

Machine Learning In Finance

Introduction

AI/ML in Finance

Benefits of using AI/ML in finance

- Financial operations (e.g., trades) are based on **pre-defined rules**. Automation via AI/ML **reduces costs and increases speed**.
- Financial decisions (e.g., granting a loan, rating a bond, making investments, etc.) require quick, fact-based judgment calls. AI/ML can help **automate decision-making**.
- AI/ML algorithms make fact-based, objective decisions and will **always comply with laws and regulation** (if programmed to do so).
- AI/ML **do not require economic theory**, they **use data to detect patterns**.

A simple ML example

- **Task:** Predict the salaries of people from their age.^a
- **Size of the data sample:** $n = 30$. Divide the data sample into a
 - ▶ Training set $n = 10$
 - ▶ Validation set $n = 10$
 - ▶ Test set $n = 10$
- **Models:** $y = \text{salary}$, $x = \text{age}$, Parameters = a, b_1
 - ▶ Linear ($Y = a + b_1X$)
 - ▶ Polynomial of order 2 ($Y = a + b_1X + b_2X^2$)
 - ▶ Polynomial of order 5 ($Y = a + b_1X + b_2X^2 + b_3X^3 + b_4X^4 + b_5X^5$)

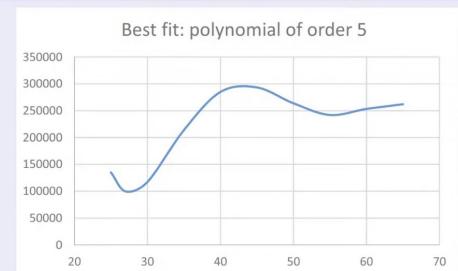
^aThe example and data are taken from [John Hull's website](#).

A simple ML example



Source: Hull (2020).

A simple ML example

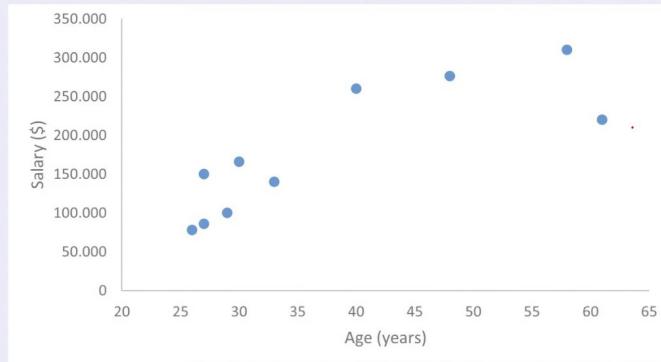


Source: Hull (2020).

● A simple ML example

- The polynomial model of order 5 yields the **best fit**. However, its ups/downs seem unrealistic thus indicating a **potential overfitting**.
- The linear model, in contrast, appears to be oversimplistic. It could suffer from **underfitting**.
- **Solution:** backtest the model using the second, so-called validation set (i.e., check whether the model generalizes well to a validation data set).

● A simple ML example: Validation set



Source: Hull (2020).

● A simple ML example

Root mean square error (rmse) of the three different model candidates:

	Linear	Polynomial (order 2)	Polynomial (order 5)	
Training set	49,731	32,932	12,902	rmse
Validation set	49,990	33,554	38,794	rmse
Difference	259	622	25,892	
Overfitting?	No	No	Yes	

Difference in rmse of Training set & Validation set.

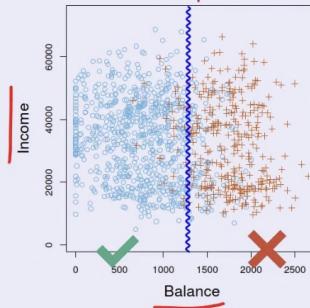
⇒ Polynomial model of order 2 produces the best fit w/o overfitting the data.

● A simple ML example

- How **accurate** is the chosen model?
- $rmse = 32,932$ (training set) or $rmse = 33,554$ (validation set)?
- Answer: none of the two. Accuracy should not be measured based on data sets that were used to choose/validate a model. i.e. **test set**.
- $rmse = 34,273$ (**test set!**)

A second ML example: Classification

- Task: Predict the probability of credit card defaults based on annual income and monthly credit card balance.



Source: James, Witten, Hastie, and Tibshirani (2013). Defaults in orange, non-defaults in blue.

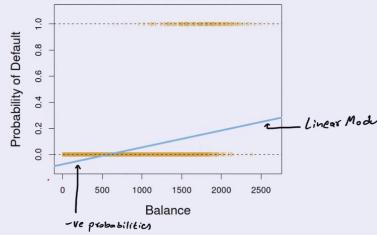
A second ML example



Source: James, Witten, Hastie, and Tibshirani (2013).

A second ML example

- One possible solution: fit a linear model $\text{Default} = a + b \times \text{Income}$ to the data. Problem: probabilities can be negative!

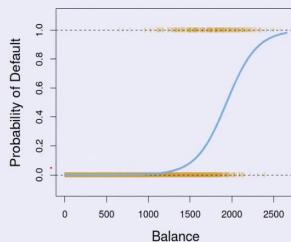


Source: James, Witten, Hastie, and Tibshirani (2013).

A second ML example

- Better solution: fit a logistic model

$$\log \left(\frac{\Pr(\text{Default}=1|\text{Income})}{1-\Pr(\text{Default}=1|\text{Income})} \right) = a + b \cdot \text{Income}$$



Source: James, Witten, Hastie, and Tibshirani (2013).

Data Sources

● Data sources/types I

- Financial data are **diverse** and come at different levels of **complexity**.
- Data can be structured, unstructured, come at low/high frequency, can be publicly available or private, etc.
- Financial data can be complemented with “**alternative data**”, i.e., data from non-financial sources that nevertheless contain information on financial operations (e.g., satellite images).
- More data sources lead to more data; more data in turn leads to the need for AI/ML methods to process this data.
- Different data sources and types make **data preprocessing** necessary.

● Data sources/types II

Types of financial data

Fundamental data	Market data	Analytics	Alternative data
Balance sheet items	Price/yield/implied volatility	Analyst recommendations	Images
Income statement items	Volume	Credit ratings	Google searches
Macro variables	Dividend/coupon	Earnings expectations	Twitter/chats
...	Open interest	News sentiment	Metadata
	Quotes/cancellations

● Data sources/types III

In addition to paid sources like Bloomberg / Compustat / EIKON / Datastream / etc. there are also (lower quality) free of charge sources like yahoo finance. Furthermore, there are many more open data sets available on the internet that can be downloaded for free. These include

- open data set finders like **kaggle** or the **UCI Machine Learning Repository** that contain a wide variety of different data sets,
- public government data sets from **data.gov** (data from multiple US government agencies) and the **EU Open Data Portal** (data published by EU institutions),
- economics data sets from **World Bank Open Data** (economic and development indicators),
- We will learn how to import different kinds of data like csv- and pdf-files, image data, etc. in practical applications throughout the lecture.



Data Preprocessing

● Data preprocessing I

- Financial (and other) data are often
 - **incomplete**: e.g., missing values or attributes of interest, containing only aggregated data
 - **noisy**: e.g., containing errors and outliers
 - **inconsistent**: e.g., containing discrepancies in codes or names.
- Data preprocessing is about **resolving these issues** and **transforming raw data** into an understandable format.
- Data preprocessing is key to a good model performance and most often consumes more time than actually specifying and fitting a model.
- This typically involves two steps. First **understanding** the data, then **preparing** the data before actually employing some or more ML models.

● Data preprocessing III

Data preparation

- The main tasks according to Lesmeister (2015) are:
 - Selecting the data
 - Cleaning the data
 - Constructing the data
 - Integrating the data
 - Formatting the data
- The goal of these tasks is to get the data ready to use as input in the algorithms. This includes merging data from different sources, feature engineering, and further transformations.
- Some algorithms require categorical variables like Yes/No to be in the form 1/0 or to be formatted as factors (in R).
- The data might be split into training, test and validation set.

● Data preprocessing V

The statistical programming language R

- All operations (exploring and preprocessing the data, fitting and validating the model, etc.) will be performed in the **statistical programming language R**.
- R is available as **free software** under the terms of the Free Software Foundation's GNU General Public License in source code form. We recommend the usage of RStudio Desktop as IDE which can be downloaded freely (see [here](#)).
- Alternatively you can use one of the RStudio servers provided by the University Computing Center (see [here](#)).

● Example: German Credit Data

● Practical example I

Importing the data and first explorative analysis

```

1 # Importing the data
2 germanCredit = read.table("http://archive.ics.uci.edu/ml/
   machine?learning?databases/statlog/german/german.data")
3
4 dim(germanCredit)
5 [1] 1000  21  #1000 observations of 21 variables
6
7 #reduce number of features to make them fit on one slide,
8 #21st feature is the response/target variable (credit rating
   (good/bad))
9 germanCredit<-germanCredit[,c(1:5, 21)] ← All rows & 6 columns 1 to 5 & 21.

```

● Data preprocessing II

Data understanding

- The main tasks according to Lesmeister (2015) are:
 - Collecting the data
 - Describing the data
 - Exploring the data
 - Verifying the data quality
- These tasks are performed to make sure that the data are adequate to meet the analytical goal.
- Furthermore, by exploring the data, determining its sparseness, and identifying missing values one gets a better idea of which learning methods might be appropriate.
- Verifying the **data quality** is critical. Understanding who collects the data and how it is collected can help to identify incomplete or erroneous data. It is also possible that rules for collecting the data change over time which might lead to structural breaks in the data.

● Data preprocessing IV

- Data understanding and data preparation are often not performed separately but are interrelated tasks.
- As an **example** we will have a closer look on the **German credit data set** by Dua and Graff (2017) on the following slides.
- The data are readily available via the UCI Machine Learning Repository along with a **description of the data**.

● Data preprocessing VI

The statistical programming language R (cont.)

- R is a **language and environment for statistical computing and graphics** and one of the major languages for data science.
- The language provides a wide variety of statistical (linear and nonlinear modeling, classical statistical tests, time-series analysis, classification, clustering, etc.) and graphical tools.
- R is **highly extensible** via packages. There is a wide variety of R packages available for machine learning and many modern methods in statistical learning are implemented in R.
- Of course there are other alternatives, most notably the programming language **Python** with libraries like scikit-learn and keras/tensorflow.

dim() → Gives dimension of data structure.

● Practical example II

Importing the data and first explorative analysis (cont.)

```
1 #Exploring the structure of the data
2 str(germanCredit)
3
4 'data.frame': 1000 obs. of 6 variables:
5   $ V1 : Factor w/ 4 levels "A11", "A12", "A13",...: 1 2 4 1 1 4
6     4 2 4 2 ...
7   $ V2 : int 6 48 12 42 24 36 24 36 12 30 ...
8   $ V3 : Factor w/ 5 levels "A30", "A31", "A32",...: 5 3 5 3 4 3
9     3 3 3 5 ...
10  $ V4 : Factor w/ 10 levels "A40", "A41", "A410",...: 5 5 8 4 1
11    8 4 2 5 1 ...
12  $ V5 : int 1169 5951 2096 7882 4870 9055 2835 6948 3059
13    5234 ...
14  $ V21: int 1 2 1 1 2 1 1 1 1 2 ...
```

● Practical example III

Importing the data and first explorative analysis (cont.)

The `str()` command gave us valuable information about the **structure** of the data set:

- Our data set consists of 1000 observations of 6 variables.
- The features are of type factor or of type integer.
- **Factors** represent categorical or ordinal variables and have the advantage that category labels are stored only once (requires less RAM and enables faster computations). The different values a factor variable can assume are referred to as **levels**. For example, the first variable ("V1") has the levels "A11", "A12", "A13", and "A14" which are represented by the integers 1, 2, 3, 4.
- We can see the first observations of all variables.

● Practical example IV

● Practical example V

A **description of the variables** and the values they can assume are provided [here](#), e.g.:

Attribute 1 ("V1")

Status of existing checking account
A11 : ... < 0 DM ↪ *less than 0 Deutsche mark*.
A12 : 0 <= ... < 200 DM
A13 : ... >= 200 DM / salary assignments for at least 1 year
A14 : no checking account

Attribute 2 ("V2")

Duration in month

● Practical example VI

And most importantly:

Attribute 21 ("V21")

Rating
1: Good
2: Bad

● Practical example VII

Based on the documentation we can assign **meaningful variable names**. Furthermore, we transform the `data.frame` into a `tibble` which provides an enhanced `print()` method and has further advantages over standard `data.frames` (see [here](#)).

Adjust variable names

```
1 #use the documentation to assign column names
2 colnames(germanCredit)<-c("chkAcctStat", "duration", "
3   credHist", "purpose", "amount", "rating")
4
5 #View some records
6 #Therefore, convert data.frame to tibble for more convenient
7   printing functionality
8 library(tibble)
9 germanCredit<-as.tibble(germanCredit)
```

● Practical example VIII

Show first records

```

1 print(germanCredit)
2
3 #A tibble: 1,000 x 6
4 chkdAcctStat duration credHist purpose amount rating
5 <fct>     <int> <fct>    <fct>    <int>    <int>
6 1 A11        6   A34   A43    1169      1
7 2 A12       48   A32   A43    5951      2
8 3 A14       12   A34   A46    2096      1
9 4 A11       42   A32   A42    7882      1
10 5 A11      24   A33   A40    4870      2
11 6 A14       36   A32   A46    9055      1
12 7 A14      24   A32   A42    2835      1
13 8 A12       36   A32   A41    6948      1
14 9 A14       12   A32   A43    3059      1
15 10 A12      30   A34   A40    5234      2
16 #... with 990 more rows

```

● Practical example IX

As a first analysis, we have a look at some **summary statistics** via the `summary()` command. The output of the function depends on the variable type.

Summary statistics

```
1 summary(germanCredit)
```

chkdAcctStat	duration	credHist	purpose	amount	rating
A11:274	Min. : 4.0	A30: 40	A43 :280	Min. : 250	Min. :1.0
A12:269	1st Qu.:12.0	A31: 49	A40 :234	1st Qu.: 1366	1st Qu.:1.0
A13: 63	Median :18.0	A32:530	A42 :181	Median : 2320	Median :1.0
A14:394	Mean :20.9	A33: 88	A41 :103	Mean : 3271	Mean :1.3
	3rd Qu.:24.0	A34:293	A49 : 97	3rd Qu.: 3972	3rd Qu.:2.0
	Max. :72.0		A46 : 50	Max. :18424	Max. :2.0
	(Other): 55				

● Practical example X

We now transform the integer variables into numeric (float) ones to facilitate later processing of the data. For this and further operations we rely on the `dplyr` package. According to the official [website](#), "dplyr is a grammar of data manipulation, providing a consistent set of verbs that help you solve the most common data manipulation challenges:

- `mutate()` adds new variables that are functions of existing variables.
- `select()` picks variables based on their names.
- `filter()` picks cases based on their values.
- `summarise()` reduces multiple values down to a single summary.
- `arrange()` changes the ordering of the rows."

We will use these operations in the practical applications throughout the lecture.

`dplyr` is part of `tidyverse`, a collection of R packages specifically designed for data science where all packages share an underlying design philosophy, grammar, and data structure (tibble).



● Practical example XI

Adjust variable types

```

1 #transform integer variables to numeric ones
2 #therefore, use dplyr-style notation with the pipe operator
3 germanCredit<- germanCredit %>%
4   mutate_if(is.integer, as.numeric)
5
6 #transform rating into a factor
7 germanCredit<- germanCredit %>%
8   mutate(rating=factor(rating, labels =c("Good",
9     "Bad")))
10
11 #Show percentage of good/bad rated credit applications
12 table(germanCredit$rating)/nrow(germanCredit)
13 Good  Bad
14 0.7  0.3

```

● Practical example XII

- On the following slides we will shortly analyze some of the variables.
- In particular, we want to gain an understanding of the **predictive power** of the covariates on the dependent variable credit rating.
- This can be done via the ratio **weight of evidence (WOE)**.

Weight of evidence

WOE encodes the relationship of a (categorical) predictor variable with a **binary target variable**. It originated in the finance industry to help separate "good" from "bad" risks (loan default). By now, it has been in use for more than forty years and is still most widely known in the finance and insurance industry.

● Practical example XIII

Weight of evidence (cont.)

- The ratio is defined as:

$$WOE := \log \left(\frac{\text{numberOfNonEvents}(Good)}{\text{numberOfEvents}(Bad)} \right),$$

where event means a default.

- Ratios of non events to events close to 1 indicate that the corresponding category has no predictive power on the target variable. This corresponds to a value near zero after applying the logarithm.
- One should be careful in deriving conclusions from the ratio as giving a loan to a defaulting customer is worse than not giving a loan to a non-defaulting (potential) customer.

● Practical example XIV

Weight of evidence (cont.)

- WOE is calculated in different **groups** that are formed based on the covariate of interest.
- For categorical variables these might be the categories or a pooling of multiple (smaller) subcategories.
- For **continuous variables** one has to create **bins** based on thresholds.
- Each category/bin should contain at least 5 % of all observations to avoid results to be driven by noise.

● Practical example XV

Checking account status : Calculating WOE

```

1 #Attribute 1: (qualitative)
2 #Status of existing checking account
3 #A11 : ... < 0 DM
4 #A12 : 0 <= ... < 200 DM
5 #A13 : ... >= 200 DM / salary assignments for at least 1 year
6 #A14 : no checking account
7
8 #Calculate WOE and some further simple ratios
9 tmpStats<-germanCredit %>%
10   select(chkAcctStat, rating) %>%
11   group_by(chkAcctStat) %>%
12   summarize (pctOfTotalObs = length(rating)/nrow(germanCredit),
13             goodRate=mean(rating=="Good"),
14             WOE=log(sum(rating=="Good")/sum(rating=="Bad")))

```

● Practical example XVI

Checking account status (cont.)

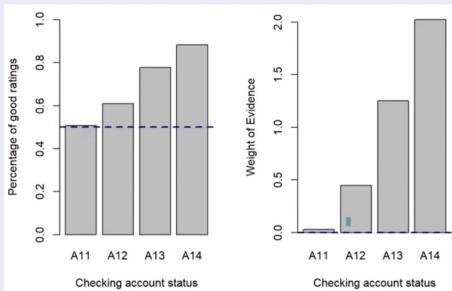
```

1 print(tmpStats)
2
3 #A tibble: 4 x 4
4 chkAcctStat pctOfTotalObs goodRate      WOE
5 <fct>          <dbl>     <dbl>     <dbl>
6 1 A11           0.274     0.507  0.0292
7 2 A12           0.269     0.610  0.446
8 3 A13           0.063     0.778  1.25
9 4 A14           0.394     0.883  2.02

```

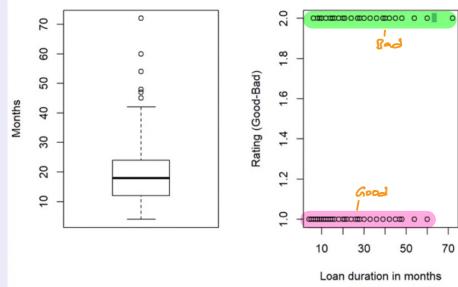
● Practical example XVII

Credit rating ~ checking account status



● Practical example XVIII

Loan duration



● Practical example XIX

When calculating the WOE separately for each unique value of loan duration, some groups only contain very few observations. We therefore group the loan durations on a yearly basis and recalculate WOE.

WOE for loan duration

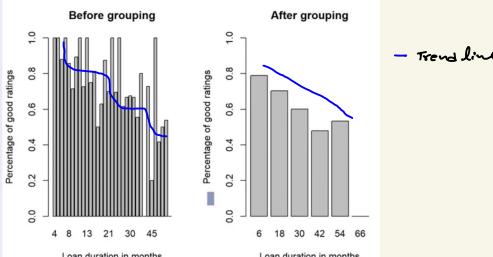
```

1 #Calculation for each unique value of loan duration
2 tmpStats<-germanCredit %>%
3   select(duration, rating) %>%
4   group_by(duration) %>%
5   summarize(goodRate=mean(rating=="Good"))
6
7 #Form larger groups on a yearly basis
8 tmpStats2<-germanCredit %>%
9   select(duration, rating) %>%
10  mutate (interval=cut(duration, breaks=12*(0:6)))
11
12 group_by(interval) %>%
13   summarize (goodRate=mean(rating=="Good"))

```

● Practical example XX

Loan duration ~ rating



Practical example XXI

We finally have a look at the credit history variable. For brevity, we do not consider further variables and leave this as an exercise.

Credit history

```
1 #A30 : no credits taken/all credits paid back duly
2 #A31 : all credits at this bank paid back duly
3 #A32 : existing credits paid back duly till now
4 #A33 : delay in paying off in the past
5 #A34 : critical account/other credits existing (not at this
       bank)
6
7 #Calculate some statistics for the five groups
8 tmpStats<-germanCredit %>%
9   select(credHist, rating) %>%
10  group_by(credHist) %>%
11  summarize (pctOfTotalObs = length(rating)/nrow(
12    germanCredit), goodRate=mean(rating=="Good")
13    , WOE=log(sum(rating=="Good")/sum(rating=="Bad")))
14 )
```

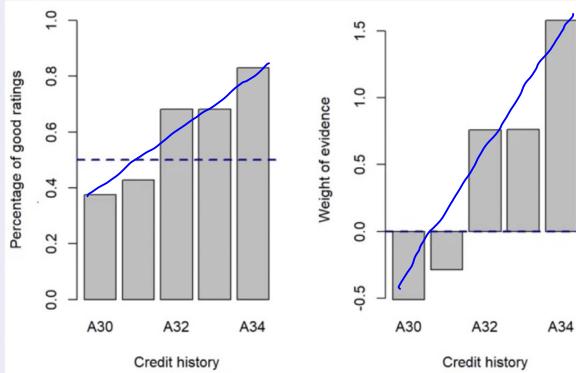
● Practical example XXII

Credit history (cont.)

```
1 print(tmpStats)
2
3 #A tibble: 5 x 4
4 credHist pctOfTotalObs goodRate      WOE
5 <fct>      <dbl>     <dbl>     <dbl>
6 1 A30        0.04     0.375 -0.511
7 2 A31        0.049    0.429 -0.288
8 3 A32        0.53     0.681  0.759
9 4 A33        0.088    0.682  0.762
10 5 A34       0.293    0.829  1.58
```

● Practical example XXIII

Rating ~ credit history



A34 is highly predictive
& has high explanatory
power for explaining bad
credit rating.

Data Generation

● Data generation I

- Many machine learning techniques rely on **synthetic data** for training and testing.
- Synthetic data is **artificial data created by algorithms** that mirror statistical properties of the original data while not revealing information regarding real people.
- The main purpose of using synthetic data is **preserving privacy, creating training data** for machine learning algorithms and **testing systems**.
- Many training data sets are highly imbalanced (e.g., credit card fraud data) making classification tasks difficult. In these cases, synthetic data generation is particularly important for building accurate machine learning models.

● Data generation II

Classification of synthetic data according to Aggarwal and Yu (2008)

- **Fully synthetic** data: The data is completely synthetic and does not contain any original data entries. Therefore, the (joint) density of the original data is estimated and random values are sampled from the estimated density function. The data has **strong privacy protection** but the **truthfulness** of the data is lost.
- **Partially synthetic** data: The method replaces values of **selected attributes** that have a high risk of disclosure with synthetic data. This technique can also be useful for imputing missing values in the data. Disclosure risk is higher than in fully synthetic data.
- **Hybrid synthetic** data: The data set is generated using both original and synthetic data. Each record in the original data set is replaced by the nearest record in the synthetic data. The method combines good privacy protection with high utility at the cost of more memory and processing time.

Fitting statistical distribution to original data & then sampling synthetic data from that fitted statistical distribution.

● Data generation III

- There are different approaches for generating synthetic data.
- Three broad concepts are
 - (i) Generating data from a **known distribution**
 - (ii) **Fitting a distribution** to real data and generating from this distribution
 - (iii) Using deep learning

● Data generation IV

(i) Generating data from a known distribution

- Applicable in cases where no real data exist but the researcher/practitioner has a **comprehensive understanding of the distribution** of the data of interest (e.g., normal, Student-t, exponential, chi-squared, etc.)
- One can then simulate random data from this a-priori chosen distribution.
- Utility of the data depends on the degree of knowledge about the underlying distribution of the data

• Data generation V

(ii) Fitting a distribution to real data

- If there is real data, one can determine the best fit distribution chosen from a given family of distributions.
- Synthetic data can now be generated from the best fit distribution via Monte Carlo methods.
- The quality of the generated data depends on the selection of the family of distributions (e.g., only normal distributions are considered) as well as on the goodness of fit of the final distribution.
- Data can also be modeled by machine learning models such as decision trees providing an approximation to non-classical distributions (e.g., multi-modal distributions). In these cases, **overfitting** might be an issue.

