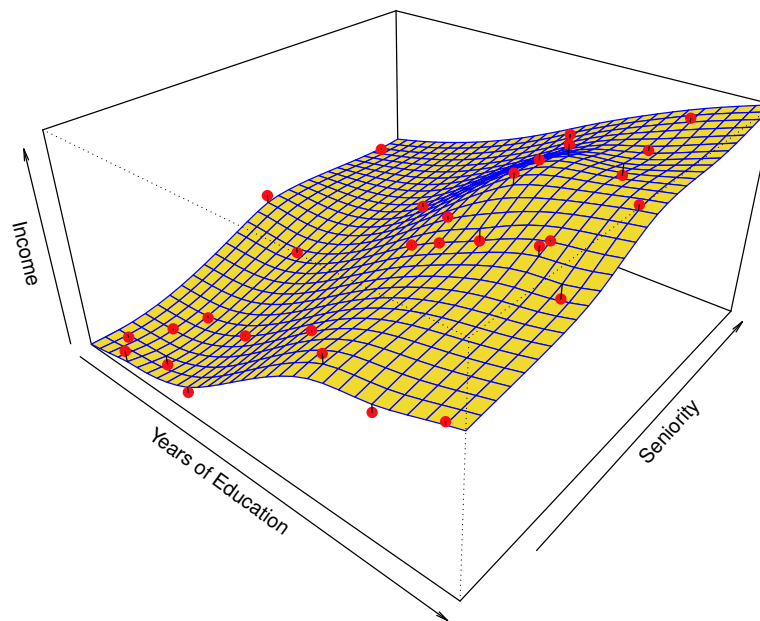
# From too simple to too complex

$f$



$\hat{f}_1$

$\hat{f}_2$

$\hat{f}_3$

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Interpretability vs. flexibility



Interpretability (y-axis: Low to High) vs. Flexibility (x-axis: Low to High)

- Subset Selection
- Lasso
- Least Squares
- Generalized Additive Models
- Trees
- Bagging, Boosting
- Support Vector Machines

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Supervised vs. unsupervised learning



supervised learning

unsupervised learning

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Typical machine learning workflow

Lang, M., Schratz, P., Binder, M., Pfisterer, F., Richter, J., Reich, N. G., & Bischl, B. (2020). *Mlr3 book*. https://mlr3book.mlr-org.com