• Data generation VI

(iii) Using deep learning

- Deep generative models such as **Variational Autoencoder (VAE)** or **Generative Adversarial Network (GAN)** can create synthetic data.
- VAE is an **unsupervised** method where the original data is compressed by an **encoder** into a more compact structure. The **decoder** then generates a representation of the original data from the compressed data. The system is trained to minimize the differences between output and original data.

• Data generation VII

Using deep learning (cont.)

- In the **GAN** model **two separate networks** are trained iteratively. The first network (the **generator**) takes random sample data to create a synthetic data set. The second network (the **discriminator**) compares the synthetic data with a real data set.
- The generator network is trained to make discriminating between the generated and the real data for the discriminator network as hard as possible.
- The method is frequently used for generating image data, see **here** for a demonstration.

Simple Linear Regression

● Introduction I

- We start our discussion of AI/ML methods in finance with simple models for **statistical learning**.
- Stat. learning refers to the set of tools used for *understanding* data.
- Broadly speaking, we can differentiate between the following types of models:
 - ▶ Supervised algorithms
 - ▶ Unsupervised algorithms

● Introduction II

Supervised vs. unsupervised learning

- **Supervised statistical learning:** Estimate/Predict an output based on inputs.
 - ▶ Example 1: Predict Y (wage) based on X_1 (education), X_2 (gender),... (called a **regression problem**)
 - ▶ Example 2: Predict ups/downs of the S&P 500 ($Y = 0, 1$) based on X_1 , X_2 ,... (called a **classification problem**)
- **Unsupervised statistical learning:** Group input observations.
 - ▶ Example 1: Group customers based on inputs X_1 (age), X_2 (income),... (called a **clustering problem**)

● Simple linear regression I

Starting point: **linear regression** as the foundation for many “modern” supervised statistical learning approaches.

Simple linear regression

- Predict a quantitative response Y on the basis of a single predictor variable X :

$$Y \approx \beta_0 + \beta_1 X \quad (1)$$

- Y is “regressed” on X .
- Two unknown parameters/coefficients β_0 (intercept) and β_1 (slope) are estimated by minimizing the residual sum of squares (RSS).
- Result:

$$\hat{y} \approx \hat{\beta}_0 + \hat{\beta}_1 x \quad (2)$$

● Simple linear regression II

How to estimate?

Simple linear regression

- Let $\hat{y}_i = \hat{\beta}_0 + \hat{\beta}_1 x_i$ be the prediction for Y based on the i th observation of X . Then

$$e_i = y_i - \hat{y}_i \quad (3)$$

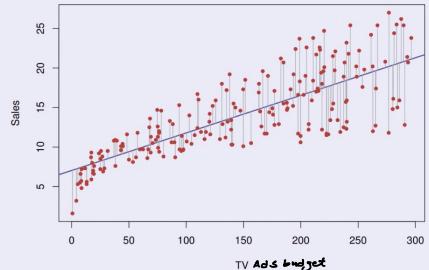
is called the i th **residual** (the difference between the i th observed and predicted values of the response variable).

- The **residual sum of squares (RSS)** is then defined as

$$RSS = \sum_{i=1}^n e_i^2 \quad (4)$$

Simple linear regression III

Residuals



Source: James, Witten, Hastie, and Tibshirani (2013). Data: <https://www.statlearning.com/s/Advertising.csv>

Simple linear regression IV

How to estimate?

Simple linear regression

Minimization of the RSS leads to the well-known OLS estimator:

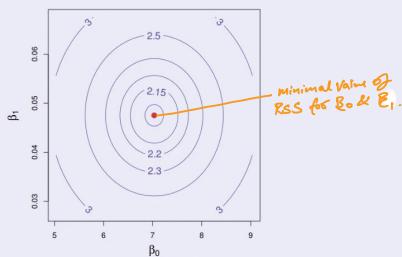
$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (5)$$

$$\hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}. \quad (6)$$

with \bar{x} and \bar{y} being the sample means of X and Y .

Simple linear regression V

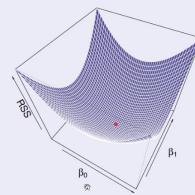
Contour plot of the RSS (response: sales; predictor: TV)



Source: James, Witten, Hastie, and Tibshirani (2013). Data: <https://www.statlearning.com/s/Advertising.csv>

Simple linear regression VI

Three-dimensional plot of the RSS (response: sales; predictor: TV)

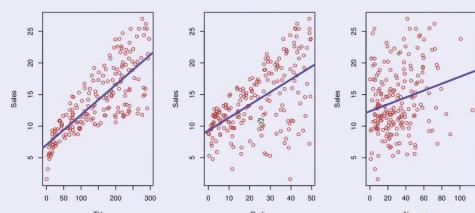


Source: James, Witten, Hastie, and Tibshirani (2013). Data: <https://www.statlearning.com/s/Advertising.csv>

Simple linear regression VII

Resulting OLS regression lines

- Y ...sales; X_1, X_2, X_3 ...advertising budgets.



Source: James, Witten, Hastie, and Tibshirani (2013). Data: <https://www.statlearning.com/s/Advertising.csv>

TV & Sales looks to have linear regression
Not so much for Newspaper & sales.

● Simple linear regression VIII

How to assess the **accuracy** of the **coefficient estimates** assuming the true relation is $Y = \beta_0 + \beta_1 X + \varepsilon$?

S.e. of the OLS coefficients

The OLS coefficients are estimated from a sample (x_i, y_i) . Thus, the accuracy of the estimates can be assessed by looking at the **standard errors of the estimators**.

$$\text{SE}(\hat{\beta}_0)^2 = \sigma^2 \left[\frac{1}{n} + \frac{\bar{x}^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \right] \quad (7)$$

and

$$\text{SE}(\hat{\beta}_1)^2 = \frac{\sigma^2}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (8)$$

with $\sigma^2 = \text{Var}(\varepsilon)$.

● Simple linear regression IX

How to estimate σ^2 ?

S.e. of the OLS coefficients

Estimate σ via the **residual standard error**:

$$\text{RSE} = \sqrt{\text{RSS}/(n-2)} \quad (9)$$

● Simple linear regression X

What to do with $\text{SE}(\hat{\beta}_0)$ (or more precisely $\text{SE}(\hat{\beta}_0)$)?

Confidence intervals for the OLS coefficients

The 95%-c.i. for β_0 and β_1 are approximately given by

$$\hat{\beta}_0 \pm 2 \cdot \text{SE}(\hat{\beta}_0) \quad \text{and} \quad \hat{\beta}_1 \pm 2 \cdot \text{SE}(\hat{\beta}_1) \quad (10)$$

● Simple linear regression XI

What to do with $\text{SE}(\hat{\beta}_0)$ (or more precisely $\text{SE}(\hat{\beta}_0)$)?

Hypothesis testing for the OLS coefficients

We test the hypothesis

$$H_0 : \beta_1 = 0 \text{ versus } H_1 : \beta_1 \neq 0$$

using the **t-statistic**

$$t = \frac{\hat{\beta}_1 - 0}{\text{SE}(\hat{\beta}_1)} \quad (11)$$

which measures the no. of standard deviations that $\hat{\beta}_1$ is away from 0.

Simple linear regression XII

How to assess the **accuracy** of the **model itself** assuming that we rejected $H_0 : \beta_1 = 0$?

Residual Standard Error

Recall that associated with each individual observation (x_i, y_i) is an error term ε/e_i . Even if the coefficients were known, these error terms prevent us from perfectly predicting Y from X . Estimate the s.e. of ε via the RSE

$$\text{RSE} = \sqrt{\frac{1}{n-2} \text{RSS}} = \sqrt{\frac{1}{n-2} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (12)$$

● Simple linear regression XIII

How to assess the **accuracy** of the **model itself** assuming that we rejected $H_0 : \beta_1 = 0$?

RSE provides an **absolute** measure of lack of fit of the model (it is given in terms of units of Y).

Problem: RSE depends on the scale of Y !

R^2

Solution: R^2 , the proportion of variance explained by the model.

$$R^2 = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}} \quad (13)$$

with $\text{TSS} = \sum (y_i - \bar{y})^2$.

Note: $R^2 = r^2$ with $r = \text{Cor}(X, Y)$.

Multiple Linear Regression

● Multiple linear regression I

Simple linear regression is a useful approach for predicting a response **on the basis of a single predictor variable**. In practice: more than one predictor (in fact, in ML often thousands!). Possible approaches:

- Estimate several simple regressions. Problem: correlations among predictors will bias the results (over-/underestimate the effect of single predictors).
 - Estimate a **multivariate/multiple linear regression**.

● Multiple linear regression II

Multiple linear regression

- Predict a quantitative response Y on the basis of several predictor variable X_1, X_2, \dots :

$$Y \approx \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \cdots + \beta_p X_p \quad (14)$$

- Y is “regressed” on the set X_1, \dots, X_p .
 - $p+1$ unknown parameters/coefficients β_0 (intercept) and β_i (slopes) are estimated by again minimizing the residual sum of squares (RSS):
use OLS to estimate coefficients β_i :

$$\hat{\beta} = (X'X)^{-1} X'Y \quad \text{where } X' \text{ is } X \text{ transpose.}$$

Vector of Coefficients

- **Result:**

$$\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x_1 + \hat{\beta}_2 x_2 + \cdots + \hat{\beta}_p x_p. \quad (15)$$

● Multiple linear regression IV

Multiple linear regression in R

Multiple regression coefficient estimates from a regression of product sales on advertising budget, community income level, average age, and education.

```
1 > library(MASS)    → library for Regression.  
2 > library(ISLR)   → carsales data set.  
3 > lm.fit=lm(Sales~Advertising+Income+Age+Education, data=  
           Carsales) → fitting linear Model.
```

● Multiple linear regression V

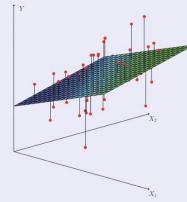
Multiple linear regression in R (cont.)

```
1 > summary(lm.fit)
2
3 Call:
4 lm(formula = Sales ~ Advertising + Income + Age + Education,
5      data = Carseats)
6
7 Residuals:
8   Min     1Q Median     3Q    Max
9 -8.1676 -1.8869 -0.1805  1.6827  8.5494
```

● Multiple linear regression III

Resulting OLS regression plane

- Regression plane with two predictors.



Source: James, Witten, Hastie, and Tibshirani (2013)

- If we have more predictors/features, we get linear hyperplane.

● Multiple linear regression VI

Multiple linear regression in R (cont.)

```

 1 Coefficients: (Intercept) Advertising Income Age Education
 2 Estimate Std. Error t value Pr(>|t|) 
 3 < 2e-16 *** 0.110168 0.019798 5.565 4.85e-08 ***
 4 < 2e-16 *** 0.013485 0.004710 2.863 0.004424 **
 5 < 2e-16 *** -0.040058 0.008109 -4.940 1.16e-06 ***
 6 < 2e-16 *** -0.036805 0.050237 -0.733 0.46422
 7 < 2e-16 *** 
 8 --- 
 9 Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' '
10 
11 Residual standard error: 2.624 on 395 degrees of freedom
12 Multiple R-squared: 0.1453, Adjusted R-squared: 0.1366
13 F-statistic: 16.78 on 4 and 395 DF, p-value: 1.022e-12

```

● Multiple linear regression VII

Where to go next?

- Is at least one of the predictors X_1, X_2, \dots, X_p useful in predicting the response? (\rightarrow F-statistic).
- Do we need all or only a subset of the predictors? (\rightarrow variable selection).
- How well does the model fit the data? (\rightarrow RSE and R^2).
- Given a set of predictor values, what response value should we predict, and how accurate is our prediction? (\rightarrow forecasting).

● Multiple linear regression VIII

Is at least one predictor useful?

F-statistic

We test the null hypothesis,

$$H_0 : \beta_1 = \beta_2 = \dots = \beta_p = 0$$

versus the alternative

$$H_a : \text{at least one } \beta_j \text{ is non-zero.}$$

This hypothesis is tested using the following F-statistic:

$$F = \frac{(TSS - RSS)/p}{RSS/(n - p - 1)}. \quad (16)$$

No relationship between the response and predictors: F-statistic close to 1.
 H_a is true: F greater than 1.

● Multiple linear regression IX

Which predictors are significant?

t-statistics / p-values

- Estimate F-statistics for models from which individual predictors are omitted.
- Calculate t-statistics for each predictor (the square of each t-statistic is the corresponding F-statistic).

Caution!

If $p > n$ or even $p \gg n$ then there are more coefficients β_j to estimate than observations from which to estimate them.

- We cannot fit the regression using (multiple) OLS.
- We cannot use the F-statistic.

→ When no. of Variables is larger than no. of Observations, e.g.: 100 variables & 10 observations → more coefficients β_i to estimate than observations from which to estimate them.

● Multiple linear regression X (Variable Selection)

Which predictors are powerful?

Variable selection

Process of determining which predictors are associated with the response, and fitting a single model that involves only the significant predictors.

Alternative approaches:

- Compare the fit of all models and calculate measures of the model's fit (e.g., Akaike's information criterion/AIC, Bayesian information criterion/BIC, adjusted R^2 , etc.). Becomes infeasible as the number of different models for p variables is 2^p . (e.g., for $p = 30$ we have 1,073,741,824 different models.)
- Use an automated and efficient procedure to select a subset of models.

→ Estimate all possible model (infeasible as variables increase)

● Multiple linear regression XI

Three classical approaches for Variable selection:

(i) Forward selection

- Start with the **null model** (contains an intercept only).
- Estimate p simple linear regressions and add to the null model the variable from the regression with the lowest RSS.
- Continue adding variables in this manner until some stopping rule is fulfilled.

● Multiple linear regression XIII

Three classical approaches for Variable selection:

(ii) Mixed selection

- Start with the **null model** (contains an intercept only).
- Estimate p simple linear regressions and add to the null model the variable from the regression with the lowest RSS.
- Remove variables that turn insignificant when new variables are added to the model. (because of correlations between predictors)
- Continue adding and removing variables in this manner until some stopping rule is fulfilled.

● Multiple linear regression XV

Variable selection in R (cont.)

```

1 > summary(Credit)
2
3   Education      Gender     Student    Married
4   Min. : 5.00   Male:193   No :360   No :155
5   1st Qu.:11.00 Female:207   Yes: 40   Yes:245
6   Median :14.00
7   Mean  :13.45
8   3rd Qu.:16.00
9   Max.  :20.00
10
11  Ethnicity      Balance response variable
12  African American: 99  Min.   : 0.00
13  Asian           :102  1st Qu.: 68.75
14  Caucasian       :199  Median : 459.50
15                           Mean   : 520.01
16                           3rd Qu.: 863.00
17                           Max.   :1999.00

```

● Multiple linear regression XVII

Variable selection in R (cont.)

```

1 > summary(lm(Balance ~ Student, data=Credit))
2
3 [...]
4
5 Residual standard error: 444.6 on 398 degrees of freedom
6 Multiple R-squared:  0.06709, Adjusted R-squared:  0.06475
7 F-statistic: 28.62 on 1 and 398 DF, p-value: 1.488e-07

```

Closer to 7% of total variance is explained by our Model.

● Multiple linear regression XII

Three classical approaches for Variable selection:

(iii) Backward selection

- Start with the model that contains all p variables.
- Remove the variable with the largest p-value.
- Continue removing variables in this manner until some stopping rule is fulfilled.

● Multiple linear regression XIV

Variable selection in R

```

1 > summary(Credit) -> credit card data.
2
3   ID          Income      Limit
4   Min.   : 1.0   Min.   :10.35   Min.   : 855
5   1st Qu.:100.8 1st Qu.:21.01   1st Qu.: 3088
6   Median :200.5 Median :33.12   Median : 4622
7   Mean   :200.5 Mean   :45.22   Mean   : 4736
8   3rd Qu.:300.2 3rd Qu.:57.47   3rd Qu.: 5873
9   Max.   :400.0  Max.   :186.63   Max.   :13913
10
11  Rating      Cards      Age
12  Min.   : 93.0  Min.   :1.000   Min.   :23.00
13  1st Qu.:247.2 1st Qu.:2.000   1st Qu.:41.75
14  Median :344.0 Median :3.000   Median :56.00
15  Mean   :354.9 Mean   :2.958   Mean   :55.67
16  3rd Qu.:437.2 3rd Qu.:4.000   3rd Qu.:70.00
17  Max.   :982.0  Max.   :9.000   Max.   :98.00

```

● Multiple linear regression XVI

Variable selection in R (cont.) *Simple Regression
(Choosing only Student Variable)*

```

1 > summary(lm(Balance ~ Student, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Student, data = Credit)
5
6 Residuals:
7   Min.   1Q   Median   3Q   Max.
8 -876.82 -458.82 -40.87 341.88 1518.63
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) 480.37    23.43   20.50 < 2e-16 ***
13 StudentYes 396.46    74.10   5.35 1.14e-07 ***
14 Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1
1

```

● Multiple linear regression XVIII

Variable selection in R (cont.) *Simple Regression
(Choosing Credit Card limit Variable)*

```

1 > summary(lm(Balance ~ Limit, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Limit, data = Credit)
5
6 Residuals:
7   Min.   1Q   Median   3Q   Max.
8 -676.95 -141.87 -11.55 134.11 776.44
9
10 Coefficients:
11             Estimate Std. Error t value Pr(>|t|)
12 (Intercept) -2.928e+02 2.668e+01 -10.97 < 2e-16 ***
13 Limit        1.716e-01 5.066e-03 33.88 < 2e-16 ***
14 Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1
1

```

● Multiple linear regression XIX

Variable selection in R (cont.)

```

1 > summary(lm(Balance ~ Limit, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Limit, data = Credit)
5
6 [...]
7
8 Residual standard error: 233.6 on 398 degrees of freedom
9 Multiple R-squared:  0.7425, Adjusted R-squared:  0.7419
10 F-statistic: 1148 on 1 and 398 DF, p-value: < 2.2e-16

```

Limit Variable selected to have much more explanatory power when it comes to Balance.

● Multiple linear regression XXI

Variable selection in R (cont.)

```

1 > summary(lm(Balance ~ Income, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Income, data = Credit)
5
6 [...]
7
8 Residual standard error: 407.9 on 398 degrees of freedom
9 Multiple R-squared:  0.215, Adjusted R-squared:  0.213
10 F-statistic: 109 on 1 and 398 DF, p-value: < 2.2e-16

```

Should Select Variable Income for Regression.

● Multiple linear regression XXII

Variable selection in R (cont.) *(Multiple Regression Using 3 Variables Income, Student & Limit)*

```

1 > summary(lm(Balance ~ Income+Student+Limit, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Income + Student + Limit, data =
  Credit)
5 Residuals:
6   Min    1Q Median    3Q   Max
7 -186.60 -78.43 -17.33  62.65 314.88
8 Coefficients:
9             Estimate Std. Error t value Pr(>|t|)
10 (Intercept) -4.323e+02 1.243e+01 -34.78 <2e-16 ***
11 Income      -7.902e+00 2.432e-01 -32.49 <2e-16 ***
12 StudentYes  4.270e+02 1.742e+01  24.52 <2e-16 ***
13 Limit        2.675e-01 3.712e-03  72.07 <2e-16 ***

```

All variables highly significant.

● Multiple linear regression XXIV

Variable selection in R (cont.)

```

1 Residual standard error: 104.4 on 396 degrees of freedom
2 Multiple R-squared:  0.9488, Adjusted R-squared:  0.9484
3 F-statistic: 2447 on 3 and 396 DF, p-value: < 2.2e-16
4
5 > AIC(lm(Balance ~ Income+Student+Limit, data=Credit))
6 [1] 4859.793

```

AIC is also much lower than before.

● Multiple linear regression XX

Variable selection in R (cont.) *Simple Regression (Choosing Important Variable)*

```

1 > summary(lm(Balance ~ Income, data=Credit))
2
3 Call:
4 lm(formula = Balance ~ Income, data = Credit)
5
6 Residuals:
7   Min    1Q Median    3Q   Max
8 -803.64 -348.99 -54.42 331.75 1100.25
9
10 Coefficients:
11              Estimate Std. Error t value Pr(>|t|)
12 (Intercept) 246.5148 33.1993 7.425 6.9e-13 ***
13 Income       6.0484  0.5794 10.440 < 2e-16 ***
14 Signif. codes: 0 *** 0.001 ** 0.01 * 0.05 . 0.1 .

```

high T-value

● Multiple linear regression XXII

Variable selection in R (cont.) *Akaike Information Criterion (lower AIC => Better Model fit)*

```

1 > AIC(lm(Balance ~ Student, data=Credit))
2 [1] 6016.933 (3rd Best)
3 > AIC(lm(Balance ~ Limit, data=Credit))
4 [1] 5501.983 (Best)
5 > AIC(lm(Balance ~ Income, data=Credit))
6 [1] 5947.894 (2nd Best)

```

● Multiple linear regression XXV

How to assess the model fit?

Adjusted R^2 *(Use this for Multiple Regression)*

Adjust

$$R^2 = 1 - \frac{\text{RSS}}{\text{TSS}}$$

by the number of variables d used in the model:

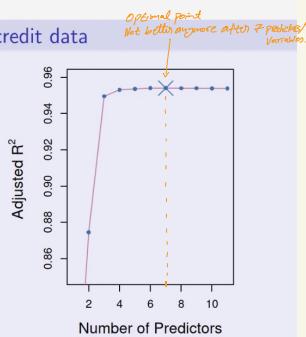
$$\text{Adjusted } R^2 = 1 - \frac{\text{RSS}/(n-d-1)}{\text{TSS}/(n-1)} \quad (17)$$

Intuition: adding noise variables to a well fitted model will decrease RSS only by a little bit. Therefore, the adjusted R^2 punishes the model for every included unnecessary predictor, to avoid overfitting.

Almost 95%

● Multiple linear regression XXVI

Adjusted R^2 for the credit data



Source: James, Witten, Hastie, and Tibshirani (2013).

● Multiple linear regression XXVII

● Multiple linear regression XXVII

How to **predict** from the model?

Prediction/Forecasting

Use (15) and predict values \hat{y} from given values x_i and estimated coefficients $\hat{\beta}_i$.

Bias will stem from

- Estimation error / reducible error: imprecise estimates of $\hat{\beta}_i$. (Only have finite data)
- Model bias: using a linear regression function instead of $y = f(x)$. (Front=)
- Irreducible error: the error term.

How to **extend** the multiple linear regression model?

Qualitative predictors

Some variables are **qualitative** rather than quantitative/metric. For example: *age* and *gender* in the credit card data. How to deal with qualitative predictors?

- Predictor with only two levels: use a **dummy variable**:

$$x_i = \begin{cases} 1, & \text{if } i\text{th person is female} \\ 0, & \text{if } i\text{th person is male} \end{cases}$$

- Predictor with k levels: use $k - 1$ dummy variables.
- Caution: the interpretation of the estimated coefficients changes!

e.g.: Age: 3 dummy variables & 4 levels.

$$\begin{aligned} x_1 &= 1 \quad \text{for } \text{age} \leq 20 \\ x_2 &= 1 \quad \text{for } 20 < \text{age} \leq 40 \\ x_3 &= 1 \quad \text{for } 40 < \text{age} \leq 60 \end{aligned}$$

● Multiple linear regression XXIX

How to **extend** the multiple linear regression model? Even further.

Interaction terms

Predictor X_1 increases/decreases the effect of a different predictor X_2 on Y .

- Use a so-called **interaction term**:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_1 X_2 + \varepsilon.$$

↑ response
↑ interaction term

- Again caution: interpretation of β_1 , β_2 , and β_3 are different!
- If you include the interaction term $X_1 X_2$, one must include X_1 and X_2 as well.
- Particularly interesting: interaction between a qualitative and a quantitative predictor.

● Multiple linear regression XXX

How to **extend** the multiple linear regression model?

Polynomial regression

Predictor X_1 influences Y **non-linearly**.

- One possible solution: Use **polynomial regression**:

$$Y = \beta_0 + \beta_1 X_1 + \beta_2 X_1^2 + \varepsilon.$$

- Including non-linear terms does not change the fact that this is still a linear model!
- Using higher polynomials (just as adding predictors) can lead to **overfitting**.