# Learn and predict on new data



Lang, M., Schratz, P., Binder, M., Pfisterer, F., Richter, J., Reich, N. G., & Bischl, B. (2020). *Mlr3 book*. https://mlr3book.mlr-org.com
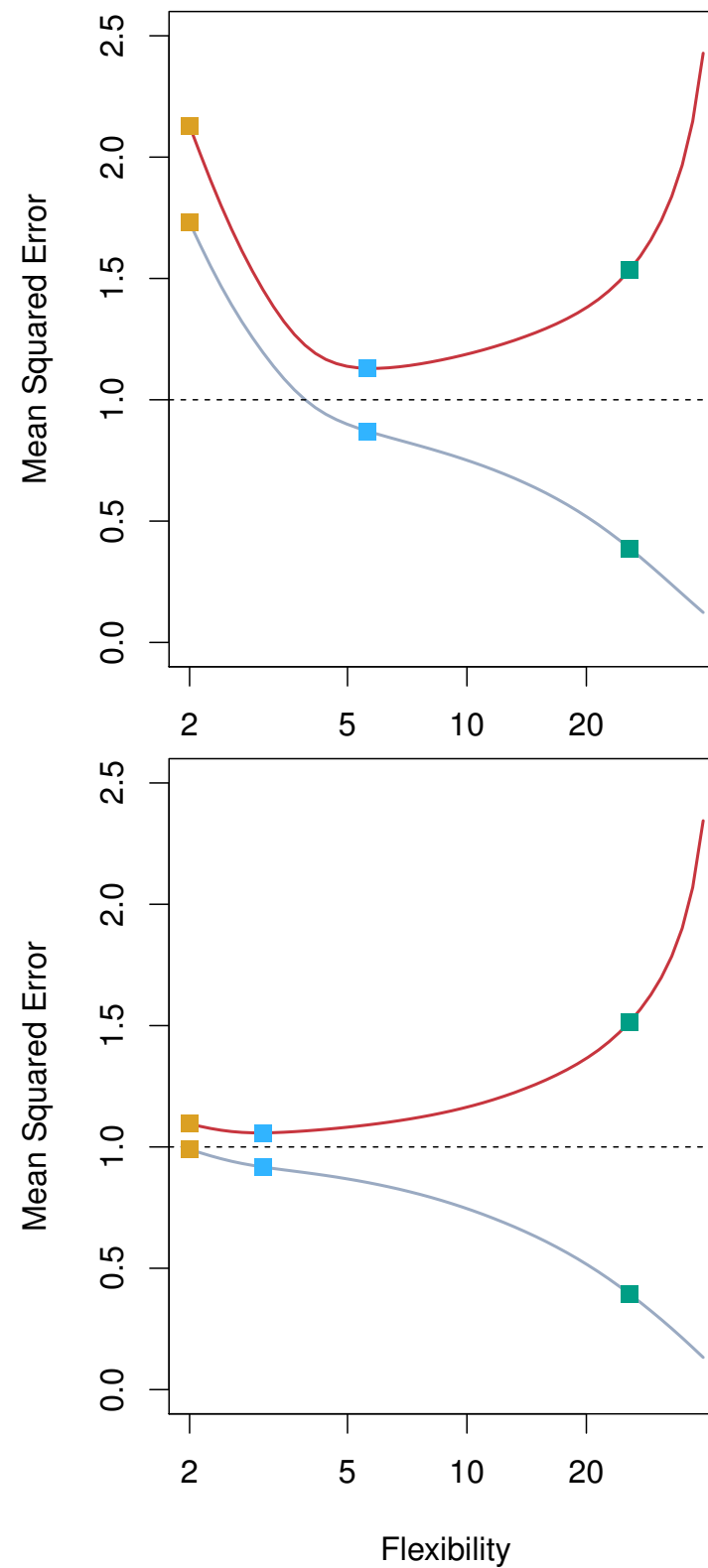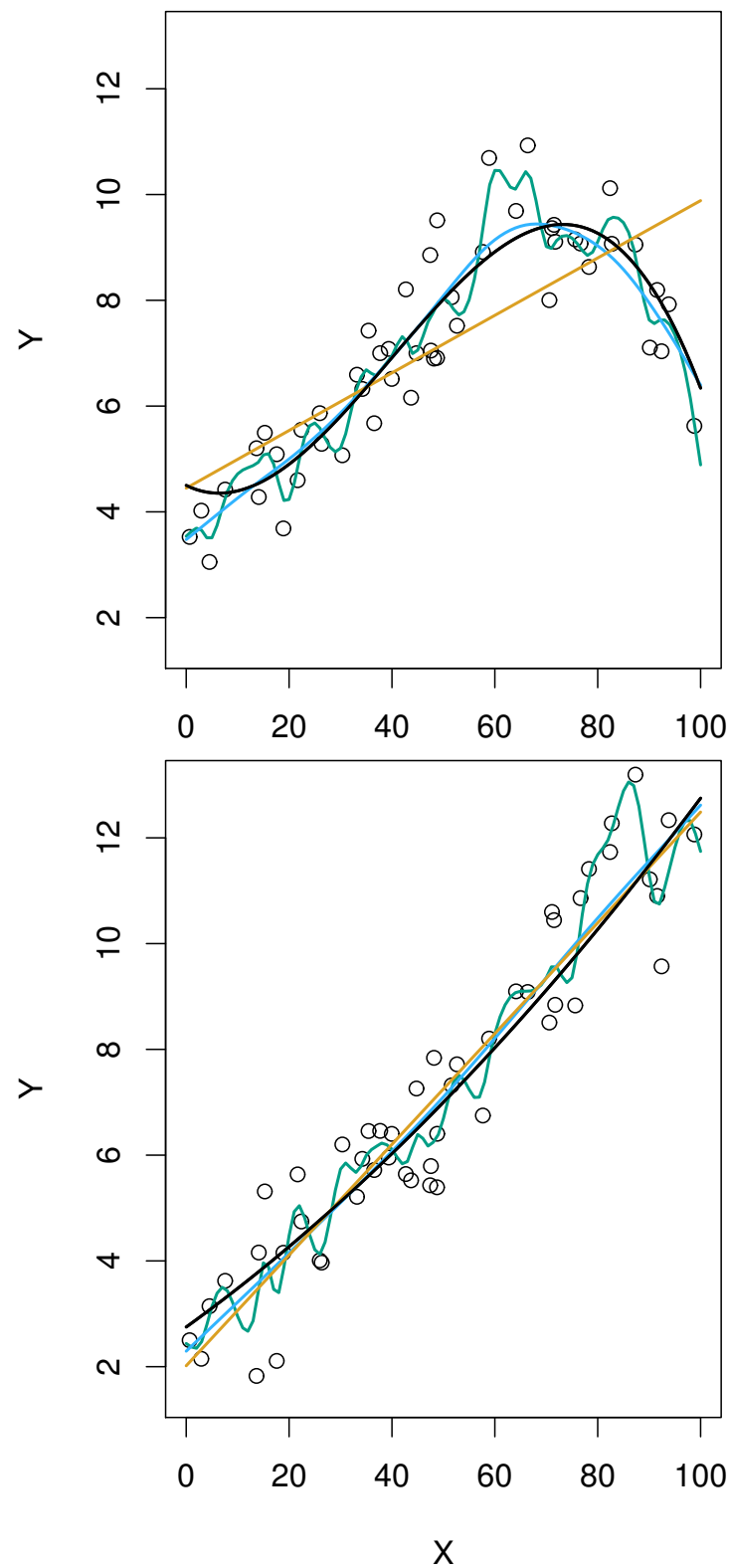
# Model evaluation