● Multiple linear regression XXXI

Polynomial regression in R (last slide XXX example)

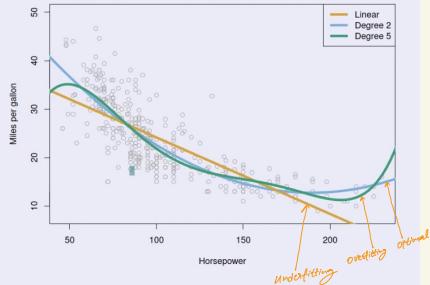
```

1 > library(ISLR)
2 > summary(lm(formula = mpg ~ horsepower + I(horsepower^2),
   data = Auto))
3
4 Call:
5 lm(formula = mpg ~ horsepower + I(horsepower^2), data = Auto)
6
7 Residuals:
8   Min    1Q  Median    3Q   Max
9 -14.7135 -2.5943 -0.0859  2.2868 15.8961
  
```

↑ need to use this form

● Multiple linear regression XXXII

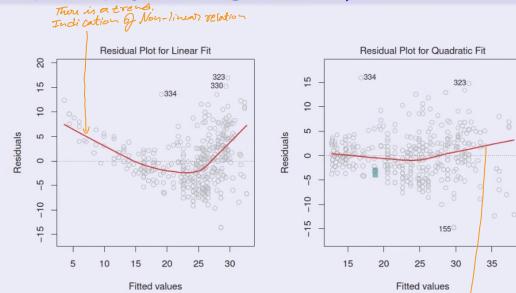
Polynomial regression of mpg on horsepower



Source: James, Witten, Hastie, and Tibshirani (2013).

● Multiple linear regression XXXV

Residual plot of polynomial regressions (a)



Source: James, Witten, Hastie, and Tibshirani (2013).

No clear trend in residuals

Quadratic Model is better than linear one as residuals are almost constant

From these 2 plots we imply
 → Non-linear relationship
 → Quadratic Model better.

● Multiple linear regression XXXII

Polynomial regression in R (cont.)

```

1 Coefficients:
2                               Estimate Std. Error t value Pr(>|t|)
3 (Intercept)           56.900097  0.8004268  71.60   <2e-16 ***
4 horsepower          -0.4661896 0.0311246 -14.98   <2e-16 ***
5 I(horsepower^2)      0.0012305 0.0001221 10.08   <2e-16 ***
6
7 Residual standard error: 4.374 on 389 degrees of freedom
8 Multiple R-squared:  0.6876, Adjusted R-squared:  0.686
9 F-statistic: 428 on 2 and 389 DF, p-value: < 2.2e-16
  
```

↑ highly significant

● Multiple linear regression XXXIV

What else can happen? in Multiple Linear Regression.

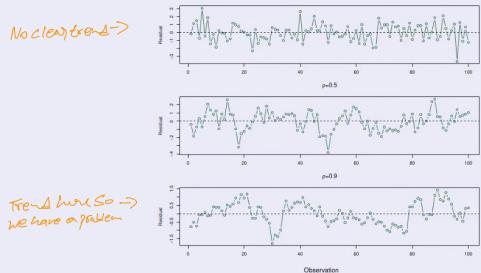
Potential problems

- (a) Non-linearity of the response-predictor relationships. ⇒ Residual plots
- (b) Correlation of error terms. ⇒ Causal inference / good experimental design.
- (c) Non-constant variance of error terms. ⇒ e.g., Weighted Least Squares
- (d) Outliers (extremal y_i). ⇒ detection, winsorization or exclusion.
- (e) High-leverage points (extremal x_i). ⇒ detection, winsorization or exclusion.
- (f) Collinearity. ⇒ Compute VIF, drop/orthogonalize collinear predictors.

↑ i.e., since high correlation means one of the variables is isolated.

● Multiple linear regression XXXVI

Residual plots, simulated data with different $\rho(\varepsilon_i, \varepsilon_{i-1})$ (b)

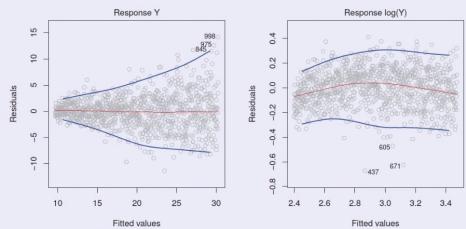


Source: James, Witten, Hastie, and Tibshirani (2013).

Simulated data with different correlations between adjacent points & with adjacent residuals.

● Multiple linear regression XXXVII

Residual plots with Y and $\log(Y)$ as responses *(*C)



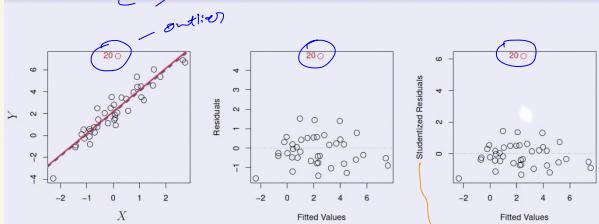
Source: James, Witten, Hastie, and Tibshirani (2013).

For the left plot even though the mean stays constant 0, the variance is increasing. By using $\log(Y)$, the right plot shows variance is constant, so sometimes using \log solves this problem.

● Multiple linear regression XXXVIII

Outliers are extreme values of Y .

Outliers *(*D)



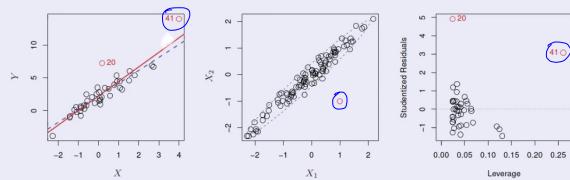
Source: James, Witten, Hastie, and Tibshirani (2013).

Including outlier doesn't change our Linear Model. However, it will lead to increase in RSE & Adjusted R^2 will decrease. So, the Model fit will get worse.

Multiple linear regression XXXIX

High leverage points are extreme values for X .

High leverage points *(*E)

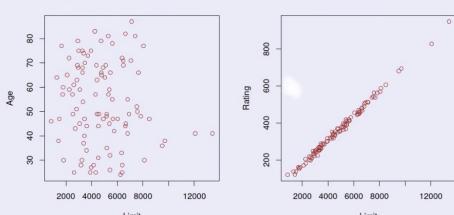


Source: James, Witten, Hastie, and Tibshirani (2013).

High leverage points changes Regression line. One single point can effect the estimation.

● Multiple linear regression XXXX

Collinearity *(*F)



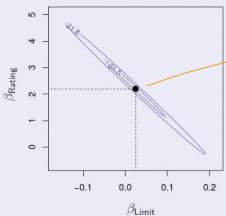
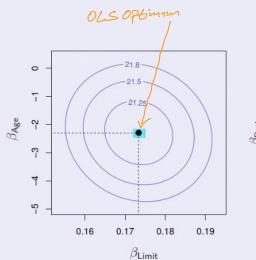
Source: James, Witten, Hastie, and Tibshirani (2013).

left plot has no correlation.

Right plot is strongly positively correlated & including Limit & Rating in Regression model will lead to problems. Coefficient estimates are strongly effected due to Numerical instability.

Multiple linear regression XXXXI

Collinearity



Contour plot here is nasty! is much difficult to find the OLS estimate.

Source: James, Witten, Hastie, and Tibshirani (2013).

Penalised Regression

Penalized regression I

- Recall the problem of model/variable selection.
- One approach: **subset variable selection** (forward/backward stepwise selection).
- Alternative approach: **Shrinkage methods/penalized regression**

Penalized regression

In penalized regression models, *all p* predictors are used in one model and the coefficient estimates are constrained/regularized (in other words, the different algorithms *shrink* the coefficient estimates towards zero).

→ All variables are included.
Most β_i are close to zero but not equal to zero.

Penalized regression II

Ridge Regression

Instead of minimizing

$$RSS = \sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 \quad (\text{OLS})$$

OutEstimate: $\hat{\beta}_0 = \bar{y}_i$

In OLS we minimize RSS.

one minimizes

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p \beta_j^2 = RSS + \lambda \sum_{j=1}^p \beta_j^2 \quad (\text{Ridge Regression}) \quad (18)$$

with $\lambda \geq 0$ being a tuning parameter.

In Ridge Regression we still minimizes RSS but we include the additional component $\lambda \sum_{j=1}^p \beta_j^2$. Note how j starts at 1 so we do not include the intercept term β_0 .

As we minimize this, we are penalizing coefficients (β_j) that are not equal to 0.

so, as some point coefficients tends to 0 it's governed by tuning parameter λ . For $\lambda = 0$, regular OLS.

Ridge Regression I

- As with OLS, Ridge Regression seeks to find a model that fits the data well (i.e., minimizes the RSS).
- The so-called **shrinkage penalty** $\lambda \sum_{j=1}^p \beta_j^2$ is small when the parameters β_j ($j = 1, \dots, p$) are close to zero.
- The parameter λ controls this process of *shrinking* the parameters.

Ridge Regression - coefficients

For each value of the tuning parameter λ , Ridge regression will produce a set of coefficients $\hat{\beta}_{\lambda}^R$.

$\lambda = 0$: $\hat{\beta}_{\lambda}^R = \hat{\beta}^{\text{OLS}}$.

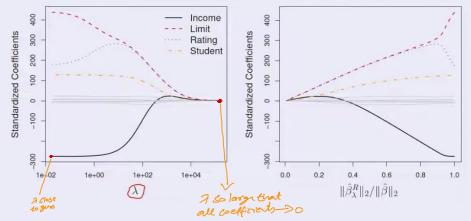
$\lambda \rightarrow \infty$: $\hat{\beta}_{\lambda}^R = 0$.

$\hat{\beta}_{\lambda}^R$ = Ridge Regression coefficient

In Ridge Regression, we repeat it for different values of λ . We get different coefficients for different values of λ .

Ridge Regression II

Ridge Regression - Credit Data



meaning when you rescale variable then the OLS coefficient will be scaled by that factor.

Ridge Regression - Standardized predictors

- OLS: coefficient estimates are scale equivariant.
- Ridge regression: coefficients can change when predictors are multiplied by a constant.
- ⇒ Use standardized predictors!

$$\tilde{x}_{ij} = \frac{x_{ij}}{\sqrt{\frac{1}{n} \sum_{i=1}^n (x_{ij} - \bar{x}_j)^2}}$$

Predictor Values for predictors
sample std.dev

- Standardized predictors have a s.d. of one.

Ridge Regression IV

Ridge Regression - Advantages over OLS

- Ridge regression's advantage over OLS: *bias-variance trade-off*.
- Recall the (*training*) *mean squared error*

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

based on training data observations $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$.

- We are interested in fitting a model that yields a low **test MSE** for a previously unseen test observation (x_0, y_0) .

If we decrease Variance, it would increase Bias & vice versa. Can't reduce both at the same time.

By moving from Blue Curve to green Curve we are decreasing the Bias but increases Variance. (No Bias in Red Curve as it fits well but increases Variance as we run model on a different data set, we get a model that no longer fits the datapoints)

Bias vs. Variance I

Bias-Variance trade-off

One can show that

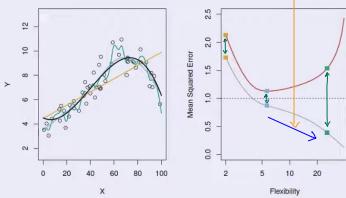
$$\begin{aligned} E(y_0 - \hat{f}(x_0))^2 &= \text{Var}(\hat{f}(x_0)) + [\text{Bias}(\hat{f}(x_0))]^2 + \text{Var}(\varepsilon), \text{ i.e.} \\ \text{expected test MSE} &= \text{Variance of } \hat{f}(x_0) + \text{squared bias of } \hat{f}(x_0) \\ &\quad + \text{variance of the error term.} \end{aligned}$$

To minimize the test MSE, we need to achieve low variance AND low bias.

- Variance** refers to the amount by which \hat{f} would change if we estimated it using a different training data set.
- Bias** refers to the error that is introduced by using a model to approximate for the unknown function f (e.g., using a *linear model*).

Bias vs. Variance II

Bias-Variance trade-off

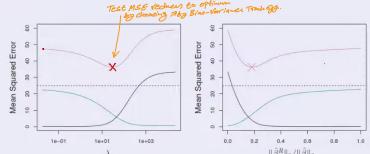


Simulated data (circles); true f : black curve; linear model in orange; smoothing splines in blue and green. Training MSE: grey curve; Test MSE: red curve. Source: James, Witten, Hastie, and Tibshirani (2013).

Ridge Regression V

Ridge Regression - Advantages over OLS

- Ridge regression: as λ increases, the flexibility of the ridge regression fit decreases and the bias increases.



Squared bias (black), variance (green), and test MSE (purple) (simulated data). Source: James, Witten, Hastie, and Tibshirani (2013).

- Thus, we can use Ridge regression to minimize the test mean squared error.

Ridge Regression VI

Ridge Regression - Advantages over OLS

- If $p > n$, then the OLS estimates do not have a unique solution.
- Ridge regression also has **computational advantages** over OLS combined with best subset selection (i.e., estimating 2^p models).
- However: ridge regression will include all **p predictors in the final model!** & almost always the coefficients will all be not equal to zero. Hence, we move on to LASSO.

More predictors (p) than observations (n).

LASSO I

- Ridge regression will always generate a model that includes relevant AND irrelevant predictors. This can make the interpretation of the model difficult as irrelevant predictors will never be excluded. Remedy: the LASSO.

LASSO (least absolute shrinkage and selection operator)

The **LASSO coefficients** $\hat{\beta}_\lambda^L$ minimize the following quantity:

$$\sum_{i=1}^n \left(y_i - \beta_0 - \sum_{j=1}^p \beta_j x_{ij} \right)^2 + \lambda \sum_{j=1}^p |\beta_j| = RSS + \lambda \sum_{j=1}^p |\beta_j| \quad (\text{LASSO}) \quad (19)$$

with $\lambda \geq 0$ being a tuning parameter.

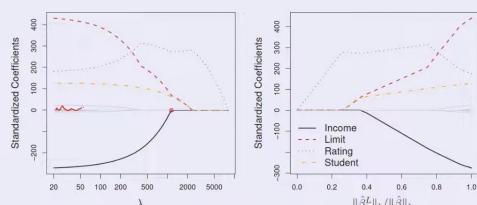
LASSO II

- In short: LASSO uses the ℓ_1 norm whereas Ridge regression uses the ℓ_2 norm to penalize too high coefficients β_j .
- Advantage of LASSO: using the ℓ_1 norm forces some coefficients to be **exactly equal to zero** when λ increases.
- Hence, LASSO performs **variable selection** and fitted models are **sparse** and thus **much easier to interpret**.

Do not include all predictors in the model.

LASSO III

LASSO coefficients - credit data



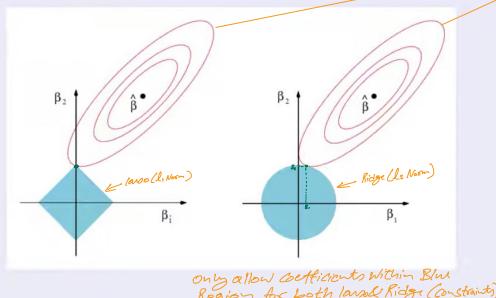
Source: James, Witten, Hastie, and Tibshirani (2013).

In Ridge Regression as λ increases the coefficients slowly tends to 0 but here it goes to 0 very fast.

LASSO IV

Why does LASSO force **some** coefficients to be zero?

Error/MSE and constraint functions



Source: James, Witten, Hastie, and Tibshirani (2013).

- Contour plots of OLS/less.
 OLS coefficients (minimal pt. if we don't have any constraints)
- The intersection of contours with constraint ℓ_1 is the optimal point. So for Ridge we will always get two β_1 & β_2 which are non-zero. In Lasso we can set $\beta_1=0$ & β_2 non-zero. This is why lasso forces **some** of the coefficients to be zero because of shape of blue region.

LASSO V

Which performs best, Ridge or LASSO?

When to use LASSO

- Neither ridge regression nor the **LASSO** will universally dominate the other.
- LASSO: a (rel.) small number of predictors have non-zero coefficients.
- Ridge: response is a function of many predictors, all with coefficients of roughly equal size.

When to USE?

Use Lasso when we expect small no. of predictors to have non-zero coefficients. For ex. with 100 predictors & you expect only 5-10 predictors to have non-zero coefficients.

Elastic Net & Cross-Validation

● Elastic Net

Can one combine LASSO and Ridge Regression? Yes:

Elastic net regularization

The estimates for the **Elastic net regularization** $\hat{\beta}_\lambda^{EN}$ are given by

$$\hat{\beta}_\lambda^{EN} = \arg \min_{\beta} RSS + \lambda_1 \sum_{j=1}^p |\beta_j| + \lambda_2 \sum_{j=1}^p \beta_j^2 \quad (20)$$

with $\lambda_1, \lambda_2 \geq 0$ being tuning parameters.

Combining Lasso & Ridge gives
Elastic Net.

● How to choose the tuning parameters?

How should we **choose the tuning parameters**?

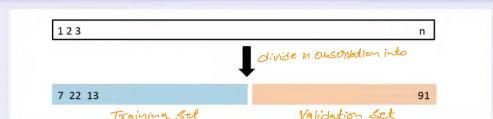
- Recall the distinction between the **test error rate** and the **training error rate**.
- If a sufficiently large test data set is available: training and test error rate are easy to calculate.
- Often not the case (i.e., no or only few test observations are available)!
- Remedy:

Validation Set Approach

Randomly divide the available set of observations into two parts: a training set and a validation/hold-out set. The model is fit on the training set, and the fitted model is used to predict the responses for the observations in the validation set.

● Validation Set Approach

Validation Set Approach



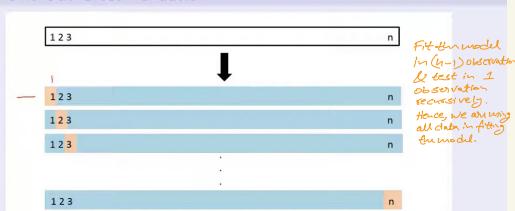
Training set (shown in blue), validation set (shown in beige). Source: James, Witten, Hastie, and Tibshirani (2013).

Drawbacks:

- The validation estimate of the test error rate can be highly variable (which observations are randomly chosen into the training/validation sets?).
- Only a subset of the available data are used for fitting the model. (We are leaving out most data & not using it for training)
(Better to do cross-validation)

● Cross-Validation I To calculate Test Error Rate

Leave-One-Out Cross-Validation

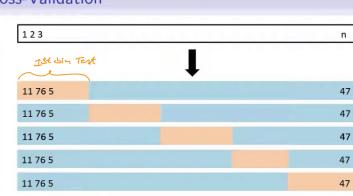


Training set (shown in blue), validation set (shown in beige). The test error is estimated by averaging the n resulting MSEs.

Source: James, Witten, Hastie, and Tibshirani (2013).

Cross-Validation II

k -fold Cross-Validation



Training set (shown in blue), validation set (shown in beige). The test error is estimated by averaging the k resulting MSEs.

Source: James, Witten, Hastie, and Tibshirani (2013).

● Application

- On the following slides we will apply **penalized linear regression** *(Lasso, Ridge & ElasticNet)* models to the task of forecasting equity returns (one day ahead Apple returns).
- Therefore, we rely on the *RiskAnalytics* R package to download some of the data.
- Information on the package and the data can be found [here](#).
- Additionally, we rely on a list of tickers of Nasdaq listed companies which can be downloaded as a csv file from [here](#).

● Application Of PLS (Penalized Linear Regression Models)

● Import and preprocess data I

Import data

```
1 #RiskAnalytics cannot be installed directly from CRAN for
  newer versions of R. Therefore, use the following two
  lines of code:
2   library(devtools)
3   #install_github("lborke/RiskAnalytics")
4
5 library(RiskAnalytics)
6 library(snow)
7 library(quantmod)
8
9 #file path to the downloaded csv file with nasdaq tickers:
10 pathNasdaq<- "data/nasdaqTicker.csv"
11 companyList<-read.csv(pathNasdaq)
12 #order firms in descending order by market capitalization
13 companyList<-companyList[order(companyList$Market.Cap,
  decreasing=TRUE),]
```

↳ loading Packages

● Import and preprocess data III

Variables contained in macroData

```
1 #^VIX: implied volatility index
2 #GSPC: S&P500 returns
3 #IYR: iShares Dow Jones US Real Estate returns
4 #3MTC: Changes in the three-month Treasury bill rate
5 #Yield: Changes in the slope of the yield curve
  corresponding to the yield spread between the ten-year
  Treasury rate and the three month bill rate
6 #Credit: Changes in the credit spread between BAA-rated
  bonds and the Treasury rate
7 #see https://www.wiwi.hu-berlin.de/de/forschung/irtg/results/discussion-papers/discussion-papers-2017-1/sfb649dp2017-003-1.pdf for more details
8
9 colnames(macroData)<-c("VIX", "SP500", "realEstate", "TR3M",
  "yield", "credit") #assign new variable names
```

● Import and preprocess data II

Import data (cont.)

```
1 #download return data for the 10 firms with the highest
  market capitalization within companyList
2 stockData <- as.data.frame(get.yahoo.data(companyList, max-
  comp.num = 10, from.date = "2010-01-01", to.date="2020-12-31")$returns_final)
3
4 macroData <- as.data.frame(get.macro.data(from.date = "
  2010-01-01", to.date = "2020-12-31", plot.macros=FALSE))
5
6 #manually add the corresponding date to each row
7 #for more information on this proceeding see the source code
  at https://github.com/lborke/RiskAnalytics/tree/master/R
8 dates<-rownames(as.matrix(getSymbols(Symbols = "^VIX",
  src = "yahoo", env = NULL, from = "2010-01-01", to =
  "2020-12-31", warnings = FALSE)[, 6]))
9 dates<-dates[-c(1, length(dates))]
10 rownames(macroData)<-dates #assign dates
```

● Import and preprocess data IV

Preprocess data

```
1 macroData$date<-as.Date(rownames(macroData)) #add date
  column
2 stockData$date<-as.Date(rownames(stockData))
3
4 #join return and macro data at the respective date via an
  inner join
5 completeData<-merge(stockData, macroData, by="date")
```

● Import and preprocess data V

Preprocess data (cont.)

```
1 library(dplyr)
2 completeData<-as_tibble(completeData)
3
4 print(completeData,n=5) Print first 5 rows.
5 # A tibble: 2,766 x 17
6   date       AAPL      MSFT      AMZN      GOOG      GOOGL
7   <date>     <dbl>     <dbl>     <dbl>     <dbl>     <dbl>
8   1 2010-01-05  0.00173  3.23e-4  0.00588 -0.00441 -0.00441
9   2 2010-01-06 -0.0160  -6.16e-3  0.0183  -0.0255  -0.0255
10  3 2010-01-07 -0.00185 -1.05e-2  -0.0172 -0.0236  -0.0236
11  4 2010-01-08  0.00663  6.87e-3  0.0267  0.0132  0.0132
12  5 2010-01-11 -0.00886 -1.28e-2  -0.0243 -0.00151 -0.00151
13 # ... with 2,761 more rows, and 11 more variables:
14 #   TSM <dbl>, V <dbl>, JPM <dbl>, NJN <dbl>, WMT <dbl>,
15 #   VIX <dbl>, SP500 <dbl>, realEstate <dbl>, TR3M <dbl>,
16 #   yield <dbl>, credit <dbl>
```

● Import and preprocess data VI

Preprocess data (cont.)

```

1 #Macro variables are extracted from completeData because the
  dates in the original macroData object do not coincide
  with the dates in stockData
2 macroData<-completeData %>% select("VIX","credit") #same
  dates as stockData
3
4 print(round(cor(macroData),3)) #correlation matrix
5 VIX      VIX  SP500  realEstate   TR3M  yield  credit
6 VIX    1.000 -0.163   -0.148 -0.013 -0.026  0.276
7 SP500   -0.163  1.000    0.790  0.026  0.006 -0.015
8 realEstate -0.148  0.790    1.000  0.028  0.006 -0.001
9 TR3M     -0.013  0.026    0.028  1.000  0.049 -0.030
10 yield    -0.026  0.006    0.006  0.049  1.000  0.308
11 credit   0.276 -0.015   -0.001 -0.030  0.308  1.000

```

● Import and preprocess data VIII

Create training and test set

```

1 library(lubridate)
2 #Choose observations before 2019 as training data and
  observations afterwards as test data
3 trainingData<-completeData %>%
4   filter(year(date)<=2018) %>%
5   select(AAPL:GOOGL, JPM, JNJ, VIX:credit)
6 #add one day ahead apple returns to be forecasted later on
7 trainingData$forecast<-c(trainingData$AAPL[-1],0)
8 trainingData<-trainingData[-nrow(trainingData),]
9
10 testData<-completeData %>%
11   filter(year(date)>2018) %>%
12   select(AAPL:GOOGL, JPM, JNJ, VIX:credit)
13 #add one day ahead apple returns in the test set, too
14 testData$forecast=c(testData$AAPL[-1],0)
15 testData<-testData[-nrow(testData),]

```

● Import and preprocess data X

Auxiliary functions for calculation of R^2 and rMSE

- We will evaluate the predictive performance of the various models based on the out-of-sample R^2 as well as the rMSE.
- Therefore, we introduce the following auxiliary functions:

```

1 rSquared<-function(pred, actual){
2   rss<-sum((pred-actual)^2)
3   tss<-sum((actual-mean(actual))^2)
4   rsq<-1-rss/tss
5   return(rsq)
6 }
7
8 rmse<-function(pred, actual){
9   return(sqrt(mean((pred-actual)^2)))
10 }

```

● Simple linear model II

Simple linear model (cont.)

```

1 ...
2 realEstate  0.012901  0.013165  0.980  0.3272
3 TR3M       -0.005056  0.008193  -0.617  0.5372
4 yield       0.002831  0.002025  1.398  0.1623
5 credit      -0.002966  0.002524  -1.175  0.2400
6 ---
7 Signif. codes:  0 '***' '0.001' '**' 0.01 '*' 0.05 '.' 0.1
8
9 Residual standard error: 0.01619 on 2248 degrees of freedom
10 Multiple R-squared:  0.007153, Adjusted R-squared:
  0.001412
11 F-statistic: 1.246 on 13 and 2248 DF, p-value: 0.2397

```

Both Multiple Adjusted R^2 is very small.
Linear Model is Not So Good.