# Mean Square Error (MSE)

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

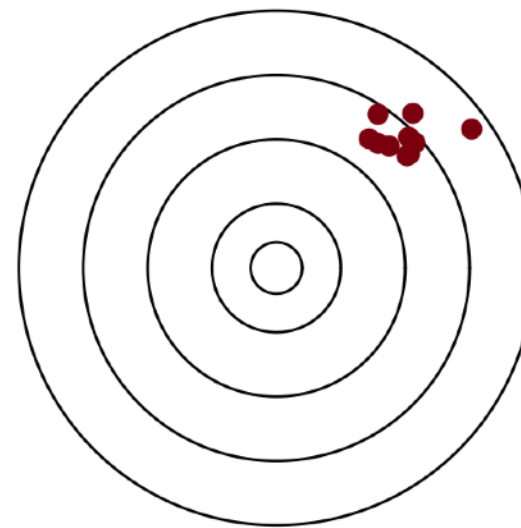$$MSE_{Test} = \frac{1}{n} \sum (y_0 - \hat{x}_o))^2$$

# MSE Train ≠ MSE Test

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Variance vs. bias



variance

bias

# MSE = Bias + Variance

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf
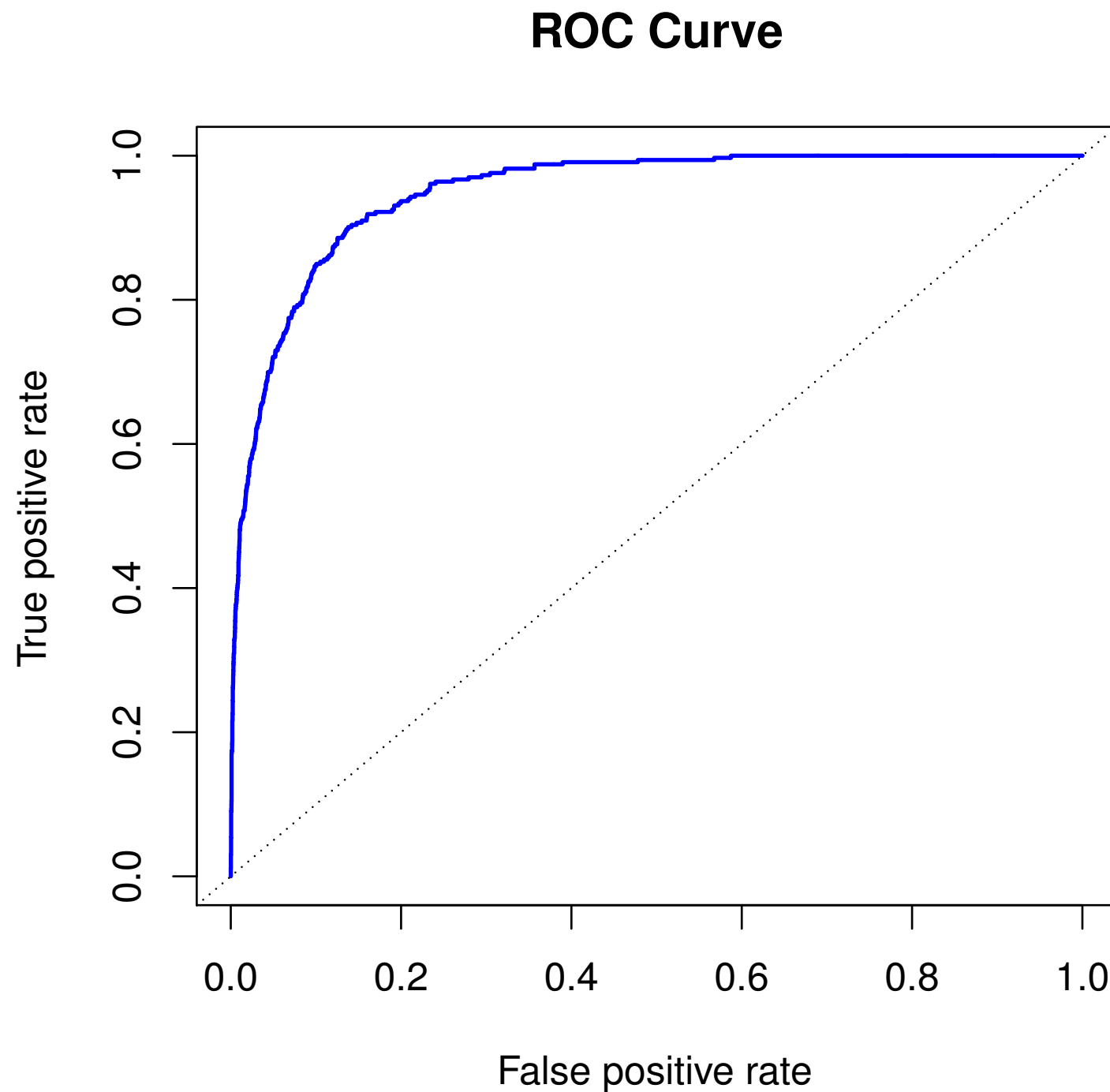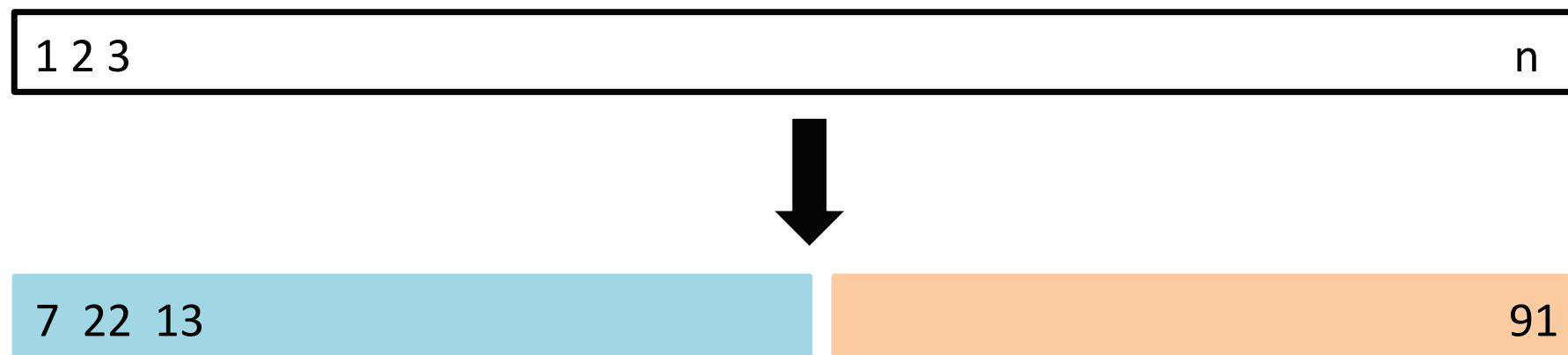
# Model evaluation in classification models