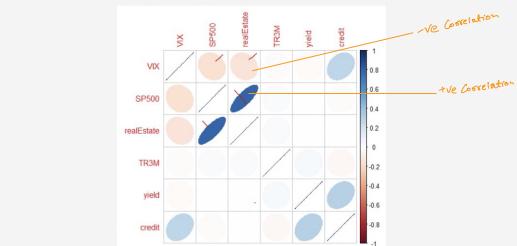
● Import and preprocess data VII

Visualization of pairwise correlations

```

1 library(corrplot)
2 corrplot(cor(macroData), method="ellipse")

```



● Import and preprocess data IX

Create training and test set

```

1 #Some of the models need dependent and explaining variables
  separately in matrix form
2
3 xTrain<-as.matrix(select(trainingData, -forecast)) #remove
  forecast from trainingData
4 yTrain<-trainingData$forecast
5
6 xTest<-as.matrix(select(testData, -forecast))
7 yTest<-testData$forecast

```

● Simple linear model I (OLS)

Simple linear model

```

1 # forecast future Apple returns with past returns of the
  stock itself and various other variables
2 linModel<-lm(forecast~., data=trainingData)
3
4 summary(linModel)
5 Coefficients:
6 Estimate Std. Error t value Pr(>|t|)
7 (Intercept)  -0.018068  0.013017 -1.388  0.1653
8 AAPL         0.036055  0.026788  1.346  0.1785
9 MSFT        -0.052560  0.033369 -1.575  0.1154
10 AMZN        -0.020489  0.022059 -0.929  0.3531
11 GOOG        -0.100961  0.0228608 -0.442  0.6588
12 GOOGL       0.118844  0.0228915  0.519  0.6037
13 JPM          -0.031266  0.0348800 -0.896  0.3702
14 JNJ          -0.085395  0.0486111 -1.757  0.0791
15 VIX          0.006921  0.005037  1.374  0.1696
16 SP500        0.020498  0.024484  0.837  0.4026
17 ...

```

JNJ stock return has slightly more statistical significance in explaining future Apple stock return.

● Simple linear model III

Out-of-sample performance

```

1 pred<-predict(linModel, newdata=testData)
2
3 rSquared(pred=pred, actual=yTest)
4   -0.01376846
5
6 rmse(pred=pred, actual=yTest)
7   0.02405451

```

Model doesn't perform better in Test data & R^2 was because -ve

● Choosing a value for the tuning parameters I

- The penalized linear regression models rely on one (ridge, LASSO) or two (elastic net) tuning parameters that control the amount of regularization.
- We want to choose the tuning parameter such that it generalizes well to new data.
- This will be done via cross-validation.

● Ridge regression I

Determining the regularization parameter

```

1 library(glmnet)
2 set.seed(2021) # for reproducibility
3
4 #Find optimal lambda parameter via cross-validation:
5 cvRidge<-cv.glmnet(x=xTrain, yTrain, nfolds=10, nlambda=100,
6   alpha=0, lambda.min.ratio=0.0001)
7 #nfolds: number of folds in cross validation
8 #nlambda: number of regularization parameters to be tested
9 #alpha: determines the weighting to be used, 0 in case of
10 ridge regression, 1 gives LASSO
11 #lambda.min.ratio: ratio of smallest lambda to largest
12 #lambda considered
13 #for further details see https://glmnet.stanford.edu/
14   articles/glmnet.html for further details
15
16 optimalLambda<-cvRidge$lambda.min
17 optimalLambda
18 0.455634 ← optimal λ

```

● Ridge regression II

Fitting the model

```

1 #fit ridge regression model with the lambda determined on
2   the previous slide
3 ridgeReg<-glmnet(x=xTrain, yTrain, alpha=0, lambda=
4   optimalLambda)
5
6 coef(ridgeReg) #coefficients of the penalized linear model
7 (Intercept) 6.308607e-04
8 AAPL 9.831866e-04
9 MSFT -8.763122e-04
10 AMZN -3.372965e-04
11 GOOG 2.089550e-04
12 GOOGL 2.573408e-04
13 ...
14 TR3M -1.985538e-04
15 yield 9.563131e-05
16 credit -7.623069e-06

```

● Ridge regression IV

Effect of high/low lambda (cont.)

Using low λ we get some
asols coefficients

```

1 coef(glmnet(x=xTrain, yTrain, alpha=0, lambda=lowLambda))
2
3 (Intercept) -0.018050381
4 AAPL 0.036021284
5 MSFT -0.052962852
6 AMZN -0.020486622
7 GOOG -0.057921967
8 GOOGL 0.076092514
9 ...
10 TR3M -0.005064605
11 yield 0.002834675
12 credit -0.002962296

```

● Choosing a value for the tuning parameters II

Cross-validation

- To perform k -fold cross-validation, we split the training data into k parts or "folds".
- One usually chooses $k = 5$ or $k = 10$.
- k -fold cross-validation then considers training on all but the k th part and then validating the model performance on that k th part, iterating over k .
- Cross-validation cannot only be used for selecting hyperparameters, but also for assessing the performance of a model.



● Ridge regression III

Effect of high/low lambda

Asymptotically we have:

- $\lambda \rightarrow 0 : \beta_{\text{ridge}} \rightarrow \beta_{\text{OLS}}$
- $\lambda \rightarrow \infty : \beta_{\text{ridge}} \rightarrow 0$

```

1 highLambda<-10
2 lowLambda<-0.0001
3 coef(glmnet(x=xTrain, yTrain, alpha=0, lambda=highLambda))
4 (Intercept) 7.729305e-04
5 AAPL 4.554773e-05
6 MSFT -4.106883e-05
7 AMZN -1.584365e-05
8 GOOG 9.446636e-06
9 GOOGL 1.164191e-05
10 ...
11 TR3M -9.432911e-06
12 yield 4.532762e-06
13 credit -2.308404e-07

```

↑
Using high λ = 10

● Ridge regression V

In and out-of-sample performance

```

1 #in-sample R^2
2 predTrain<-predict(ridgeReg, newx=xTrain)
3 rSquared(pred=predTrain, actual=yTrain)
4 0.0003141296 → Much lower than OLS model, so better
5 #reminder: in-sample R^2 simple linear model = 0.007153
6
7 #out-of-sample performance
8 predTest<-predict(ridgeReg, newx=xTest)
9 rSquared(pred=predTest, actual=yTest)
10 -0.005109216 → slightly better than OLS.
11 #reminder: out-of-sample R^2 simple linear model = -0.0138
12 rmse(pred=predTest, actual=yTest)
13 0.02395156
14 #reminder: out-of-sample RMSE simple linear model = 0.0241

```

● Ridge regression VI

Summary

- We determined the hyperparameter λ via 10-fold cross-validation.
- Using this λ we obtained a linear model with coefficients that are smaller (in terms of their absolute value) than the corresponding coefficients in the simple OLS model. The degree of penalization of large coefficients is governed by λ .
- The OLS model achieves by construction the highest **in-sample R^2** . However, the ridge regression model might generalize better to previously unseen data (**out-of-sample**). This is true in our application where the ridge regression model achieved a (slightly) higher out-of-sample R^2 and a lower MSE compared to the OLS model.

Since the parameter λ governs the training process it is called **hyperparameter** instead of **parameter**.

● LASSO regression I

Determining the regularization parameter

```

1 library(glmnet)
2 set.seed(2021) # for reproducibility
3
4 cvLASSO<-cv.glmnet(x=xTrain, yTrain, nlambda=100, alpha=1,
                      lambda.min.ratio=0.0001)
5 #parameters are the same as for ridge regression except for
#alpha
6 #alpha: determines the weighting to be used, 0 for ridge and
# 1 for LASSO regression
7
8 optimalLambda<-cvLASSO$lambda.min
9 optimalLambda
10 0.0005000578

```

● LASSO regression III

In and out-of-sample performance

```

1 #in-sample R^2
2 predTrain<-predict(LASSOReg, newx=xTrain)
3 rSquared(pred=predTrain, actual=yTrain)
4 0
5
6 #out-of-sample performance
7 predTest<-predict(LASSOReg, newx=xTest)
8 rSquared(pred=predTest, actual=yTest)
9 -0.005150357
10 #reminder: out-of-sample R^2 simple linear model = -0.0138
11 rmse(pred=predTest, actual=yTest)
12 0.02395205
13 #reminder: out-of-sample RMSE simple linear model = 0.0241

```

● Elastic net I

Determining the regularization parameters

```

1 library(caret)
2 set.seed(2021) #for reproducibility
3 #specify trainingControl object that specifies how repeated
#cross validation will take place
4 trainingControl<-trainControl(method="repeatedcv", number=5,
                                 repeats=10, search="random")
5
6 #builds the model in which a range of possible alpha and
#lambda values are tested and the optimum values are
#selected.
7 #tuneLength specifies the number of different alpha and
#lambda combinations to be tested
8 elasticNet<-train(forecast~., data=trainingData, method="glmnet",
                     tuneLength=10, trControl=trainingControl)
9
10 elasticNet$bestTune #Optimal parameters
11      alpha    lambda
12      0.8154215 5.488304

```

● LASSO regression II

Fitting the model

```

1 #fit LASSO regression model with the lambda determined on
#the previous slide to the training data
2 LASSOReg<-glmnet(x=xTrain, yTrain, alpha=1, lambda=
#optimalLambda)
3
4 coef(LASSOReg)
5 (Intercept) 0.0007795497
6 AAPL 0.0000000000
7 MSFT .
8 AMZN .
9 GOOG .
10 GOOGL .
11 ...
12 TR3M .
13 yield .
14 credit .

```

● Elastic net II

In and out-of-sample performance

```

1 #in sample predictions based on the optimal parameters from
#the previous slide:
2 predTrain<-predict(elasticNet, xTrain)
3 rSquared(pred=predTrain, actual=yTrain) #in-sample R^2
4 0
5
6 #out-of-sample predictions based on the optimal parameters
#from the previous slide:
7 predTest<-predict(elasticNet, xTest)
8 rSquared(pred=predTest, actual=yTest)
9 -0.005150357
10 #reminder: out-of-sample R^2 simple linear model = -0.0138
11 rmse(pred=predTest, actual=yTest)
12 0.02395205
13 #reminder: out-of-sample RMSE simple linear model = 0.0241

```

Bayes & KNN Classifier

● Model accuracy - classification I

Recall the classification problem. How can we assess the model accuracy of classifiers (in contrast to regression methods)?

Training error rate (qualitative response)

The **training error rate**, i.e. the proportion of mistakes that are made if we apply our estimate \hat{f} to the training observations, is defined as

$$\frac{1}{n} \sum_{i=1}^n I(y_i \neq \hat{y}_i) \quad \text{Indicator function}\ I(y_i \neq \hat{y}_i) \text{ equals 1 when } y_i \neq \hat{y}_i \quad (21)$$

with \hat{y}_i being the predicted class label for the i th observation.

● Bayes classifier (Most Basic Classifier)

Bayes classifier

The **Bayes classifier** is the classifier that **assigns each observation to the most likely class, given its predictor values**: Assign a test observation with predictor vector x_0 to the class j for which the cond. probability

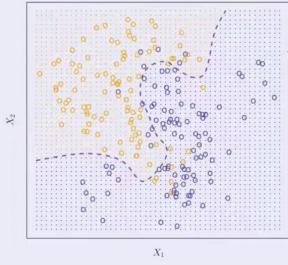
$$Pr(Y = j | X = x_0)$$

is largest.

- The test error rate given is minimized, on average, by the Bayes classifier.
- Problem: Distribution of Y given X is not known! So only good in theory but not in practice as conditional prob. not known.

● k-nearest neighbors II

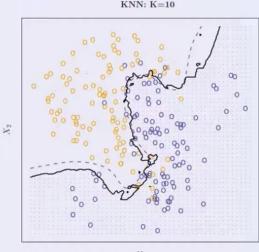
Classification - Bayes decision boundary



Simulated data with two classes. Bayes decision boundary in purple. Source: James, Witten, Hastie, and Tibshirani (2013).

● k-nearest neighbors IV

kNN vs. Bayes decision boundaries



Source: James, Witten, Hastie, and Tibshirani (2013).

● Model accuracy - classification II

How can we assess the model accuracy of classifiers **when applied to new data?**

Test error rate (qualitative response)

The **test error rate**, i.e. the proportion of mistakes that are made if we apply our estimate \hat{f} to the **test** observations (x_0, y_0) , is defined as

$$\text{Ave}(I(y_0 \neq \hat{y}_0)), \quad (22)$$

i.e., the percentage of incorrectly classified test observations.

● k-nearest neighbors I (use this since Bayes can't be used in practice)

k-nearest neighbors (kNN) classifier

Starting point: a positive integer k and a test observation x_0
Procedure:

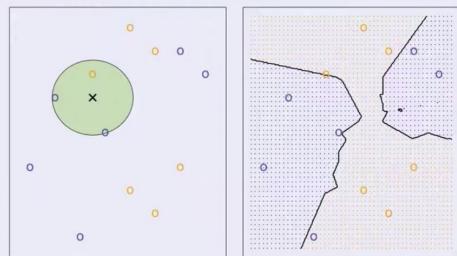
- Identify the k points in the training data that are closest to x_0 , represented by \mathcal{N}_0 .
- Estimate the conditional probability for class j as the fraction of points in \mathcal{N}_0 whose response values equal j :

$$Pr(Y = j | X = x_0) = \frac{1}{k} \sum_{i \in \mathcal{N}_0} I(y_i = j). \quad (23)$$

- Apply Bayes rule and classify the test observation x_0 to the class with the largest probability.

● k-nearest neighbors III

kNN - Simulated data $k = 3$

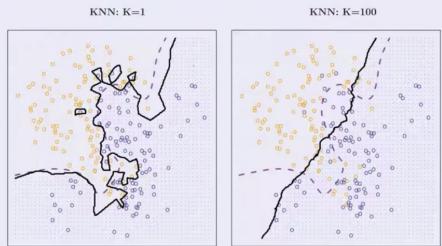


Source: James, Witten, Hastie, and Tibshirani (2013).

In Finance, Classification is normally used in Credit Risk Management. For eg. Classify Good & Bad loans, defaulting Customers vs Nondefaulters, Customers prone to switch insurance.

● k-nearest neighbors V

kNN decision boundaries for $k = 1$ and $k = 100$



Source: James, Witten, Hastie, and Tibshirani (2013).

K is a hyperparameter. Needs to be chosen before applying algorithm via cross-validation.
 $K=1 \Rightarrow$ More wiggly
 $K=100 \Rightarrow$ Very coarse even (linear) decision boundary.
 Bias-Variance tradeoff in play.

● Maximal Margin Classifier

● Support Vector Machines I

- **SVMs:** usually the summary of three distinct methods: **the maximal margin classifier**, **the support vector classifier**, and **the support vector machine**.
- Like kNN, it is used for **classification**.

Hyperplane (in p -dimensional)

A **hyperplane** is a flat affine (linear) subspace of dimension $p - 1$.
 $p = 2$: a line; $p = 3$: a plane. A hyperplane divides the p -dimensional space into two halves and is defined by the equation

$$\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_p X_p = 0 \quad \text{Hyperplane Eq. 1}$$

● Support Vector Machines III

- Observations
- Observations in
- **Input:** $n \times p$ data matrix X of n p -dimensional training observations

$$x_1 = \begin{pmatrix} x_{11} \\ \vdots \\ x_{1p} \end{pmatrix}, \dots, x_n = \begin{pmatrix} x_{n1} \\ \vdots \\ x_{np} \end{pmatrix}$$

with each observation belonging to one class, i.e., $y_1, \dots, y_n \in \{-1; 1\}$
 where -1 represents one class and 1 the other class, and a test observation with a p -vector of observed features $x^* = (x_1^*, \dots, x_p^*)^T$

- **Output:** classification of x^* using a **separating hyperplane**.

● Support Vector Machines V

Problem: if **one** separating hyperplane exists, then an **infinite** number of such hyperplanes exists. Which one should we use?

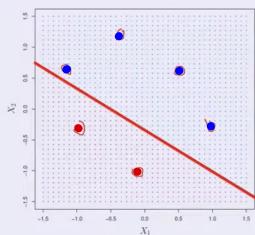
Maximal Margin Classifier

Natural choice: the **maximal margin hyperplane**, which is the separating hyperplane that is farthest from the training observations.

- Compute the (perpendicular) distance from each training observation to a given separating hyperplane.
- The smallest such distance is known as the margin.
- The maximal margin hyperplane is the separating hyperplane for which the margin is largest. S.t. distance from blue & red pts are equal

● Support Vector Machines II

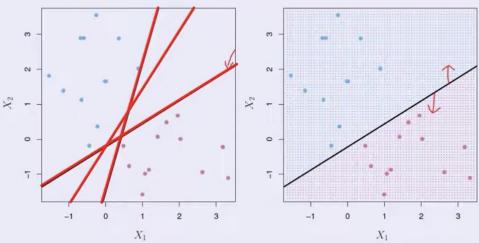
Hyperplane



Hyperplane defined by $1 + 2X_1 + 3X_2 = 0$. Source: James, Witten, Hastie, and Tibshirani (2013).

● Support Vector Machines IV

Classifier based on a separating hyperplane

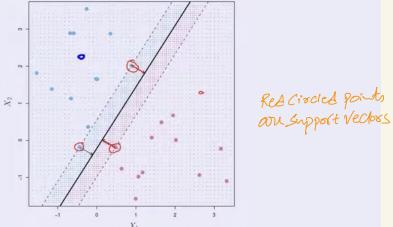


Source: James, Witten, Hastie, and Tibshirani (2013).

Maximum Margin classifier only depends on points called **Support vectors** so adding points that are not support vectors doesn't change Maximal margin hyperplane.

● Support Vector Machines VI

Maximal margin hyperplane/classifier

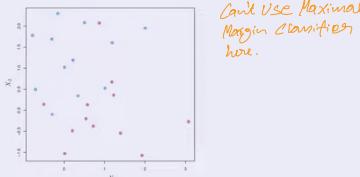


The observations on the dashed lines that are equidistant from the separating hyperplane are the **support vectors**. Source: James, Witten, Hastie, and Tibshirani (2013).

● Support Vector Machines VIII

If a **separating hyperplane** exists, use the Maximal Margin Classifier. But is this always the case? No.

Non-separable case



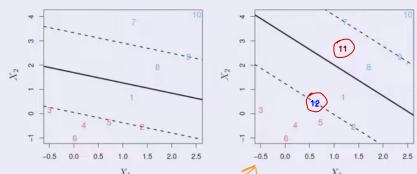
Source: James, Witten, Hastie, and Tibshirani (2013).

Support Vector Machines

● Support Vector Machines X

Solution: Using a classifier that allows for some observations to be incorrectly classified.

Support Vector Classifier



Source: James, Witten, Hastie, and Tibshirani (2013)

With addition of two points 11 (red) & 12 (blue), the decision boundary should change that drastically by using Support Vector Classifier. 11 is incorrectly classified.

● Support Vector Machines VII

Construction of the Maximal Margin Classifier

The **Maximal Margin hyperplane** is the solution to the optimization problem

$$\arg \max_{\beta_0, \beta_1, \dots, \beta_p, M} M \quad (24)$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1, \quad (25)$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M, \forall i = 1, \dots, n. \quad (26)$$

Constraint (26): ensures that each observation will be on the correct side of the hyperplane (with some buffer/margin M).

Constraint (25) ensures that each observation is at least a distance M from the hyperplane.

$$y_i = 1 \text{ or } -1$$

● Support Vector Machines IX

If a **separating hyperplane** exists, should we always use a classifier based on it? No. Since it's very sensitive to new observations

Sensitivity of the Maximal Margin Classifier



Source: James, Witten, Hastie, and Tibshirani (2013).

● Support Vector Machines XI

Construction of the Support Vector Classifier

The **Support Vector hyperplane** is the solution to the problem

$$\arg \max_{\beta_0, \beta_1, \dots, \beta_p, \epsilon_1, \dots, \epsilon_n, M} M \quad (27)$$

$$\text{subject to } \sum_{j=1}^p \beta_j^2 = 1, \quad (28)$$

$$y_i(\beta_0 + \beta_1 x_{i1} + \beta_2 x_{i2} + \dots + \beta_p x_{ip}) \geq M(1 - \epsilon_i), \quad (29)$$

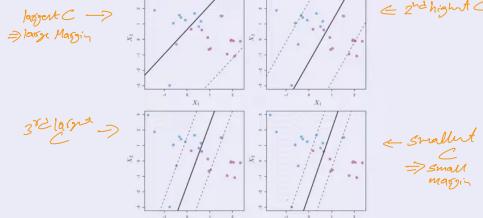
$$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C, \quad (30)$$

with C being a nonnegative tuning parameter chosen via cross-validation. M is the width of the margin that we want to maximize. ϵ_i are so-called **slack variables** that allow observations to be on the wrong side of the margin or hyperplane.

When C is larger, there is a high tolerance for observations to be incorrectly classified. Hence, Margin will be large.

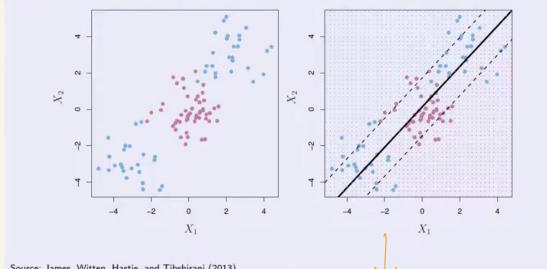
Support Vector Machines XII

Support Vector Classifier - Influence of the tuning parameter C



Support Vector Machines XIII

Is a linear classifier always warranted?



Requires Non-linear classification.

Support Vector Machines XIV

The Support Vector Classifier is a *linear* classifier. How can we (automatically) convert it to a classifier with *non-linear decision bounds*?

Classifier with non-linear decision bounds

$$\arg \max_{\beta_0, \beta_{11}, \beta_{12}, \dots, \beta_{p1}, \beta_{p2}, \epsilon_1, \dots, \epsilon_n, M} M \quad (31)$$

$$\text{subject to } y_i \left(\beta_0 + \sum_{j=1}^p \beta_{j1} x_{ij} + \sum_{j=1}^p \beta_{j2} x_{ij}^2 \right) \geq M(1 - \epsilon_i), \quad (32)$$

$$\epsilon_i \geq 0, \sum_{i=1}^n \epsilon_i \leq C, \sum_{j=1}^p \sum_{k=1}^2 \beta_{jk}^2 = 1. \quad (33)$$

Previously we had a linear hyperplane, we now allow a polynomial fn. This makes decision boundary to be polynomial.

This is one way of doing it, however it might not be enough & we need more Non-linearity.

Support Vector Machines XV

Recall the standard (or euclidean) inner product:

$$\langle x_i, x_{i'} \rangle = \sum_{j=1}^p x_{ij} x_{i'j}.$$

Linear Support Vector Classifier - Alternate representation

The linear Support Vector Classifier can be represented as

$$\text{hyperplane} \rightarrow f(x) = \beta_0 + \sum_{i=1}^n \alpha_i \langle x, x_i \rangle \quad (34)$$

with n parameters α_i . To estimate these and β_0 , we only need the $n(n-1)/2$ inner products $\langle x, x_i \rangle$ between all pairs of training observations. Moreover, α_i is nonzero only for the support vectors in the solution (i.e., usually a much smaller number of parameters will suffice).

Support Vector Machines XVI

Kernel - examples

Linear Kernel:

$$K(x_i, x_{i'}) = \sum_{j=1}^p x_{ij} x_{i'j}. \quad (35)$$

two vectors x_i & $x_{i'}$

Polynomial Kernel of degree d :

$$K(x_i, x_{i'}) = \left(1 + \sum_{j=1}^p x_{ij} x_{i'j} \right)^d \quad (36)$$

Radial Kernel:

$$K(x_i, x_{i'}) = \exp \left(-\gamma \sum_{j=1}^p (x_{ij} - x_{i'j})^2 \right), (\gamma > 0). \quad (37)$$

Support Vector Machines XVI

Why do we need a different representation of SV classifier?
⇒ Because we can extend it with kernel.

How can we extend the linear SV classifier? Replace the inner product with a more general Kernel.

Kernel

A Kernel K (in the context of machine learning) is a function that quantifies the similarity of two observations.

← different representation of SV classifier.

If we have two observations x_i & $x_{i'}$, a kernel measures the similarity between x_i & $x_{i'}$ (could be distance between two points). In case of inner product it measures distance.

After extending SV classifier with Kernel (& substituting linear kernel with Non-linear kernel) we get Support Vector Machine.

Support Vector Machines XVIII

Support Vector Machine

When using a *non-linear kernel* in combination with the Support Vector Classifier, the resulting classifier is known as a **Support Vector Machine**.

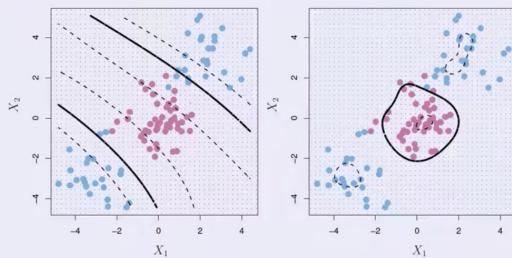
$$f(x) = \beta_0 + \sum_{i \in S} \alpha_i K(x, x_i) \quad (38)$$

\leftarrow we don't need all n parameters as α_i will be non-zero only for support vectors.

with S being the set of support vectors (i.e., the training observations for which α_i is nonzero).

● Support Vector Machines XIX

Support Vector Machines - Polynomial and Radial Kernels



Source: James, Witten, Hastie, and Tibshirani (2013).

← Much finer way of classification with SVM.

Application Of SVMs

● Application I

Which customers are prone to terminate their contract.

- On the following slides we will apply the kNN model and SVMs to a dataset of credit card customers to predict customer churn.
- The data are available at [kaggle](#). There you can find various (Python) notebooks featuring a wide variety of different methods for accomplishing this task.
- However, in this application we only use the kNN model and SVMs ([linear](#) and with a radial basis function (RBF) kernel).

● Application III

The dataset

- The dataset contains information on more than 10,000 customers including age, salary, credit card limit, and other features.
- The data are [unbalanced](#) with only about 16 % of the customers having churned. This complicates training and interpretation of the predictive performance of the models.
- It might be more favorable to contact a customer who is not about to churn than not contacting a customer who is.
- In this lecture, our primary focus is on the models. We will therefore not specifically address the class imbalance problem.

● Application II

Motivation

- The manager responsible for the credit card customers is worried about an increasing fraction of these customers quitting their contracts.
- To slow down customer churn, she wants to contact these customers proactively in order to change their decision.
- Therefore, she needs a prediction of who is going to terminate the contract.

Class imbalance → Many more instances of some classes than others.

● Application IV

For fitting the kNN model and the SVMs (as well as the elastic net model from the previous section) we rely on the [caret R-package](#) by Kuhn (2020). Before we work with the data we will therefore shortly introduce the package.

caret

- Caret is short for [classification and regression training](#).
- It [summarizes activities](#) related to model development in a [streamlined process](#). It allows you to test different models with very little changes to the code and offers nearly automatic cross validation and parameter tuning.
- The package provides a [consistent modeling syntax](#). For example, by simply changing the method argument one can easily change the underlying model. In total, the package provides access to [more than 200 different models](#).

● Application V

caret (cont.)

- Note that behind the scenes the package is [not performing the modeling itself](#). Instead, it relies on various other specialized R-packages.
- For example, when using the method `lm` caret uses the `lm()` function from the `stats` package to estimate a linear regression model.
- When performing kNN classification, caret relies on the `class` package by Venables and Ripley (2002).
- For a complete list of models caret provides access to see [here](#).

● Import and preprocess data I

Import data

```

1 # download bankChurners.csv from here: https://www.kaggle.com/
  /sakshigoyal7/credit-card-customers/download
2 #then import csv-file
3 library(readr)
4 bankChurners <- read_csv("data/BankChurners.csv")
5 #drop two last columns (according to dataset description at
  https://www.kaggle.com/sakshigoyal7/credit-card-
  customers)
6 bankChurners<-bankChurners[,1:(ncol(bankChurners)-2)]

```

● Import and preprocess data III

Variable description

- CLIENTNUM: unique identifier for the customer holding the account
- Attrition_Flag: if the account is closed then 1 else 0 (\rightarrow Customer Churned/Not)
- Dependent count: number of dependents (Count of age etc)
- ...
- A description of all variables can be found [here](#).

● Import and preprocess data V

Explorative analysis (cont.)

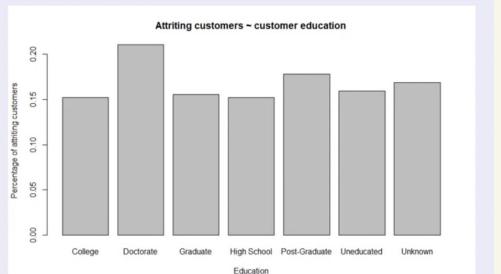
```

1 # Analyze churning customers by age
2 library(dplyr)
3 percAgeGrouped<-bankChurners %>%
4   mutate(interval=as.numeric(cut(Customer.Age, breaks=10*
  (2:8)))) %>%
5   group_by(interval) %>%
6   summarize(percentageAttrited=mean(Attrition_Flag==""
  Attrited Customer"))
7
8 plot(x=10*(2:7), y=percAgeGrouped$percentageAttrited, xlab="Customer age", ylab="Percentage of attriting customers",
  type="b")

```

● Import and preprocess data VII

Explorative analysis (cont.)



In this case since doctorate influences attriting customers, we could use one dummy variable & encode one for doctorate & zero for rest.



● Import and preprocess data II

Explorative analysis

```

1 print(bankChurners,n=7)
2
3 # A tibble: 10,127 x 21
4   CLIENTNUM Attrition_Flag   Customer_Age Gender
5   <dbl> <chr>           <dbl> <chr>
6   1 768805383 Existing Customer    45 M
7   2 818770008 Existing Customer    49 F
8   3 713982108 Existing Customer    51 M
9   4 769911858 Existing Customer    40 F
10  5 709106358 Existing Customer    40 M
11  6 713061558 Existing Customer    44 M
12  7 810347208 Existing Customer    51 M
13 # ... with 10,120 more rows, and 17 more variables:

```

● Import and preprocess data IV

Explorative analysis (cont.)

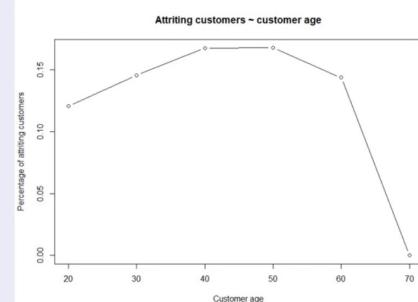
```

1 any(is.na(bankChurners))  $\rightarrow$  Check for missing values
2 FALSE # no missing values
3
4 any(duplicated(bankChurners))  $\rightarrow$  Check for duplicates
5 FALSE # no duplicates
6
7 # class imbalance  $\rightarrow$  Checking class imbalance but we are not solving this in this example
8 round(table(bankChurners$Attrition_Flag)/nrow(bankChurners),
  3)
9
10 Attrited Customer Existing Customer
11 0.161          0.839
 $\downarrow 16\% \text{ Customer Churned}$   $\uparrow 83\% \text{ Remaining Customer}$ 

```

● Import and preprocess data VI

Explorative analysis (cont.)



● Import and preprocess data VIII

Categorical features

- Customer education is a **categorical feature** and cannot simply be encoded in a numeric variable as is sometimes possible for ordinal variables.
- Customer education might be encoded via 6 binary **dummy variables**. However, this raises the dimensionality of the data significantly. This is an especially serious problem for non-parametric methods like kNN classification because of the **curse of dimensionality**.
- For simplicity, we therefore drop all non-numeric variables. In practice, one should carefully decide for each (categorical) variable if it should be included or not.

● Import and preprocess data IX

Create training and test set

- We remove all non-numeric features and the column CLIENTNUM (no predictive value).
- Furthermore, we create the training and test set by randomly including 80 % of observations in the training and 20 % of observations in the test set.

```
1 set.seed(2021) #for reproducibility
2 index<-sample.int(n=2, size=nrow(bankChurners), prob=c
  (0.8,0.2), replace=TRUE)
3 table(index)
4      1     2
5 8061 2066
6
7 xTrain<-bankChurners %>%
8   select (-c("CLIENTNUM", "Attrition_Flag")) %>%
9   select(where(is.numeric)) %>% filter(index==1)
```

● Import and preprocess data XI

Create training and test set (cont.)

By default, kNN is based on the **Euclidean distance**. To ensure that all features contribute equally to the measured distances we **scale the data** based on the minimum and maximum values of the training data.

```
1 minVec<-apply(xTrain,2,min) # determine per-column minimum
2 maxVec<-apply(xTrain,2,max) # determine per-column maximum
3
4 # scale variables
5 xTrainScaled<-as.data.frame(t(apply(xTrain,1,function(row){(
  row-minVec)/(maxVec-minVec)})))
6
7 xTestScaled<-as.data.frame(t(apply(xTest,1,function(row){(
  row-minVec)/(maxVec-minVec)})))
```

● Import and preprocess data XII

Create training and test set (cont.)

All variables in the **training set** now have minimum 0 and maximum 1. As scaling is based on data from the training set, the same is not necessarily true for the test set:

+

```
1 # calculate maximum per column --> apply()
2 # then provide summary statistics and round results to four
  decimal places
3 round(summary(apply(X=xTrainScaled, MARGIN=2, FUN=max)),4)
4   Min. 1st Qu. Median Mean 3rd Qu. Max.
5       1       1       1       1       1
6
7 round(summary(apply(X=xTestScaled, MARGIN=2, FUN=max)),4)
8   Min. 1st Qu. Median Mean 3rd Qu. Max.
9       0.7741  0.9614  1.0000  0.9633  1.0000  1.0000
```

● kNN classification I

Introduction

- For our kNN classification model we rely on the Euclidean distance. Many other **distance metrics** are possible (e.g., Manhattan distance, cosine distance, etc.).
- The **number of neighbors** is a hyperparameter. We determine an appropriate value via 10-fold cross validation.
- Computations can be performed in **parallel** to speed up the estimation process.
- We determine the optimal parameter and evaluate a models' predictive performance based on the metrics **accuracy** and **Cohen's kappa** (more details later).

● Import and preprocess data X

Create training and test set (cont.)

```
1 xTest<-bankChurners %>%
2   select (-c("CLIENTNUM", "Attrition_Flag")) %>%
3   select(where(is.numeric)) %>%
4   filter(index==2)
5
6 yTrain<-as.factor(bankChurners$Attrition_Flag[index==1])
7 yTest<-as.factor(bankChurners$Attrition_Flag[index==2])
```

bankChurners & Attrition_Flag is our response variable.

kNN classification II

Hyperparameter tuning

```

1 library(caret)
2 trControl <- trainControl(method = "cv", number = 10)
3 set.seed(2021) #for reproducibility
4
5 # perform cross validation in parallel
6 # if you encounter any problems with the PSOCK-cluster, the
7 # corresponding lines can be commented out.
8 library(doParallel)
9 cl <- makePSOCKcluster(6) #number of cores/threads to use
10 registerDoParallel(cl)
11
12 kNNModel<-train(x=xTrainScaled, y=yTrain, method="kNN",
13 tuneGrid=expand.grid(k=1:10), trControl=trControl,
14 metric="Accuracy") # consider 1, 2, ..., 10 neighbors
15
16 stopCluster(cl)

```

cross-validation

tuning parameters

kNN classification IV

Predictive performance

```

1 # predict class labels for the test set
2 yPredKnn<-predict(kNNModel,newdata=xTestScaled)
3
4 # print confusion matrix (see the next slide) confMatrix
5 print(confusionMatrix(data=yPredKnn, reference=yTest))

```

kNN classification VI

Predictive performance (cont.)

- **Accuracy** is defined as the percentage of correctly classified observations. In our example, we obtain an accuracy of

$$90.08 \% = \frac{189 + 1672}{189 + 1672 + 53 + 152} = \frac{1861}{2066}.$$

- **No Information Rate (NIR)** is defined as the largest proportion of the observed classes. In our example this corresponds to the proportion of existing customers

$$83.49 \% = \frac{53 + 1672}{2066}.$$

NIR is the highest accuracy which can be achieved by a constant prediction.

kNN classification VIII

Predictive performance (cont.)

- When comparing two models, a higher kappa signals a better predictive performance. The maximum value is 1.
- There is, however, no standardized way of interpreting its values.
- A negative value of kappa signals that the model's predictions are worse than predicting by chance.

kNN classification III

Model summary

```

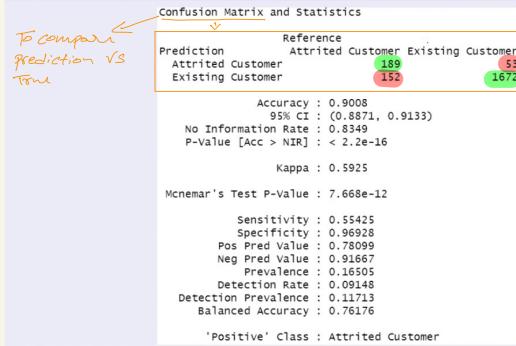
1 print(kNNModel) #optimal number of neighbors is the same for
2 both accuracy and Cohen's kappa. In general, this is
3 not necessarily the case.
4
5 Resampling results across tuning parameters:
6   k Accuracy Kappa
7   1  0.8988959  0.6070618
8   ...
9   6  0.9131598  0.6398903
10  7  0.9160140  0.6506041
11  8  0.9150200  0.6423632
12  ...
13 10  0.9125397  0.6288506
14
15 Accuracy was used to select the optimal model using the
16 largest value.
17
18 The final value used for the model was k = 6.

```

kNN classification V

Bad for prediction Accuracy
Good for prediction Accuracy

Predictive performance (cont.)



kNN classification VII

Predictive performance (cont.)

- **Cohen's kappa** is a measure of a classifier's performance relative to how well the model would have performed simply by chance. Therefore, one compares the accuracy of the model to the hypothetical probability of an agreement by chance.
- In our example, 83.49 % (NIR) of observations belong to existing customers while 16.51 % belong to attrited customers. The kNN model classifies $\frac{152+1672}{2066} = 88.29\%$ as existing and $\frac{189+53}{2066} = 11.71\%$ as attrited customers. The probability for an agreement by chance p_e can thus be calculated as $p_e = 83.49 \% \cdot 88.29 \% + 16.51 \% \cdot 11.71 \% = 75.65\%$. Cohen's kappa κ is now defined as

$$\kappa := \frac{\text{accuracy} - p_e}{1 - p_e} = \frac{90.08 \% - 75.65 \%}{1 - 75.65 \%} = 0.5926.$$

● kNN classification IX

Predictive performance (cont.)

- Error costs of positives (in our example attrited customers) and negatives (existing customers) are usually different (think, e.g., of credit default or Covid-19 tests). In this context, sensitivity and specificity are more informative than accuracy.
- Sensitivity (recall or true positive rate) is defined as the number of correct positive predictions divided by the total number of positives (55.43 % in our example).
- Specificity (true negative rate) is defined as the number of correct negative predictions divided by the total number of negatives (96.93 %).
- Our kNN classifier is good in predicting customers that do not terminate their credit card contract but bad in predicting customers that do.

Possible reasons, possible problems? *due to unbalanced data*

● SVM I

We now use SVMs for the same classification task.

Introduction

- We will fit a linear SVM as well as a SVM with a RBF kernel to the training data to predict labels for the test data.
- For the linear SVM (without kernel) we determine the most appropriate value for the hyperparameter C, also known as *Cost*, out of the set (grid) $\{2^{-5}, 2^{-4}, \dots, 2^4, 2^5\}$. Again, this can be done in parallel.
- The SVM with RBF kernel has two hyperparameters (*sigma* and *C*). We do not provide a set of possible values but let caret try 10 reasonable parameter combinations that it chooses itself (*tuneLength* parameter).

● SVM III

Model summary (linear SVM)

```
1 print(svmLinear)
2 Resampling results across tuning parameters:
3   C      Accuracy   Kappa
4   0.03125  0.9017494  0.5792921
5   0.06250  0.9010053  0.5806835
6   0.12500  0.9011295  0.5831063
7   0.25000  0.9006332  0.5821980
8   ...
9   32.00000 0.9003851  0.5823553
10
11 Accuracy was used to select the optimal model using the
12 largest value.
13 The final value used for the model was C = 0.03125.
```

In this case, choosing C with regard to accuracy leads to a different value than choosing it with regard to Cohen's kappa.

● SVM V

Model summary (SVM with RBF kernel)

```
1 print(svmRadial)
2
3 Resampling results across tuning parameters:
4   C      Accuracy   Kappa
5   0.25  0.9175042  0.6405062
6   ...
7   4.00  0.9405779  0.7624900
8   8.00  0.9420669  0.7711532
9   16.00 0.9410748  0.7703549
10
11 128.00 0.9304059  0.7378950
12
13 Tuning parameter 'sigma' was held constant at a value of
14 0.05209594. Accuracy was used to select the optimal
15 model using the largest value. The final values used for
16 the model were sigma = 0.05209594 and C = 8.
```

● SVM II

Hyperparameter tuning (linear SVM)

```
1 set.seed(2021)
2 trControl <- trainControl(method = "cv", number = 5) #
3   # only use 5-fold cross validation to speed up estimation
4 grid<-expand.grid(C=2^{(-5:5)})
5
6 cl <- makePSOCKcluster(6) #number of cores/threads
7 registerDoParallel(cl)
8
9 svmLinear <- train(x=xTrainScaled, y=yTrain, method =
10   "svmLinear", trControl=trControl, tuneGrid = grid)
11
12 stopCluster(cl)
```

● SVM IV

Hyperparameter tuning (SVM with RBF kernel)

```
1 set.seed(2021)
2 cl <- makePSOCKcluster(6) #number of cores/threads
3 registerDoParallel(cl)
4
5 svmRadial <- train(x=xTrainScaled, y=yTrain, method =
6   "svmRadial", trControl=trControl, tuneLength = 10)
7
8 stopCluster(cl)
```

● SVM VI

Predictive performance

```
1 yPredSvmLinear<-predict(svmLinear,newdata=xTestScaled)
2 yPredSvmRadial<-predict(svmRadial,newdata=xTestScaled)
3
4 # output on the next slides
5 print(confusionMatrix(data=yPredSvmLinear, reference=yTest))
6 print(confusionMatrix(data=yPredSvmRadial, reference=yTest))
```

● SVM VII

Predictive performance (linear SVM)

Confusion Matrix and Statistics

		Reference	
Prediction	Attrited Customer	Existing Customer	Existing Customer
Attrited Customer	186	46	
Existing Customer	155	1679	
Accuracy :	0.8027		
	95% CI : (0.8891, 0.9152)		
No Information Rate	0.6349		
P-Value [Acc > NIR]	< 2.2e-16		
Kappa :	0.5951		
McNemar's Test P-Value :	2.582e-14		
Sensitivity :	0.54545		
Specificity :	0.97333		
Pos Pred Value :	0.80172		
Neg Pred Value :	0.91549		
Prevalence :	0.16505		
Detection Rate :	0.09003		
Detection Prevalence :	0.11229		
Balanced Accuracy :	0.75939		
'Positive' Class : Attrited Customer			

● SVM IX

Predictive performance (summary)

- The predictive performance of the linear SVM is very similar to that of the kNN classifier in terms of accuracy, Cohen's kappa, sensitivity, and specificity.
- The SVM with RBF kernel can improve on these results with an accuracy of 93.80 % and a kappa of 0.7647.
- This substantial increase is mainly due to an improvement in predicting attriting customers (sensitivity 75.66 %). This compares to a sensitivity of 54.55 % for the linear SVM.

● SVM VIII

Better than linear SVM & KNN
in this Example.

Predictive performance (SVM with RBF kernel)

Confusion Matrix and Statistics

		Reference	
Prediction	Attrited Customer	Existing Customer	Existing Customer
Attrited Customer	258	45	
Existing Customer	83	1680	
Accuracy :	0.938		
	95% CI : (0.9268, 0.9481)		
No Information Rate	0.8349		
P-Value [Acc > NIR]	< 2.2e-16		
Kappa :	0.7647		
McNemar's Test P-Value :	0.001074		
Sensitivity :	0.7566		
Specificity :	0.9739		
Pos Pred Value :	0.8515		
Neg Pred Value :	0.9529		
Prevalence :	0.1651		
Detection Rate :	0.1249		
Detection Prevalence :	0.1467		
Balanced Accuracy :	0.8653		
'Positive' Class : Attrited Customer			

Classification And Regression Trees

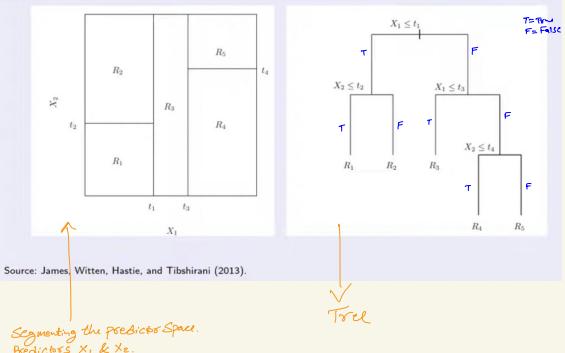
Introduction I

- Trees can be applied to both **classification** and **regression** problems.
- Trees involve stratifying or **segmenting** the predictor space into a number of simple **box-shaped regions**.
- This segmenting of the predictor space is performed based on a set of **splitting rules**. These rules can be summarized in a **tree**.

Making predictions

Introduction II

Exemplary partition of a two-dimensional feature space and the corresponding tree



- Predictions** for a specific observation are typically made by a **majority vote** (classification) or by using the **mean** or mode of the training observations (regression) in the **region** corresponding to the given observation.
- For a given region, the prediction for every observation that falls into this region is the same!

● Building a tree I

- In theory, for making predictions as on the previous slide the regions could have any shape.
- However, for simplicity and ease of interpretation the predictor space is divided into **high-dimensional rectangles** or boxes.
- These boxes are chosen such that the classification error (classification) or the squared errors of the residuals (regression) are minimized.