$$e = \frac{1}{n} \sum_{i=1}^{n} I(y_i \neq \hat{y}_i)$$

# Klassifikationsmodelle visuell beurteilen

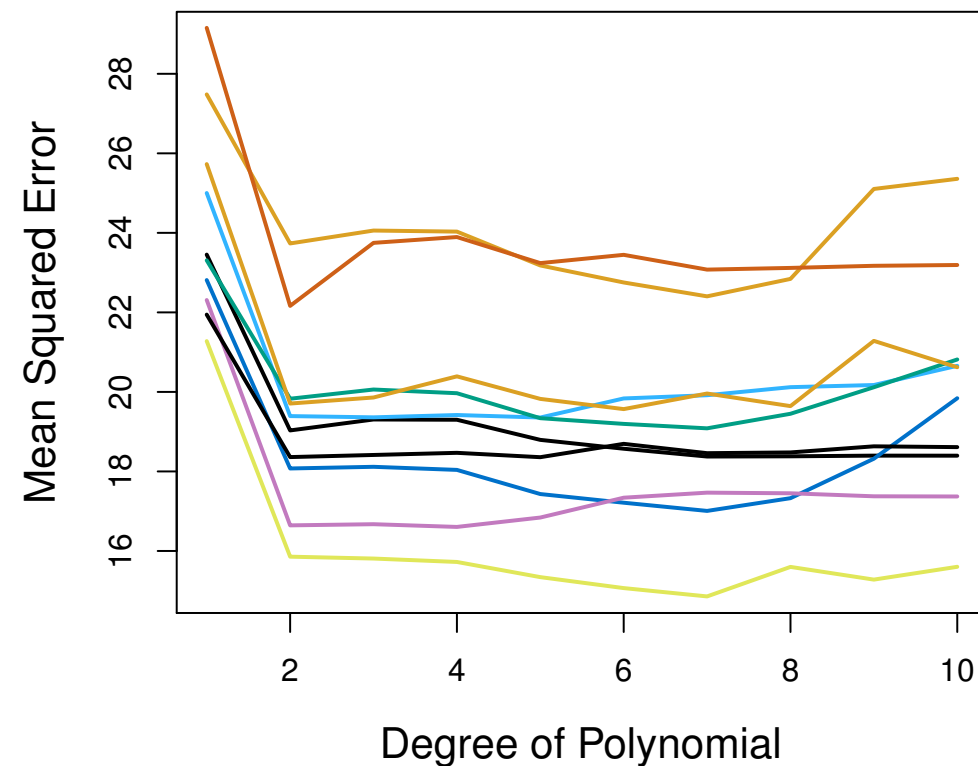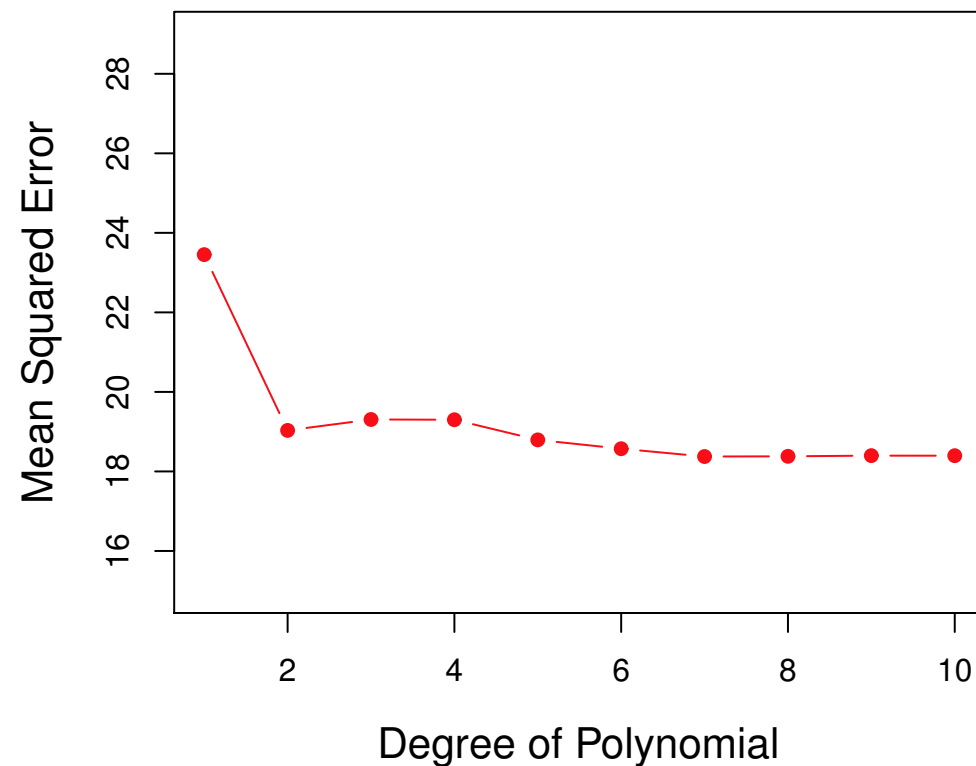

ROC Curve

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Resampling

# Train-Test split

```
1 2 3                                                                n