● Building a tree III

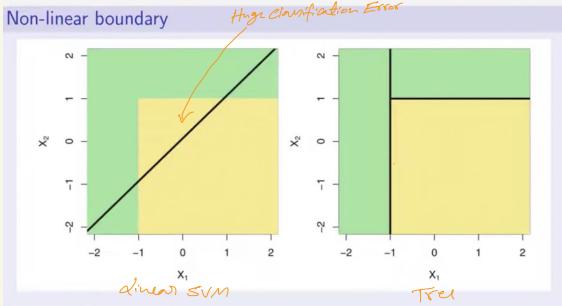
- For performing recursive binary splitting, one first selects the predictor X_j and the **cutpoint s** such that the squared errors (regression) or the classification error rate (classification) are **minimized** over the resulting regions $\{X|X_j < s\}$ and $\{X|X_j \geq s\}$. In classification tasks, sometimes also the **Gini index** or **entropy**¹ is used as a criterion for making tree splits.
- In this process, all predictors X_1, \dots, X_p and all possible values for the cutpoint s are considered.
- This process is then **repeated in each of the resulting regions**. Therefore, one again looks for the best predictor and best cutpoint to minimize the squared regression errors or the classification error rate further.

¹The Gini index is defined as $G := \sum_{k=1}^K \hat{p}_{mk}(1 - \hat{p}_{mk})$, where \hat{p}_{mk} represents the proportion of training observations in the m 'th region that are from the k 'th class. Entropy is defined as $D = -\sum_{k=1}^K \hat{p}_{mk} \log \hat{p}_{mk}$.

● Building a tree V

- A possible remedy to overfitting trees is to grow the tree only as long as the reduction in **classification or regression error exceeds some threshold**, yielding smaller trees.
- This minimum required reduction can be described by a **complexity parameter** balancing reductions in classification or regression error against the **complexity** of the model.
- In the following application, we will see that a high value of the complexity parameter leads to more shallow trees while a low value yields deeper trees.
- The optimal value of the complexity parameter can be determined via **cross-validation**.

● Trees vs. linear models II



Source: James, Witten, Hastie, and Tibshirani (2013).

● Building a tree II

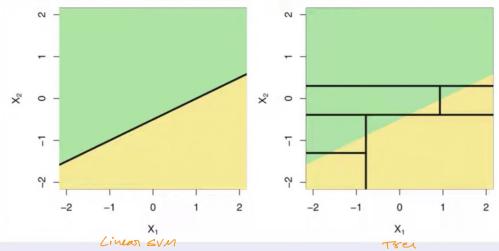
- (possible to partition the predictor space into as many boxes)
- It is computationally infeasible to consider any possible partition of the feature space into a finite number of boxes.
 - Therefore, one typically relies on a **top-down, greedy approach** that is commonly referred to as **recursive binary splitting**.
 - top-down:** The approach begins at the top of the tree where all observations belong to a single region and then successively splits the feature space.
 - greedy:** At each step the **best split** at that **particular** step is made. One does **not** look ahead and pick a split that will lead to a better tree at some **future** step.

● Building a tree IV

- This process continues until some **stopping criterion** is met.
- By adding further splits to the tree and growing the tree deeper and deeper, the squared regression errors or the classification error rate can only decrease (or stay the same).
- Consequently, the stopping criterion is crucial to avoid **overfitting** due to a tree that is too **complex**! Indeed, smaller trees with fewer splits might lead to lower variance and better interpretation at the cost of some additional bias.

● Trees vs. linear models I

Linear boundary



Source: James, Witten, Hastie, and Tibshirani (2013).

For linear boundary using a Tree leads to large bias & huge Classification Error.

● Advantages and disadvantages I

Advantages

Trees

- are **easy to explain**.
- closely mirror **human decision-making**.
- can be displayed graphically.
- can handle qualitative predictors.

● Advantages and disadvantages II

Disadvantages

- (Single) trees often exhibit an **inferior predictive accuracy** as compared to many other approaches.
- Trees can be very **non-robust**. A small change in the data might lead to a completely different tree.

By aggregating many trees via methods like bagging, random forests, or boosting, their predictive performance can be improved substantially. An example is included in the application.

● Classification Trees: Applications

● Applications

- On the following slides we will employ tree models for both classification and regression tasks.
- We will again use the credit card customer dataset from the previous lecture to predict customer churn (classification). To provide a more complete picture we will also employ boosted decision trees (a special kind of random forests).
- In the regression task, we will forecast health insurance premia. The data are available at **kaggle**.

● Classification I

We start with the classification task. Here is a short reminder of the credit card customer dataset:

Credit card customer churn prediction - reminder

- The manager responsible for the credit card customers is worried about an increasing fraction of the customers quitting their credit card services.
- To slow down customer churn she wants to proactively contact customers who are about to leave their credit card services to change their decision.
- Therefore, she needs **predictions** on **who is going to quit the services**.

● Classification II

Credit card customer churn prediction - reminder (cont.)

- The dataset contains information on more than 10,000 customers including age, salary, credit car limit, and other features.
- The data is **unbalanced** with only about 16 % of the customers having churned. This complicates training and interpretation of the accuracy of the models.
- Again, we focus on the application of the models and will not specifically address the class imbalance problem.

● Import and preprocess data I

We proceed similar to slides 220 ff.

Import data

```
1 #download bankChurners.csv from here: https://www.kaggle.com/sakshigoyal7/credit-card-customers/download
2 #then import csv-file
3
4 library(readr)
5 bankChurners <- read_csv("data/bankChurners.csv")
6
7 #drop two last columns (according to dataset description at https://www.kaggle.com/sakshigoyal7/credit-card-customers)
8 bankChurners<-bankChurners[,1:(ncol(bankChurners)-2)]
```

● Import and preprocess data II

As opposed to the KNN/SVM example, we do **not exclude categorical features**.

Create training and test set

```
1 # We randomly select 80 % of observations for the training and the remaining 20 % of observations for the test set.
2 # We exclude CLIENTNUM which is a number for identifying individual customers and has no predictive value.
3 library(dplyr)
4 set.seed(2021) #for reproducibility
5 index<-sample.int(n=2, size=nrow(bankChurners), prob=c(0.8,0.2), replace=TRUE)
6
7 xTrain<-bankChurners %>% select (-c("CLIENTNUM", "Attrition_Flag"))
8 filter(index==1)
```

● Import and preprocess data III

Create training and test set (cont.)

```

1 xTest<-bankChurners %>% select (-c("CLIENTNUM", "Attrition_
  Flag")) %>%
2   filter(index==2)
3
4 yTrain<-as.factor(bankChurners$Attrition.Flag[index==1])
5 yTest<-as.factor(bankChurners$Attrition.Flag[index==2])

```

Decision trees do **not** require feature scaling because they are not sensitive to the variance in the data.

With increasing Complexity, the accuracy from Cross Validation decreases.

● Fitting the classification tree I

Hyperparameter tuning

```

1 library(caret)
2 library(doParallel)
3 trControl <- trainControl(method = "cv", number = 10)
4 set.seed(2021) # hyperparameter selection
5
6 # if you encounter any problems with the PSOCK-cluster, the
  according lines can be commented out.
7 cl <- makePSOCKcluster(6) # number of cores/threads
8 registerDoParallel(cl)
9 treeModel<-train(x=xTrain, y=yTrain, method="rpart",
10   tuneLength=20, # training trees is
    relatively fast, so we can consider
    more possible values
11   trControl=trControl, metric="Accuracy")
12 stopCluster(cl)

```

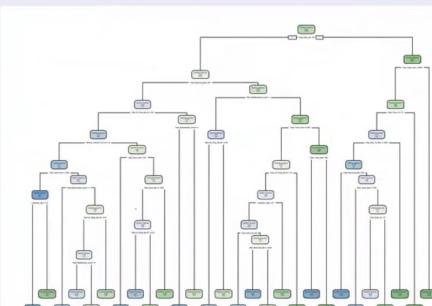
● Fitting the classification tree III

Hyperparameter tuning - the complexity parameter

- The complexity parameter **cp** in the **rpart** package (implicitly used by **caret**) is the **minimum improvement required at each node to make a further split**.
- The parameter **balances** possible **reductions in the classification error** via further splits against the **complexity** of the model (number of splits, depth).
- A **high value of cp** leads to more **shallow trees** while a **low value yields deeper trees** (that might overfit). We will exemplify this on the following slides.

● Fitting the classification tree V

Visualization of the tree (cont.)



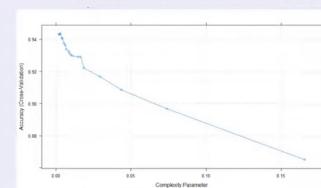
● Fitting the classification tree II

Hyperparameter tuning (cont.)

```

1 print(treeModel$bestTune)
2           cp # complexity parameter
3           0.002851218
4 plot(treeModel)

```



● Fitting the classification tree IV

The **rpart.plot** R-package provides convenient visualizations of tree models. However, it requires a tree model that was fitted by **rpart** directly (and not via **caret**). We will therefore again fit the respective model with the optimal parameters determined in the hyperparameter tuning step.

Visualization of the tree

```

1 library(rpart.plot)
2 bestParam<-treeModel$bestTune # best complexity parameter
  from cross validation
3
4 # fit tree on training set (index==1) and remove first
  column (CLIENTNUM)
5 treeModell<-rpart(Attrition_Flag ~.,
6   data=bankChurners[index==1,-1],
7   cp=bestParam)
8
9 rpart.plot(treeModell) # next slide

```

● Fitting the classification tree VI

We will now **raise the complexity parameter** to obtain a more **shallow tree**.

Shallow tree

```

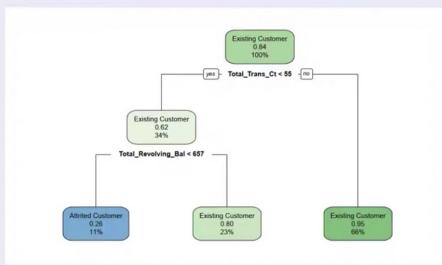
1 # multiply the complexity parameter from the current tree
2 # (bestParam) with 50
3 treeModell2<-rpart(Attrition_Flag ~.,
4   data=bankChurners[index==1,-1],
5   cp=bestParam*50)
6
7 rpart.plot(treeModell2) # see next slide

```

● Fitting the classification tree VII

Shallow tree (cont.)

Total_Trans_Ct denotes the total transaction count in the last 12 months.
Total_Revolving_Bal denotes the total revolving balance on the credit card.

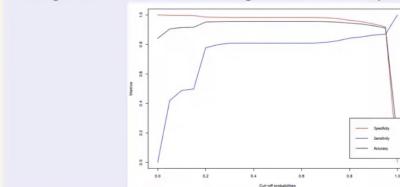


● Fitting the classification tree VIII

> Box → Class 1
 < Box → Class 2

Cut-off level

Classification in each node is performed by **majority vote**. However, it is possible and sometimes useful to choose a **cut-off level** different from 0.5. By varying the cut-off level, we can **influence sensitivity and specificity**. Accuracy closely follows specificity due to the **class imbalance** in the data. The figure is based on the training data and the deeper tree.



● Predictive performance of the classification tree I

Predictive performance

For prediction in the test set we rely on the model determined via cross validation (the deeper tree).

```

1 yPredTree<-predict(treeModel,newdata=xTest)
2
3 print(confusionMatrix(data=yPredTree, reference=yTest)) # see the next slide
  
```

● Predictive performance of the classification tree II

Predictive performance (cont.)

Confusion Matrix and Statistics

Prediction	Reference	
	Attrited Customer	Existing Customer
Attrited Customer	252	48
Existing Customer	89	1677
Accuracy : 0.9337 95% CI : (0.9221, 0.944) No Information Rate : 0.3439 P-value [Acc > NIR] : < 2.2e-16		
Kappa : 0.7472 Mcnemar's Test P-Value : 0.0006322		
Sensitivity : 0.7390 Specificity : 0.9722 Pos Pred Value : 0.8400 Neg Pred Value : 0.9496 Prevalence : 0.1551 Detection Rate : 0.1220 Detection Prevalence : 0.1452 Balanced Accuracy : 0.8556		
'Positive' Class : Attrited customer		

The predictive performance is comparable to that of the SVM with a RBF kernel, see slide 248.

● Boosted Classification Trees

Improving the predictive accuracy & predictive performance of decision & classification trees.

● Excursion I

- Deep trees are prone to **overfitting** while shallow trees might **underfit** the data.
- A possible solution to this problem is to **construct multiple (different) trees** at training time and produce forecasts based on the **mode** of the **classes (classification)** or **mean/median** of the response variable **(regression)** of the individual trees.
- This **ensemble technique** is called **random forest**.

● Excursion III

Fit boosted decision trees

```

1 library(caret)
2 library(doParallel)
3 trControl <- trainControl(method = "cv", number = 5,
   verboseIter = FALSE) # fitting boosted decision trees is
   computationally very intense. We therefore only use 5-
   fold cross validation.
4 set.seed(2021) # for reproducibility
5
6 # XGBoost can only handle numeric data. Therefore, we
   exclude all non-numeric features.
7 xTrain<-xTrain %>% select(where(is.numeric))
8 xTest<-xTest %>% select(where(is.numeric))
  
```

● Excursion II

- Random forests can be constructed by introducing **randomness** into the construction of each tree, e.g., by choosing random **subsamples**, **variables**, and **splits**, etc. This yields trees that are built **independently** of each other. This concept is often referred to as **bagging** (bootstrap aggregating).
- Another approach is to use multiple trees in a **boosting** framework where **weak learners** (shallow trees) are **combined** to yield a **stronger estimator** such that at each iteration step classification or regression error is reduced. This yields **boosted classification/regression trees** where the individual trees are **no longer built independently** of each other. For an illustration see, e.g., [here](#) or [there](#).
- On the next slides we will apply boosted decision trees to our classification problem. Therefore, we rely on the **XGBoost** algorithm.

Excursion IV

Fit boosted decision trees (cont.)

```

1 # the underlying XGBoost algorithm already uses parallel
   processing (implicitly). We therefore only employ two
   parallel processes in caret (wrapped around the implicit
   parallelism).
2 cl <- makeSOCKcluster(2)
3 registerDoParallel(cl)
4
5 # we set tuneLength=2 as caret selects 5 different
   parameters leading to 2^5 different parameter
   combinations that are considered in cross validation
6 xgBoostModel<-train(x=xTrain, y=yTrain, method="xgbTree",
   tuneLength=2, trControl=trControl, metric="Accuracy")
7
8 stopCluster(cl)

```

XGBoost

Excursion V

Fit boosted decision trees (cont.)

```

1 print(t(xgBoostModel$bestTune))
2
3 nrounds      100.0
4 max_depth    2.0
5 eta          0.4
6 gamma        0.0
7 colsample_bytree 0.8
8 min_child_weight 1.0
9 subsample    1.0

```

Excursion VI

Forecasting accuracy of boosted decision trees

```

1 yPredXgBoost<-predict(xgBoostModel,newdata=xTest)
2
3 print(confusionMatrix(data=yPredXgBoost, reference=yTest))

```

Excursion VIII

Forecasting accuracy of boosted decision trees (cont.)

- The predictive performance of the boosted decision tree model exceeds that of all previous models (decision tree, SVMs, KNN classifier).
- In particular, the sensitivity is substantially improved from 73.9 % for the single tree to 87.39 % for the boosted tree ensemble.
- By relying on the predictions from the boosted classification tree ensemble, for 298 out of 341 customers we can (rightly) predict that they are about to quit their credit card services. By, e.g., offering special price conditions to these customers, the bank manager might prevent them from quitting their contracts.
- On the other hand, we would falsely approach only 26 out of 1725 existing customers.

Excursion VII

Forecasting accuracy of boosted decision trees (cont.)

Confusion Matrix and Statistics		
		Reference
		Attrited Customer Existing Customer
Prediction	Attrited Customer	298 26
	Existing Customer	43 1699
	Accuracy :	0.9666
	95% CI :	(0.9579, 0.9739)
	No Information Rate :	0.8349
	P-Value [Acc > NIR] :	< 2e-16
	Kappa :	0.8764
Mcnemar's Test	P-Value :	0.05408
Sensitivity :	0.8739	
Specificity :	0.9849	
Pos Pred Value :	0.9198	
Neg Pred Value :	0.9753	
Prevalence :	0.1651	
Detection Rate :	0.1442	
Detection Prevalence :	0.1568	
Balanced Accuracy :	0.9294	
'Positive' class : Attrited Customer		

↙
Much better than what we have used before like KNN, SVM or decision tree.

Regression Trees

(Predicting Insurance Premium)

● Regression trees I

We now continue with regression trees.

The dataset

- For our application we rely on the insurance dataset available at [kaggle](#).
- The dataset contains information on more than 1,300 health insurance policies including four numeric features (age, bmi, number of children, expenses) and three nominal features (sex, smoker, region).
- We are interested in predicting the insurance premia (expenses).

● Regression trees III

View data

```
1 library(tibble)
2 insurance<-as_tibble(insurance)
3
4 print(insurance, n=8)
5   age sex   bmi children smoker region      expenses
6   <dbl> <chr> <dbl>    <dbl> <chr>   <dbl>
7 1 19 female 27.9     0 yes southwest 16885.
8 2 18 male   33.8     1 no  southeast 1726.
9 3 28 male   33       3 no  southeast 4449.
10 4 33 male   22.7     0 no northwest 21984.
11 5 32 male   28.9     0 no northwest 3867.
12 6 31 female 25.7     0 no  southeast 3757.
13 7 46 female 33.4     1 no  southeast 8241.
14 8 37 female 27.7     3 no northwest 7282.
15 # ... with 1,327 more rows
```

● Regression trees V

Fitting the regression tree - hyperparameter tuning

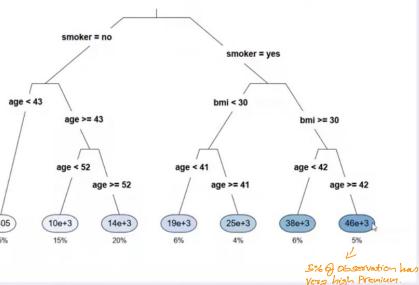
```
1 library(caret)
2 trControl <- trainControl(method = "cv", number = 10,
3   verbose=FALSE)
3 set.seed(2021) # for reproducibility
4
5 # Single trees are fitted very fast. Therefore, we do not
6 # need parallel computing and can consider various
7 # hyperparameters (tuneLength=100)
8 regressionTree<-train(x=xTrain, y=yTrain, method="rpart",
9   tuneLength=100, trControl=trControl, metric
10   ="RMSE")
11
12 print(regressionTree$bestTune) # complexity parameter
```



Can't use "Accuracy" metric which is for classification. So we will use RMSE for Regression trees.

● Regression trees VII

Visualization of the optimal tree (optimal cp)



● Regression trees II

Import data

```
1 library(readr)
2 # download data (see link on the previous slide) and save
3 # the file as insurance.csv
4 insurance <- read_csv("data/insurance.csv")
5
6 sum(duplicated(insurance))
7 1 # there seems to be a duplicate entry
8 insurance<-insurance[!duplicated(insurance),] # remove
# duplicate entry
9
10 sum(is.na(insurance))
11 0 # no missing values
12
13 summary(insurance$expenses) # Insurance Premium Summary
14   Min. 1st Qu. Median Mean 3rd Qu. Max.
15 1122    4740   9382 13270 16640 63770
```

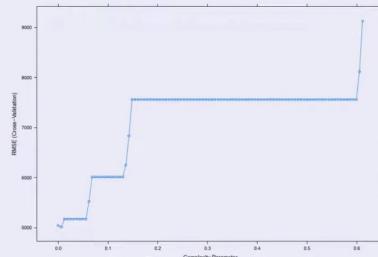
● Regression trees IV

Create training and test data

```
1 index<-sample(c(1,2), size=nrow(insurance), prob=c(0.8,0.2),
2   replace=TRUE) # sample with replacement from the set
# {1,2} with probabilities 80 % and 20 %.
3
4 xTrain<-insurance[index==1, 1:6] # exclude last column (
# expenses)
5 yTrain<-insurance$expenses[index==1]
6
7 xTest<-insurance[index==2, 1:6]
8 yTest<-insurance$expenses[index==2]
```

● Regression trees VI

Fitting the regression tree - hyperparameter tuning (cont.)



As Complexity parameter increases RMSE increases.

● Regression trees VIII

Predictive accuracy

```
1 yPredRegressionTree<-predict(regressionTree,newdata=xTest)
2 metrics<-postResample(pred=yPredRegressionTree, obs=yTest)
3
4 print(metrics)
5   RMSE Rsquared      MAE
6   4238.624 0.8779503 2909.371
7
8 # RMSE: root mean squared error
9 # Rsquared: R^2 → 87.93% In this model Close to 88% of Variance can be explained by this
# Regression Tree
10 # MAE: mean absolute error
```

Regression trees IX

Boosted regression trees

- We could also use **boosted regression trees**. However, XGBoost only works with **numeric features**. Excluding categorical variables like **smoker** leads to an **inferior model** performance ($\text{RMSE} = 11452$, $R^2 = 0.09$, $\text{MAE} = 8993$).
- Instead relying on **bagging** via the `method="treebag"` (and using all features) does not lead to an improved performance compared to the single tree model ($\text{RMSE} = 4888$, $R^2 = 0.81$, $\text{MAE} = 3045$)

These alternate Models doesn't perform well for this example.

Neural Networks: Structure

Introduction I

- The term **neural network** encompasses a **large class of models**. The main focus in this lecture is on "vanilla" neural networks which are often referred to as **single-layer perceptron** or **single hidden layer back-propagation networks**.
- We also shortly cover **deeper networks (multilayer perceptron)** as well as **convolutional neural networks**.
- This section's contents are based on Hastie, Tibshirani, and Friedman (2017) and Goodfellow, Bengio, and Courville (2016).

Introduction II

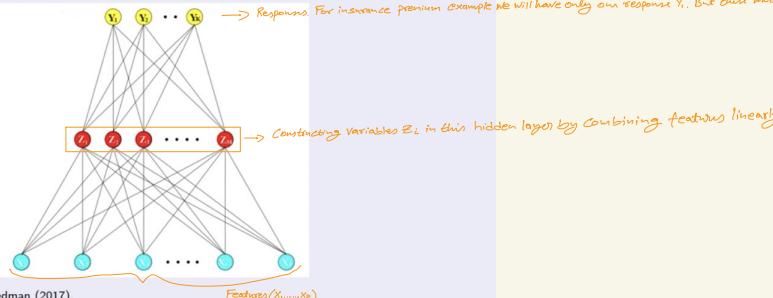
Historical trends

- Neural networks are **not new**. In fact, neural networks and deep learning date back to the 1940s and have since gone by many names.
- In the 1940s-1960s, deep learning was known as cybernetics, in the 1980s-1990s it was known as connectionism, and only since around 2006 it is commonly referred to as deep learning.
- With an **increasing amount of training data**, neural networks have become more useful.
- Over time, **models have grown in size** as both hardware and software have improved.
- Neural networks have helped to solve increasingly complicated applications with increasing accuracy.

Structure of a neural network I

The following figure shows the **network diagram** of a **single-layer perceptron**:

Single hidden layer, feed-forward neural network



Instead of using features to directly predict y using regression analysis/Tree, we would first compute the hidden variable z_i then again recombine these variables to predict y .

Structure of a neural network II

- The network structure from the previous slide is suitable for both regression and classification tasks.
- This single-layer perceptron can be seen as a **two-stage regression or classification model**, where in the intermediate layer so-called **hidden features** (Z_1, \dots, Z_M) are derived as a non-linear function of the inputs. These features are then used to model the **targets** (Y_1, \dots, Y_K).
- For **regression**, one typically sets $K = 1$ and employs only **one output node** Y_1 at the top. Nevertheless, as neural networks can handle multiple quantitative outputs at the same time, we deal with the more general case of K output nodes.
- For classification with **K classes**, the K output units model the **probability of each class**. In the training data, the $Y_k, 1 \leq k \leq K$ correspond to hot coded 0-1 variables for the k 'th class.

● Structure of a neural network III

The single-layer perceptron can be defined as follows:

Formal description

$$\begin{aligned} \text{hidden variable } & Z_m = \sigma(\alpha_{0m} + \alpha_m^T X), m = 1, \dots, M, \\ & T_k = \beta_{0k} + \beta_k^T Z, k = 1, \dots, K, \\ \text{prediction } & \hat{y}_k = f_k(X) = g_k(T), k = 1, \dots, K, \end{aligned}$$

where $Z = (Z_1, \dots, Z_M)$, and $T = (T_1, \dots, T_K)$. The so-called activation function σ is often chosen to be the sigmoid

$$\sigma_{\text{sigmoid}}(\nu) := \frac{1}{1 + e^{-\nu}} \quad \text{Non-linear fn}$$

or the rectified linear unit (ReLU) function

$$\sigma_{\text{ReLU}}(\nu) := \max(0, \nu).$$

The intercept in this case is called Bias. Regression coefficients are called weights.
For σ equal to Identity fn we get linear regression.
The linear combination of features is Non-linearly transformed to get hidden layer Z_k .
Hidden layer Z_k is again linearly combined to yield T_k & again apply non-linear fn g_k to T to get our prediction $\hat{y}_k(X)$.

● Structure of a neural network IV

Formal description (cont.)

- $f_k(X)$ is the model's prediction for the k 'th output feature Y_k ($1 \leq k \leq K$) given the input vector X .
- These predictions are constructed as a linear combination of the hidden features Z_1, \dots, Z_M to which a final transformation g_k is applied.
- For regression, often no final transformation is performed, i.e., the functions g_k are chosen to be the identity function.
- For classification, usually the softmax function

$$g_k(T) := \frac{e^{T_k}}{\sum_{l=1}^K e^{T_l}}$$

is employed. This is the same function as in the multilogit model.

Our prediction $\hat{y}_k = T_k$ for Regression.

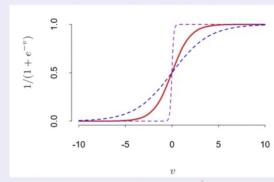
● Structure of a neural network V

- The units in the middle of the network are called **hidden units**, and the layer itself **hidden layer**, because the values Z_1, \dots, Z_M cannot be observed directly.
- Instead, these **features** are learned and modeled from the (observable) input features X_1, \dots, X_p .
- The hidden features are then used to produce predictions on the (observable) output features Y_1, \dots, Y_K .
- Note that when σ is the identity function, the entire model collapses to a linear model in the inputs. Consequently, a neural network can be thought of as a **nonlinear generalization** of the linear model.
- Deeper network architectures (multilayer perceptron) can be obtained by stacking multiple hidden layers on each other. Examples can be found in the application.
- An excellent graphical illustration of neural networks is provided in this [video](#).

For extreme case $s=10$ (Purple Curve), for small values below 0, we have no Activation, so the resulting function is zero. And for positive value, the Activation is 1.

● Structure of a neural network VI

Sigmoid function



Plot of the sigmoid function $\sigma(\nu) = \frac{1}{1+e^{-\nu}}$. s is a scaling parameter that controls the activation rate, where the "standard" sigmoid function is obtained for $s = 1$ (red curve) and the blue and the purple curve are obtained for $s = 0.5$ and $s = 10$, respectively.

Source: Hastie, Tibshirani, and Friedman (2017).

Fitting Neural Networks

• Fitting neural networks I

- The neural network model has a lot of unknown parameters that have to be fitted to the data.
- These parameters are often referred to as weights (the vectors $\alpha_m, m = 1, \dots, M$ and $\beta_k, k = 1, \dots, K$ on Slide 301) and biases (the parameters $\alpha_{0m}, m = 1, \dots, M$ and $\beta_{0k}, k = 1, \dots, K$ on the same slide). Sometimes, the weights and biases are simply referred to as weights.
- The parameters are chosen such that the model predictions fit the training data well.

in concepts

• Fitting neural networks II

- For regression, one usually relies on the sum of squared errors as a measure of fit:

$$\text{Cost fn} \longrightarrow R(\theta) = \sum_{k=1}^K \sum_{i=1}^N (y_{ik} - f_k(x_i))^2,$$

↑ response
↑ prediction

where θ is the vector of all trainable parameters and N denotes the number of training examples.

- For classification, one often uses the cross-entropy (deviance):

$$R(\theta) = - \sum_{i=1}^N \sum_{k=1}^K y_{ik} \log f_k(x_i).$$

The corresponding classifier is $G(x) := \arg \max_k f_k(x)$. That is, the class to which the highest probability is assigned, is chosen as the prediction.

• Fitting neural networks IV

- Calculating gradients based on the whole data sample as on Slide 306 and back-propagating them through the network is called batch learning. This can be costly in computational terms as datasets can be quite large.
- Therefore, one usually relies on the stochastic gradient descent (SGD) algorithm which is also referred to as mini-batch learning. The SGD algorithm only uses small random subsamples to update the network weights. As a consequence, the computational burden for each iteration step does not increase with the (total) number of training examples.
- The number of training examples in each mini-batch is referred to as batch size while a complete sweep over the entire dataset is referred to as one epoch. A neural network is typically trained over multiple epochs.

• Fitting neural networks III

- The generic approach towards minimizing $R(\theta)$ is via gradient descent.
- In the setting of neural networks, gradient descent is often referred to as back-propagation as the gradient can easily be derived by the chain rule of differentiation. This can be done in a forward and backward sweep over the network.
- For details on the back-propagation equations, we refer to Hastie, Tibshirani, and Friedman (2017, p. 396).
- An excellent graphical illustration and explanation on the training error in neural networks and the back-propagation algorithm can be found in [this](#) and [that](#) video.

Already include the Validation Set in our Training.

• Fitting neural networks V

Overfitting

- Neural networks involve numerous parameters and are therefore prone to overfitting at the global minimum of the cost function R (Slide 306).
- While there are many means to mitigate overfitting, e.g., by using smaller batch sizes that have a regularizing effect, there are some methods that explicitly address the problem of overfitting. In the application, we discuss and apply dropout and early stopping.
 - Dropout is a frequently used approach where non-output units are randomly removed from the network during training.
 - Early stopping addresses the question of how long to train a neural network by stopping training when the model performance starts to deteriorate on a validation set.

Before we see these methods in work in the application, we quickly introduce convolutional neural networks.

Convolutional Neural Networks

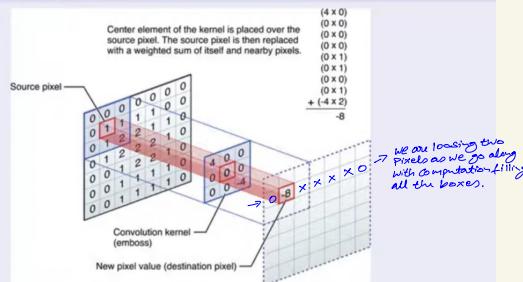
● Convolutional neural networks I

Motivation

- Convolutional neural networks (CNNs) are a **specialized** kind of neural network for processing data with a **grid-like topology** (e.g., **images**)
- CNNs have been tremendously successful in practical applications.
- The term “convolutional” indicates that CNNs use a **mathematical operation called convolution**, which is a specialized kind of **linear operation**. CNNs are thus neural networks that use convolution instead of general matrix multiplication (like in the single-layer or multilayer perceptron) in at least one of their layers.

● Convolutional neural networks II

The convolution operation



Source: Graphical representation of 2D Convolution by Jon Hoffman licensed under CC BY 4.0.

● Convolutional neural networks III

Convolution kernels are also often referred to as **filters** where each filter is meant to extract a specific feature.

Feature extraction (edge detection)



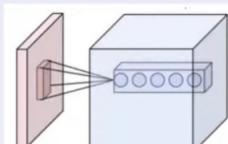
→ *Convolutional Kernel*

Source: Work by Michael Plotke licensed under CC BY-SA 3.0.

● Convolutional neural networks IV

- Typically, **many different filters** are used in each layer of a convolutional neural network in parallel. This is done to extract various different features at the same time.
- As a result, the output of a convolutional layer has the structure of a cuboid.

Output structure of a convolutional layer with multiple filters



Source: Input volume connected to a convolutional layer by Aphex34 licensed under CC BY-SA 4.0.

● Convolutional neural networks V

- Some **further architectural components and considerations** like, e.g., **pooling layers** are covered in the **application**.
- An extensive treatment of the CNN model can be found in Chapter 9 of the book by Goodfellow, Bengio, and Courville (2016).
- In the same book, you can also find information on other network architectures, e.g., **recurrent neural networks**, that are beyond the scope of this lecture.

We now turn to an **application** of both **multilayer perceptron models and CNNs**.



Neural Networks Application

● Application I

- In most companies, a variety of **documents only exists in paper form** (e.g., invoices, bank statements, printouts of static-data, business cards, etc.)
- Of course, a **scanned image can not be searched in its native state**.
- It is therefore common to **digitize printed texts** to enable **electronic editing, searching, compact storage, and online display**. The retrieved "hidden text" behind an image can also be fed directly into further machine processes, e.g., for automatically managing invoices.
- The technique behind this is called **optical character recognition (OCR)**.

● Application III

Keras

- We employ the **keras R-package** for fitting the neural networks.
- The package provides an **interface** for the open-source Python library **Keras** which in turn acts as an interface for the **TensorFlow** library. The syntax in the R-package is very similar to the Python syntax in the original library.
- Keras is one of the leading **high-level neural network APIs**.
- It has a focus on enabling **fast experimentation**. It is very **user-friendly** making it easy to quickly prototype deep learning models.
- Keras allows the training of neural networks on both **CPUs and GPUs** without changing the code.
- It supports arbitrary network architectures and is therefore appropriate for building essentially any deep learning model.

● Import and preprocess data I digitized images data

Import data

```
1 library(keras)
2
3 mnist<-dataset_mnist() # downloads the MNIST database
4 object.size(mnist)/1000000
5      219.8 bytes #size in MB
6
7 str(mnist)
8   List of 2
9     $ train:List of 2
10    ..$ x: int [1:60000, 1:28, 1:28] 0 0 0 0 0 0 0 ...
11    ..$ y: int [1:60000(1d)] 5 0 4 1 9 2 1 3 ...
12   $ test :List of 2
13     ..$ x: int [1:10000, 1:28, 1:28] 0 0 0 0 0 0 0 ...
14     ..$ y: int [1:10000(1d)] 7 2 1 0 4 1 4 9 ...
```

● Import and preprocess data III

Visualize the data (cont.)

```
1 mnist$train$y[1] #check image label
2      5 <-- response/output of 1st observation
3 image((mnist$train$x[,]),useRaster = TRUE, axes=FALSE, col
=gray.colors(10, rev=TRUE))
```



● Application II

- On the following slides we employ neural networks for the task of **handwritten digit recognition** (a **classification** problem with 10 classes).
- We start with a **multilayer perceptron (MLP)** model and discuss some **regularization techniques**.
- Finally, we also fit a **CNN** to the data.
- For the application, we rely on the **MNIST database** that is provided within the **keras R-package**. The database contains 60,000 training and 10,000 test examples of handwritten digits. More details on the dataset can be found [here](#).

● Application IV

- Before we can start to work with the data we have to **install** the (Python) **Keras and TensorFlow backend** first.
- This can mainly be done from within R. During the process, you might be asked which TensorFlow version you want to have installed (specifying "default" yields the CPU version which we will use here). Please note that installing Python Keras via the keras R-package seems not to work on the RStudio-Servers ... so please use your own computer.
- Before executing the following lines of code, please install **anaconda** with the default settings from [here](#).

Installing Keras and TensorFlow backend

```
1 library(keras)
2 ?install_keras
3 install_keras() # this will install miniconda and several
python-packages (keras, numpy, etc.)
```

● Import and preprocess data II

Visualize data

first observation in dataset A Matrix of numbers

```
1 View(mnist$train$x[,]) # 28x28 grayscale image
```

*2nd dimension for width
3rd dimension for height yields this image, this is how data is stored*

1st dimension for image

2nd dimension for width

3rd dimension for height

● Import and preprocess data IV

Create training and test set

The data are stored in a three dimensional array (image x width x height), see Slide 319. To be fed into a MLP the matrix has to be **flattened**, that is, transformed into a vector. Additionally, we transform the gray-scale values (currently between 0 and 255) to the range [0, 1].

```
1 x_train <- mnist$train$x
2 y_train <- mnist$train$y
3 x_test <- mnist$test$x
4 y_test <- mnist$test$y
5
6 # reshaping the data
7 x_train <- array_reshape(x_train, c(nrow(x_train), 784))
8 x_test <- array_reshape(x_test, c(nrow(x-test), 784))
9 # rescaling the data
10 x_train <- x_train / 255
11 x_test <- x_test / 255
```

● Import and preprocess data V

Create training and test set (cont.)

The y data (`y_train`, `y_test`) are integer vectors with values ranging from 0 to 9. For training, we **one-hot encode** the two vectors into binary class matrices. This is done via the Keras `to_categorical()` function.

```
1 str(y_train) #structure before transformation
2      int [1:60000(1d)] 5 0 4 1 9 2 1 3 1 4 ...
3
4 y_train <- to_categorical(y_train, 10)
5 str(y_train) #now a matrix
6      num [1:60000, 1:10] 0 1 0 0 0 0 0 0 0 0 ...
7
8 y_test <- to_categorical(y_test, 10)
```

5 0 4 1 9 2 ...
No Yes No No No No : Are these zero
0 1 0 0 0 0 : Encoding in binary

● Multilayer Perceptron - Application

Categorising hand written digits 0-9 using MLP (Neural Network)

● MLP I

Building a first model

- The core data structure in Keras is a model. The simplest type of model is a sequential one where (potentially different kinds of) **layers** are stacked sequentially on each other.
- We start with defining a sequential model via `keras.model_sequential()` and subsequently add layers to it.
- Instead of the object-oriented syntax in the Python Keras library (`model.add()`), the R-package uses the pipe operator (`%>%`) that we are already familiar with from the `dplyr` R-package.
- Even shallow models can exhibit hundreds of thousands of parameters. Therefore, be careful of **overfitting!**

● MLP II

Building a first model (cont.)

```
1 modell <- keras_model_sequential() #define a sequential
                                     model
2
3 # add one hidden layer with 256 neurons and the relu
   activation function
4 # The input_shape corresponds to the length of the flattened
   images (28*28=784)
5 modell %>%
6   layer_dense(units=256, activation="relu", input_shape=c
   (784)) %>%
7   layer_dense(units=10, activation="softmax") #output
   layer with one neuron for each class reflecting the
   probability for each digit
```

Hidden layer
Output layer

● MLP III

Building a first model (cont.)

```

1 # printing a model summary
2 summary(modell)
3 Model: "sequential"
4
5 Layer (type)      Output Shape     Param
6 =====
7 dense (Dense)    (None, 256)     200960
8
9 dense_1 (Dense)  (None, 10)      2570
10
11 Total params: 203,530
12 Trainable params: 203,530
13 Non-trainable params: 0
14

```

lots of parameters, hence bound to suffer from Overfitting

● MLP V

Building a first model (cont.)

Now, we can **train** the model. The results from the training steps are saved in the history object:

```

1 history <- modell %>% fit(
2   x.train, y.train,
3   epochs = 50, #epoch = one sweep over the whole dataset
4   batch_size = 64, #number of images processed after which
        the parameters are updated
5   validation_split = 0.2 #to provide an indication of the
        generalization performance of the model
6 )
```

20% Validation
20% train

● MLP VII

Making predictions from the model

```

1 predictions<-modell%>%predict_classes(x_test)
2 predictions[1] #prediction for first image in the test set
3   7
4 y.test[1,] #true value
5   0 0 0 0 0 0 0 1 0 0 ← Binary value for 7.

```



● Multilayer Perceptron : Overfitting

● MLP IX

Visualization of the training progress (cont.)

- It is obvious from the previous slide (and by comparing the accuracy (and loss) in the training and in the test set) that our first **model overfits** the data. The model ends up memorizing the training sample (with its more than 200,000 parameters) and does **not generalize well** to previously unseen data.
- While the **loss** continuously decreases on the training set, it **starts increasing from epoch 10 on the validation set** (20 % of the training sample randomly selected via the validation.split parameter).
- As a consequence, the classification accuracy on the validation set does also not further increase (as opposed to the accuracy on the training set).
- A possible solution to overfitting in neural networks are **regularization techniques such as dropout and early stopping**.

● MLP IV

Building a first model (cont.)

Next, we compile the model from the previous slides with an appropriate **loss function, optimizer, and error metric**:

```

1 modell %>% compile(
2   loss = 'categorical_crossentropy',
3   optimizer = optimizer_rmsprop(),
4   metrics = c('accuracy')
5 )

```

● MLP VI

Evaluating the model

```

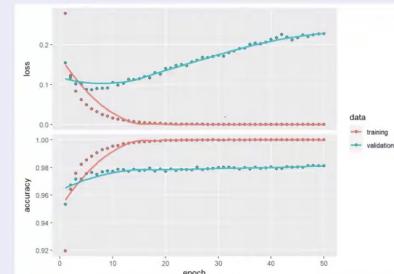
1 # Evaluate modell on the test set:
2 resl <- modell %>% evaluate(x.test, y.test)
3 print(resl) # prints categorical_crossentropy loss and
              classification accuracy (the selected loss function and
              error metric specified in the compile-stage)
4   loss accuracy
5   0.1864205 0.9816000 ← Very Accurate but this could be due
              to overfitting in training set. Needs
              to see if it generalizes well for new data.
6
7 # The model evaluates in the training set:
8 print(modell %>% evaluate(x.train, y.train))
9   loss accuracy
10  0.04544459 0.99624997

```

● MLP VIII

Visualization of the training progress

```
1 plot(history) #overfitting! In training set due to more parameters than observations.
```



Loss increases with every epoch for Validation set. So, the Model doesn't generalize well for New data.

● MLP X

Dropout

- proposed by Srivastava, Hinton, Krizhevsky, Sutskever, and Salakhutdinov (2014)
- powerful regularization method applicable to a broad family of models
- computationally inexpensive
- frequently used in the current literature

● MLP XI

Dropout (cont.)

- Dropout trains an **ensemble** of subnetworks of a given neural network.
- Therefore, non-output units are **randomly removed from the network**. This is typically achieved by multiplying the outputs of the respective neurons with zero. For each mini-batch, a different subsample of hidden units is used.
- Calculating gradients and backpropagating these through the network continues as usual.

● MLP with dropout I

We now continue by adding a **dropout layer** to model1. The **dropout rate** specifying the percentage of neurons excluded per mini-batch serves as a **hyperparameter** to the training process. Here, we choose a dropout rate of 50 %.

Specifying a regularized model

```
1 model2 <- keras_model_sequential()
2 model2 %>%
3   layer_dense(units = 256, activation = 'relu', input_shape
4     = c(784)) %>%
5   layer_dropout(rate = 0.5) %>%
6   layer_dense(units = 10, activation = 'softmax')
```

→ One extra layer added this time.

● MLP with dropout III

We continue as before ...

Compiling and fitting the regularized model

```
1 model2 %>% compile(
2   loss = 'categorical_crossentropy',
3   optimizer = optimizer_rmsprop(),
4   metrics = c('accuracy')
5 )
6
7 history <- model2 %>% fit(
8   x_train, y_train,
9   epochs = 50,
10  batch_size = 64,
11  validation_split = 0.2
12 )
```

● MLP with dropout V

- Accuracy in the validation set (slightly) increases with the number of epochs while the loss over the validation set only slightly increases → **overfitting is not a (major) issue** any more.
- While the accuracy of the regularized model (about 99 %) is lower in the training set (compared to the original model), it is higher (98.2 %) in the test set → **generalization** performance has **improved**.
- This is also reflected in a lower loss over the test set (0.111 for the regularized vs. 0.186 for the original model).

● MLP XII

Early stopping

- addresses the question of **how long to train** a neural network
- Two little training might lead to underfitting while too much training might cause overfitting (bad generalization to the test set).
- Early stopping proposes a compromise to this problem by **stopping training** at the point **when performance on a validation set starts to degrade**.
- simple, effective, and widely used

● MLP with dropout II

Model summary of the regularized model

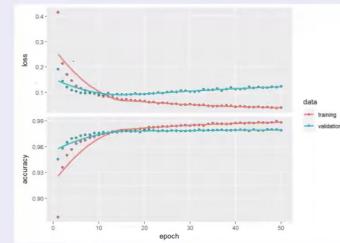
Layer (type)	Output Shape	Param
dense_2 (Dense)	(None, 256)	200960
dropout (Dropout)	(None, 256)	0
dense_3 (Dense)	(None, 10)	2570
Total params:	203,530	
Trainable params:	203,530	
Non-trainable params:	0	

Note that the model relies on the same number of parameters as before.

● MLP with dropout IV

Visualizing the training progress of the regularized model

```
1 plot(history)
```



loss in Validation set increases but not as bad as last time (model1)

Hence, it's a good way to regularize a Neural Network & to prevent from Overfitting.

Multilayer Perceptron : Deep Model

● Deep model I

- Now that overfitting is no longer a (major) issue, we can fit a **deeper model** to the data.
- Therefore, we add **two additional hidden layers** to the previous model where we reduce the number of neurons by a factor of two in each consecutive layer.
- Again, we apply dropout (with a dropout rate of 50 %) to each hidden layer.

● Deep model III

Model summary

```
1 summary(model3)
2
3      Layer (type)      Output Shape       Param
4      ========
5      dense_4 (Dense)    (None, 256)     200960
6      dropout_1 (Dropout) (None, 256)     0
7
8      dense_5 (Dense)    (None, 128)      32896
9      dropout_2 (Dropout) (None, 128)     0
10
11     dense_6 (Dense)    (None, 64)       8256
12     dropout_3 (Dropout) (None, 64)       0
13
14     dense_7 (Dense)    (None, 10)       650
15
16 Total params: 242,762
17 Trainable params: 242,762
18 Non-trainable params: 0
```

● Deep model II

Specifying a deep model

```
1 model3 <- keras_model_sequential()
2 model3 %>%
3   layer_dense(units = 256, activation = 'relu', input_shape
4               = c(784)) %>%
5   layer_dropout(rate = 0.5) %>%
6   layer_dense(units = 128, activation = 'relu') %>%
7   layer_dropout(rate = 0.5) %>%
8   layer_dense(units = 64, activation = 'relu') %>%
9   layer_dropout(rate = 0.5) %>%
10  layer_dense(units = 10, activation = 'softmax')
```

● Deep model IV

We continue as before ...

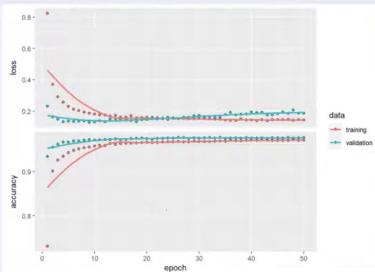
Compiling and fitting the model

```
1 model3 %>% compile(
2   loss = 'categorical_crossentropy',
3   optimizer = optimizer_rmsprop(),
4   metrics = c('accuracy')
5 )
6
7 history <- model3 %>% fit(
8   x_train, y_train,
9   epochs = 50,
10  batch_size = 64,
11  validation_split = 0.2
12 )
```

● Deep model V

Visualizing the training progress

```
1 plot(history)
```



● Deep model VI

Evaluating the model

```
1 # Evaluate the classification accuracy on the test set
2 res3 <- model3 %>% evaluate(x_test, y_test)
3 print(res3)
4
5           loss  accuracy
6 0.1864311 0.9768000
```

Multilayer Perceptron : Early Stopping

● Early stopping I

- Finally, we fit a model with **early stopping**, see Slide 335.
- Therefore, we rely on the previous model with three hidden layers.
- For illustrative purposes we reduce regularization by setting the dropout rate to 30 %. Apart from this, the model specification is the same as on Slide 343. ([modul3](#))

● Early stopping II

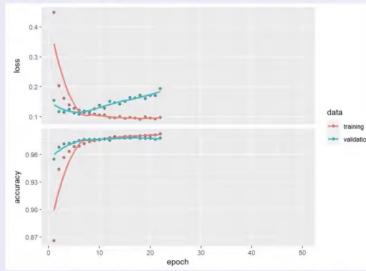
Compiling and fitting the model

```
1 # Compiling is the same as for the previous model:  
2 model4 %>% compile(  
3   loss = "categorical_crossentropy",  
4   optimizer = optimizer_rmsprop(),  
5   metrics = c('accuracy'))  
6  
7 # Training with early stopping is performed via introducing  
# a so-called callback monitoring the accuracy on the  
# validation set:  
8 history <- model4 %>% fit(  
9   x_train, y_train, epochs = 50, batch_size = 64,  
10  validation_split = 0.2,  
11  callbacks = list(callback_early_stopping(  
12    monitor = "val_accuracy",  
13    patience = 5,  
14    restore_best_weights = TRUE)))
```

● Early stopping IV

Visualizing the training progress

```
1 plot(history)
```



Evaluating the model (cont.)

- Both **deep models** do not achieve a better performance than the shallow model with dropout. **Deeper models are not better per se.**
- Neural networks need to be **carefully specified and trained**. Deeper models might perform better with more training examples, with a different dropout rate, different numbers of neurons, layers, etc.
- On the following slides, we consider a different network architecture, namely **CNNs**. This type of network is especially suited for data with a grid-like structure such as images.