```

↓

```
7  22  13                                              91
```

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Train-Test split yields highly variable Test MSE



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf
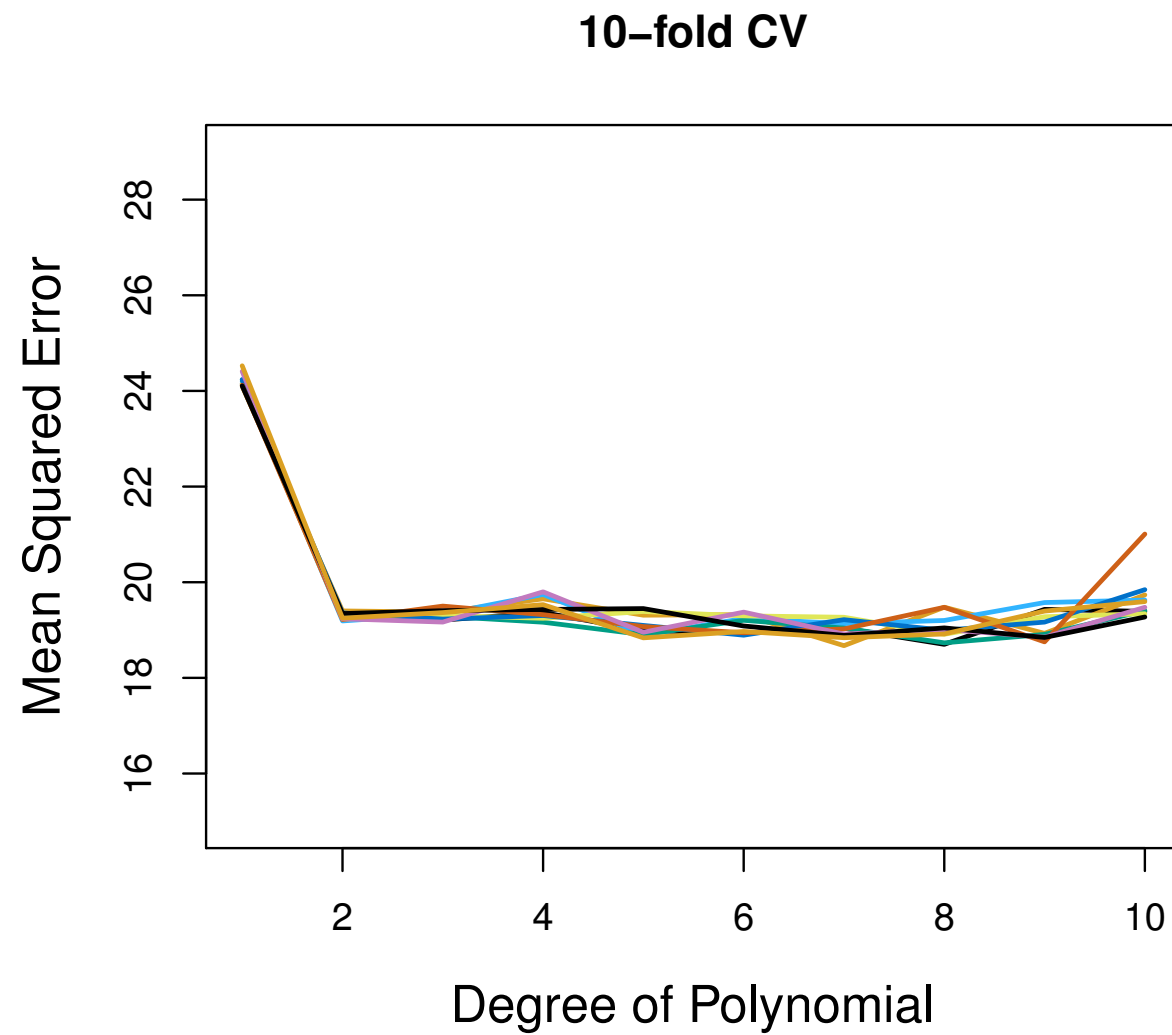
# k-fold cross-validation (k-fold cc)



$$CV_{(k)} = \frac{1}{k} \sum_{i=1}^{k} MSE_i$$

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# k-fold cross-validation (k-fold cc)



10–fold CV

# Bootstrap



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Tree-based Methods
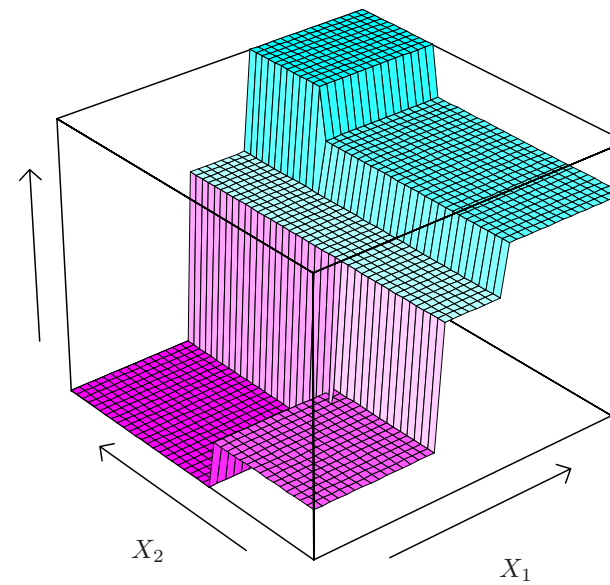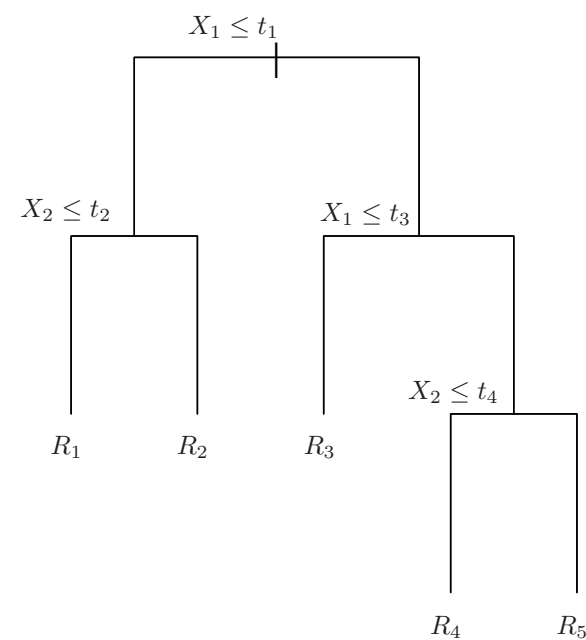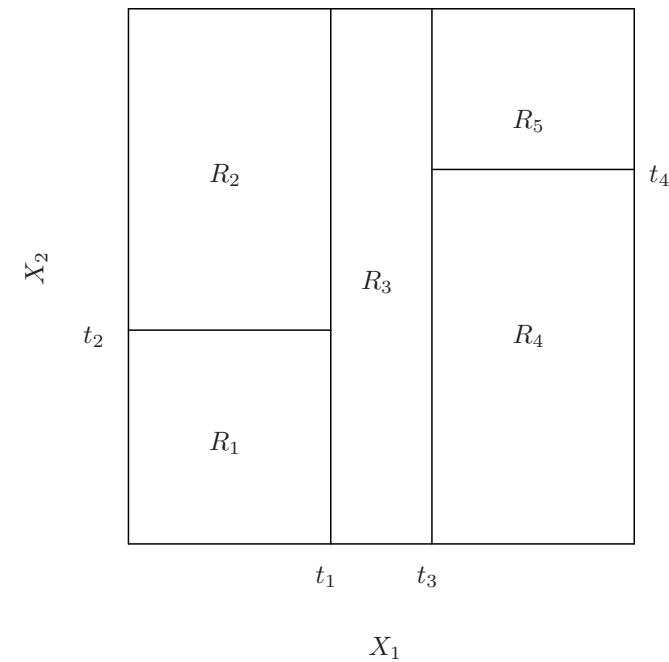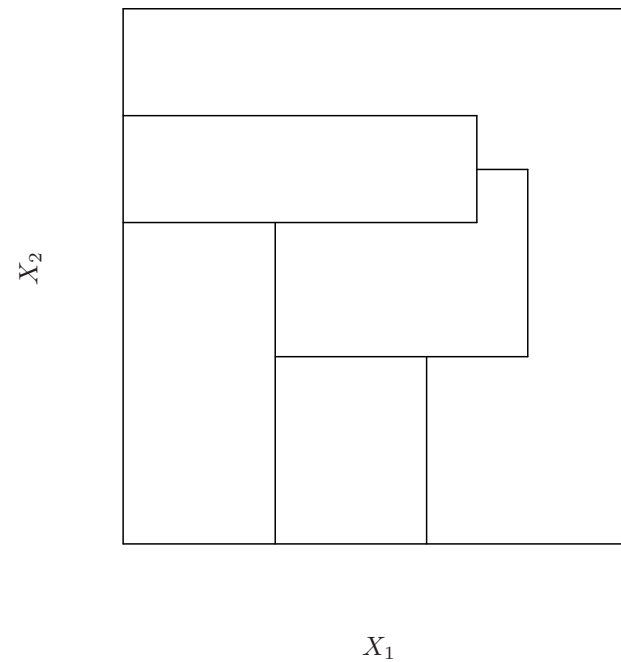
# Minimize the error, be pure

# Many Yes-No decisions yield a tree

# Different visualisations of a tree



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Bigger trees, better prediction?

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# More trees yield a less variable estimate

# Bootstrapping

in bag            out of bag

| 1 1 3 2 2 | 4 5 |

| 2 1 4 2 2 | 3 5 |

| 5 1 5 2 5 | 3 4 |

| ... | ... |

5   4
  1
    2
3

# Random Forests



tree *k*

...

tree 1

Random Forest

▶ random sample of predictors at each split

# Comparing different random forest models



James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Comparing trees, bagging and random forests

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf

# Variable importance

James, G., Witten, D., Hastie, T., & Tibshirani, R. (2013). *An introduction to statistical learning*. Springer. https://faculty.marshall.usc.edu/gareth-james/ISL/ISLR%20Seventh%20Printing.pdf