● Early stopping III

The parameters to the `callback.early_stopping()` function are explained [here](#).

- monitor: quantity to be monitored
- patience: number of epochs with no improvement after which training is stopped
- restore_best_weights: whether to restore model weights from the epoch with the best value of the monitored quantity
-

● Early stopping V

Evaluating the model

```
1 res4 <- model4 %>% evaluate(x_test, y_test)  
2  
3 print(res4)  
4 loss accuracy  
5 0.1458071 0.9781000
```

The classification accuracy has (slightly) improved compared to the previous (deep) model. However, to assess whether the improvement is due to early stopping, the lower dropout-rate, or just happened by chance one would have to perform further analyses.

Convolutional Neural Networks Application

CNN I

- For the MLP models we flattened the input data, that is, we transformed the 2D images into vectors of gray-scale values.
- The CNN, however, exploits the grid-like structure of the images. Therefore, we need the training and test set in a different structure.

Preprocessing the data

```
1 x.train <- mnist$train$x/255 # grayscale values in [0,1]
2 x.test <- mnist$test$x/255
3 dim(x.train)
4 60000 28 28 → 60000 images, 28 pixels in height & width
5 dim(x.train)<-c(dim(x.train),1) #the CNN takes images in
       three dimensions, where the last dimension typically has
       three values (the RGB channels). As we have gray-scale
       images, we only need one channel. Adding one more dimension.
6 dim(x.train)
7 60000 28 28 1 → 60000 images, 28 pixels, 28 pixels, 1 RGB Channel
8 dim(x.test)<-c(dim(x.test),1) #same transformation for test
   set
```

CNN III

- The output feature map (obtained from the various filters) is sensitive to the location of the features in the input.
- Pooling layers are an approach to down sample feature maps. They summarize the presence of features in patches of the feature map.
- Common pooling methods are average and maximum pooling. Average pooling summarizes the average presence of a feature while maximum pooling captures the most activated presence of a feature ("cat" present in the respective part of the image or not").
- In our example, we use maximum pooling over patches of size 2x2. This reduces each 28x28 feature map obtained by the convolutional layer to a 14x14 map. *↳ because in each dimension we are losing 2 pixels*.

Adding a pooling layer

```
1 modelCNN <- modelCNN %>%
2   layer_max_pooling_2d(pool_size=c(2,2))
```

CNN V

- To complete the model, we feed the outputs from the last convolutional layer (with attached pooling layer) into a dense layer to perform classification.
- The dense layer has the same structure as a layer in the MLP.
- Before the outputs from the convolutional layer can be fed into the dense layer, the 3D output has to be flattened to 1D.

Adding a dense layer

```
1 modelCNN <- modelCNN %>%
2   layer_flatten() %>%
3   layer_dense(units=10, activation="softmax")
```

CNN VII

Dimensionality of the layers

- The input data consist of 28x28 gray-scale images.
- Applying filters of size 3x3 to them ("convolution") leads to the loss of two pixels in each dimension (if we do not apply some kind of padding). This is because a filter typically starts at the upper left corner of the image with the left-hand side of the filter sitting on the far left pixels of the image. The filter then moves across the image until the far right pixels of the image are reached. This yields 26x26 feature maps.
- As in the first convolutional layer we apply 32 different filters, the output dimensionality is 26x26x32.
- By maximum pooling with a patch size of 2x2 the dimensionality is reduced to 13x13x32.

CNN II

- CNNs can automatically learn a large number of filters (32 in our case in the first convolutional layer) in parallel.
- Each filter provides highly specific features that can be detected anywhere on the input images.
- In this example, we apply filters of size 3x3.
- The following lines of code specify the whole model as sequential and add the first convolutional layer.

Adding the first convolutional layer

```
1 modelCNN <- keras_model_sequential() %>%
2   layer_conv_2d(
3     filters=32,
4     kernel_size=c(3,3),
5     activation="relu", input_shape=c(28,28,1))
```

CNN IV

We now add another convolutional layer, this time with 64 filters for detecting more detailed features in the image, followed by another maximum pooling layer.

Adding another convolution and pooling layer

```
1 modelCNN <- modelCNN %>%
2   layer_conv_2d(filters=64, kernel_size=c(3,3), activation
3     = "relu") %>%
   layer_max_pooling_2d(pool_size=c(2,2))
```

CNN VI

The whole model architecture looks like the following:

Model summary

1 summary(modelCNN)	Output Shape	Param
2 Layer (type)		
3 =====	=====	=====
4 conv2d (Conv2D)	(None, 26, 26, 32)	320
5 max_pooling2d (MaxPooling2D)	(None, 13, 13, 32)	0
6 -----	-----	-----
7 conv2d.1 (Conv2D)	(None, 11, 11, 64)	18496
8 max_pooling2d.1 (MaxPooling2D)	(None, 5, 5, 64)	0
9 -----	-----	-----
10 flatten (Flatten)	(None, 1600)	0
11 dense_2 (Dense)	(None, 10)	16010
12 =====	=====	=====
13 Total params: 34,826		
14 Trainable params: 34,826		
15 Non-trainable params: 0		

CNN VIII

Dimensionality of the layers (cont.)

- The next convolutional layer applies 64 different filters of size 3x3 to the output (tensor) of the previous layer.
- Again, the input feature map (11x11x64) loses two "pixels" in the first two dimensions, yielding an output dimensionality of 9x9x64.
- By maximum pooling with patches of size 2x2, the number of "pixels" in the first two dimensions is essentially divided by 2 yielding the output dimension of 5x5x64.
- Flattening this output tensor yields a vector of length $5 \cdot 5 \cdot 64 = 1600$.

CNN IX

Number of parameters in each layer

- A very important feature of CNNs is **parameter sharing**. This is achieved by moving each filter over the picture thereby **employing the same parameters at each location**. As a consequence, the first convolutional layer of our network only employs 320 different parameters. This compares to 200,960 parameters in the first hidden layer of the previous MLP models!
- Each 3×3 filter involves $3 \cdot 3$ parameters + 1 bias parameter. As we employ 32 filters, this yields $(3 \cdot 3 + 1) \cdot 32 = 320$ parameters.
- Pooling does not involve any trainable parameters. (The size of the patches over which pooling is performed is a hyperparameter.)

CNN X

Number of parameters in each layer (cont.)

- Each of the 64 filters of size 3×3 from the second convolutional layer is applied to all 32 feature maps from the previous layer at once.
- This yields $(3 \cdot 3 \cdot 32 + 1) \cdot 64 = 18,496$ parameters.
- Pooling and flattening do not involve (trainable) parameters.
- The last dense layer requires 10 weights for each of the 1600 neurons (for mapping to the 10 output neurons). Additionally, each output neuron requires a bias parameter. This yields $(1600 + 1) \cdot 10 = 16,010$ parameters.
- In total, we obtain **34,826 parameters** which are still substantially less parameters than are employed in the one hidden layer MLP model (**203,530 parameters**).

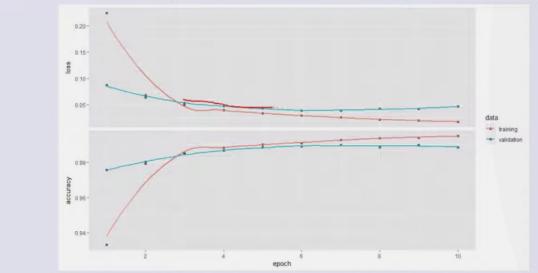
CNN XI

- Compiling and training the model is performed **analogously to the MLP models**.
- As fitting the convolutional network is **computationally very expensive**, we only perform 10 epochs.
- Our **University Computing Center** provides a **GPU cluster** which is suitable for fitting large (or many) neural networks. The GPU cluster features 6 nodes with 4 Tesla V100 GPUs and 512GB RAM as well as 23 nodes with 8 RXT2080Ti and 512GB RAM.

CNN XIII

Visualization of the training progress

```
1 plot(history)
```



CNN XIV

Evaluating the classification performance

```
1 resCNN <- modelCNN %>% evaluate(x-test, y-test)
2 print(resCNN)
3           loss   accuracy
4      0.0330491  0.9899000
```

- The CNN provides a **higher forecasting accuracy on the test set** than any of the MLP models.
- Classification accuracy could probably be improved by further tuning the model specification and by increasing the number of epochs. (↳ **so epochs**)
- This is, however, beyond the scope of this lecture.

Regulation : Regtech

● Regtech I

To comply with law & regulations

Definition

The use of AI & ML in the management of **regulatory processes** for the financial industry. Main functions are **regulatory monitoring, reporting, compliance**.

Examples:

- Anti-fraud and risk management for digital transactions (**IdentityMind Global**)
- Management of consent for customer service data (**Trunomi**)

● Regtech in Numbers² III

- Dominated by start-ups; 70% of firms are younger than 5 years
- 44,000 people employed globally
- About \$4.9.bn annual revenue
- About \$9.7.bn raised in external funding

Market Environment

Market and regulatory environment are rated to be **generally favourable** to Regtech companies by vendors in the sector. The **pace of regulatory changes has increased** after the Financial Crisis 08/09, expensive punishment for non-compliance with regulatory rulings lead to a surge in the demand for automated and reliable methods.

²The Global RegTech Industry Report by Cambridge Center for Alternative Finance and supported by EY Japan

● Regtech - Clients³ IV

The Cambridge Center for Alternative Finance in collaboration with EY has surveyed 658 Regtech firms (about 80% of firms in the sector) with the following results concerning clientèle of Regtech companies:

- 89-94% of Regtech companies target banks as clients
- 61% target insurers
- and 57% serve FinTechs

However,

- 58% also have clients outside of the financial service sector such as advisory, consultancy and regulators

This means, there is **substantial overlap between what we consider Regtech and Suptech companies**. Where many firms can be subsumed under both definitions depending on the client they are currently serving.

³The Global RegTech Industry Report by Cambridge Center for Alternative Finance and supported by EY Japan

● Regtech - Regulation VII

Regulators

Regulators have to adapt to ever new technology-enabled financial services, which may present a challenge especially for emerging and developing economies.

- Regulatory Sandboxes** - formal programs that allow certain financial services and business models which are not yet fully compliant with existing laws. The aim is to learn about the opportunities and risks that a particular innovation carries and to develop evidence-based policy.
- Innovation Hubs / Offices** - places where innovators and regulators meet to discuss solutions to sector challenges.
 - Example:** TechSprints from the UK Financial Conduct Authority (FCA)
 - So far 7 TechSprints since 2016 with an 8th planned in October 2021
 - 2-day events including industry and innovators
 - past focus areas were e.g. (1) regulatory reporting; (2) financial service and mental health; (3) anti-money laundering

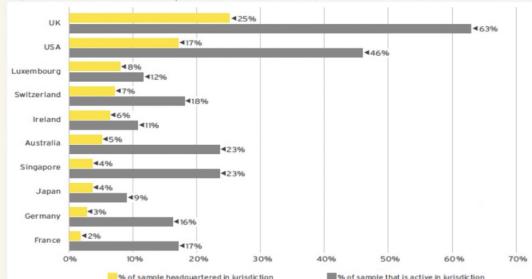
Regulator Barfins (Gesamt regulator) came up with this.
eg. crypto.

● Regtech - Market Segments II

- Profiling and Due Diligence:** Collect or integrate data from multiple internal and external sources. Aims to profile an entity, confirm their identity, or categorize them e.g. according to regulatory requirements.
- Reporting and Dashboards:** Collect or integrate data from multiple internal sources. Aims to build standardized reports for management or compliance purposes.
- Risk Analytics:** Collect or integrate data from multiple internal and external sources. Aims to assess the risk of fraud, market abuse, or other misconduct at the transaction level.
- Dynamic Compliance:** Using Machine Learning methods to facilitate and monitor regulatory changes, ensuring flexible adaptation of policies and controls.
- Market Monitoring:** Collect or integrate data from multiple external sources. Aims to match market-level adverse outcomes to regulatory or business rules. For example, poor product performance, market manipulation, etc.

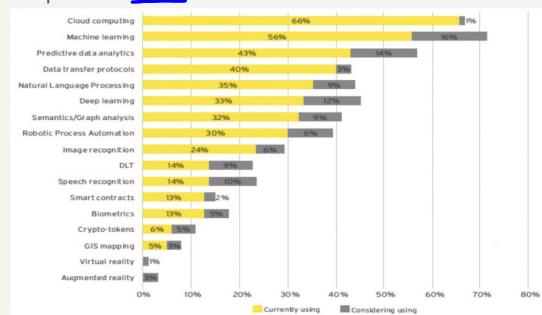
● Regtech - Markets IV

Top 10 RegTech markets, by % of firms present (headquarters or significant market share) from the CCAF report, 2019



Source: CCAF Report.

● Regtechs - Technologies and Tools used by Regtech companies VI



Source: CCAF Report.

Suptech And Systemic Risk (AI using by Regulators)

● Suptech I

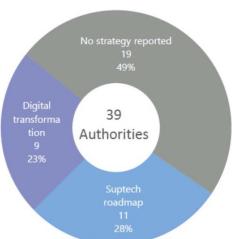
Definition

A sub-discipline of Regtech focussing on innovative technologies such as Big Data and AI supporting financial authorities in their work. The focus is on misconduct analysis, reporting, data management.

Examples:

- Collection and management of detailed data on loans in the euro area (**ECB**)
- Accelerating Suptech solutions and prototyping (**R²A**)

● Suptech - Focus areas III



Source: FSI Report.

● Systemic Risks due to AI/ML in finance IV

Drivers of systemic risk, when using AI in regulatory functions:

- AI is unable to reason about events it has not yet been exposed to. Humans can draw on a broad range of prior experience and imagination, while AI can only extrapolate from what it has seen before.
- We do not know how AI makes decisions, it is too complex for us to follow. (less transparency causes uncertainty to investors & regulators)
- AI is more likely to amplify current cycles (pro-cyclicality) than human regulators. Automation favours homogenous methodologies and standardization
- The high predictability and transparency of AI will enable individuals to bet against it.

● Endogenous risk VII

Endogenous Risk

Endogenous risk is caused by events from within the system. It starts when individual entities within the system stop acting independently but synchronize their behavior. It is very difficult to measure.

AI are trained with various games for situations including interactions between entities...

- Full information games:** In Chess, all possible moves can be known. Learning Algorithms and Deep Neural Networks succeed in these games, they assume their opponent is their clone (i.e., knows the same). They do not prepare for endogenous risk situations.
- Incomplete information games:** In Poker, the opponents cards are unknown. AI performs worse than human players, as it cannot theorize about the opponent's intentions.
- Cooperative Games:** Success is even more difficult for AI, when cooperation leads to multiple local optima (e.g. Diplomacy).

● Suptech - Survey II

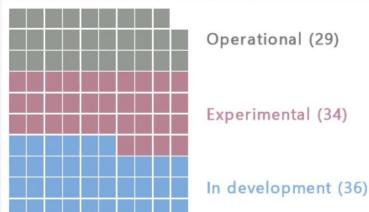
Survey

The financial stability institute (FSI) surveyed 39 financial authorities from 31 jurisdictions worldwide about the implementation of suptech strategies.

- two types of approaches were identified (which can also be combined)
 - Specific Suptech Roadmaps** based on particular needs of a department. The approach tends to be more experimental.
 - Institution-wide Digital Transformation and Data-driven Innovation programs (DT&DI)** is a broader approach encompassing large parts of the authority

● Suptech - Conclusion IV

- Suptech is still in its infancy, but is gaining momentum.
- In 2019, most suptech solutions were still experimental or in a developmental stage.
- Authorities are willing to explore the area, yet funding and infrastructure may be problematic.



Source: FSI Report.

● Exogenous risk VI

Exogenous Risk

Exogenous risk is caused by events from the outside. Much like an asteroid falling on London. It is easy to measure.

- AI is suited for the evaluation and management of exogenous risk
- It uses large data-sets, well-established statistical methods, and repeated events to train.
- Therefore, it is well suited for micro-regulation and internal risk management.

Tighten Regulation in Crisis
& loosen Regulation When Market is booming leading to another crisis.

● Ethical Considerations ⁵

Ethical risks at different levels...

- data
 - ▶ training data must be free of biases and include all relevant types of stimuli (e.g., racial or sexist biases)
 - ▶ AI should only be deployed in environments it has been trained in (i.e., no cross-use for a different purpose)
 - ▶ privacy laws must be taken into account when collecting and processing data
- algorithm
 - ▶ unethical coding influenced by biases of the developers (most developers are white males)
 - ▶ AI must be controlled for emerging biases during the life cycle
- business use
 - ▶ the purpose of using AI should be ethical in the first place
 - ▶ unintended impact of AI can be unethical

⁵Study by Thomaz et al. 2021

● Ethical Considerations

Practical Recommendations

- implementation of a general statement on the firm's intention for AI ethics ⁶
- should be an extension to the firms already existing mission or purpose statement
- implementation of internal application-specific design plans
- regular auditing of processes (flagging risks and concerns)
- records of decisions concerning ethical trade-offs for transparency

⁶Study by Thomaz et al. 2021

● Regulation of AI/ML in finance

Definition of AI according to BaFin

Artificial Intelligence is a combination of large amounts of data (Big Data), computing resources, and machine learning (ML).⁷

In machine learning, computers are given the ability to learn from data and experience on the basis of special algorithms. Compared to rule-based methods, learning takes place without the programmer specifying which results are to be derived from certain data constellations and how.

⁷BaFin (2018): Big Data trifft auf Künstliche Intelligenz: Herausforderungen und Implikationen für Aufsicht und Regulierung von Finanzdienstleistungen, URL: <https://www.bafin.de/dok/10985478>

● Regulation of AI/ML in finance

Overriding Principles

Clear responsibility of the management

- Management is responsible for company-wide strategies and guidelines or policies for the use of algorithm-based decision-making processes.
- Potentials of such processes as well as their limits and risks should be taken into account and clearly stated.
- Company-wide strategy for the use of algorithm-based decision processes should also be reflected in the IT strategy.

● Ethical Considerations

Racism and Sexism in AI

Machines in themselves have no bias, but the learning data sets and algorithms are most likely as biased as our current society is. Most of the AI developers are white males, lacking the perspectives of minorities. Consequently, **existing biases will be amplified, if we are not actively working against them.**

Organisations working against biased algorithms

- American Civil Liberties Union (**ACLU**) - fights in a wide variety of issues, but has exposed *Amazon's Rekognition* as racially biased
- Algorithmic Justice League (**AJL**) - founded in 2016 with the mission to raise public awareness to biases in AI
- several independent researchers

● Regulation of AI/ML in finance

- Regulation of artificial intelligence is still in its infancy
- Artificial intelligence and machine learning not yet clearly defined
- **European Commission** has stated in its **Digital Finance Strategy** that it intends to clarify, together with the European Supervisory Authorities, by 2024 at the latest whether and how existing financial market regulation should be applied to the use of Big Data and Artificial Intelligence (BDAI)
- **BaFin** has proposed several principles for the general use of algorithms in decision-making processes of financial firms⁸
 - ▶ Principles represent preliminary considerations for minimum supervisory requirements for the use of artificial intelligence

⁸BaFin: Big Data und künstliche Intelligenz: Prinzipien für den Einsatz von Algorithmen in Entscheidungsprozessen

● Regulation of AI/ML in finance

Algorithms

Algorithms are rules of action that are usually integrated into a computer program and solve an (optimization) problem or a class of problems.

In addition to the distinction according to the type of algorithms (how is the problem technically solved), applications of ML can also be differentiated according to result types (the basic distinction between classification, regression, and clustering) and data types (e.g., for text, language, and image data).

● Regulation of AI/ML in finance

Overriding Principles

Adequate risk and outsourcing management

- Establishment of a risk management system adapted to the use of algorithm-based decision-making processes.
- If applications are sourced from a service provider, management must also have an effective outsourcing management.
- Responsibility, reporting, and control structures must be clearly defined.
- When establishing adequate risk management, one needs to consider the risks of an algorithmic decision-making process → risk mitigating measures and processes should start exactly where a risk originates according to the polluter pays principle.

● Regulation of AI/ML in finance IX

Overriding Principles

Bias prevention

- Avoidance of a bias, i.e., the systematic distortion of results, in algorithm-based decision-making processes
- Bias must be avoided in order to:
 - be able to make business decisions based on results that are not systematically distorted,
 - exclude the possibility of systematic bias-based discrimination against individual customer groups and the resulting reputational risks.
- In accordance with the polluter pays principle, the risk of bias must be identified where it can arise, it must be analyzed and either eliminated or at least mitigated

● Regulation of AI/ML in finance X

Overriding Principles

Exclude differentiation prohibited by law

- For some financial services, it is also stipulated by law that certain characteristics may not be used for differentiation, i.e., for calculating risk and prices.
- Danger of discrimination exists if these characteristics are replaced by an approximation.
- Would be associated with increased reputational risks and legal risks.
- Companies should establish (statistical) verification processes that exclude discrimination.

● Ethics of AI in finance XI

EIOPA's Group on Digital Ethics (GDE) AI governance principles

- ① Human oversight
- ② Robustness and performance
- ③ Data governance and record keeping
- ④ Transparency and explainability
- ⑤ Fairness and non-discrimination
- ⑥ Principle of proportionality

● EIOPA's GDE AI governance principles XII

Robustness and performance

Insurance firms should assess and monitor the performance of AI systems on an on-going basis and taking due consideration of their limitations and potential shortcomings.

- Performance metrics should be adapted to the objective pursued and the nature of the data used.
- Sound data management is key to ensure the performance of AI systems.
- AI systems should produce stable outcomes over time.
- Insurance firms should develop resilient IT systems and infrastructures.

● EIOPA's GDE AI governance principles XIII

Human oversight

Insurance firms should establish adequate levels of human oversight taking into account the impact of specific AI use cases and other governance and control measures in place.

- Selection of the level of human oversight should be proportionate to the nature, scale, and complexity of the risks inherent to the specific AI use case.
- Different roles and responsibilities for the staff involved in AI processes should be clearly defined in policy documents.

● EIOPA's GDE AI governance principles XIV

Data governance and record keeping

Insurance firms should adapt the data governance and record keeping measures to the impact of specific AI use cases.

- Data used in AI models should be accurate, complete, and appropriate.
- Sound data governance should be applied throughout the AI model lifecycle.
- Data used in AI models should be handled and stored in a secure manner.
- Appropriate records of the data and the modelling methods should be kept to ensure their reproducibility/traceability.

● EIOPA's GDE AI governance principles XV

Transparency and explainability

Insurance firms should adapt the types of explanations to specific AI use cases and to the recipient stakeholders.

- Firms should adapt their explanations to the different types of stakeholders.
- Firms should strive to use explainable AI models, in particular in high-impact AI use cases.
- Data used need to be transparently communicated.

● EIOPA's GDE AI governance principles XVI

Fairness and non-discrimination

- Sound and transparent governance processes are key to ensure fairness and non-discrimination.
- Insurance firms should conduct their business in a fair manner when using AI and make reasonable efforts to take into account the outcomes of AI systems.
- Consumers not willing to share very personal and sensitive data are not strictly necessary for risks assessments should still have access to affordable insurance coverage.
- Insurance firms should respect the principle of human autonomy by developing AI systems that support consumers in their decision-making process and avoid unfair nudging practices.

● EIOPA's GDE AI governance principles XVII

Principle of proportionality

Insurance undertakings should establish the necessary governance measures that are proportionate to the nature, scale, and complexity of their operations.

- AI use case impact assessment and the governance measures should be proportionate to the potential impact of a specific AI use case on consumers and/or insurance firms.
- Insurance firms should then assess the combination of measures put in place in order to ensure an ethical and trustworthy use of AI.

● Literature XVIII

AGGARWAL, C. C., AND P. S. YU (2008): *Privacy-Preserving Data Mining: Models and Algorithms*, vol. 34 of *Advances in Database Systems*. Springer, New York, NY and Heidelberg.

DUA, D., AND C. GRAFF (2017): "UCI Machine Learning Repository."

GOODFELLOW, I., Y. Bengio, AND A. COURVILLE (2016): *Deep Learning*. MIT Press, Cambridge, MA, USA.

HASTIE, T., R. TIBSHIRANI, AND J. H. FRIEDMAN (2017): *The elements of statistical learning: Data mining, inference, and prediction*, Springer series in statistics. Springer, New York, NY, second edition, corrected at 12th printing 2017 edn.

HULL, J. C. (2020): *Machine Learning in Business*. Second Printing, Toronto, second edn.

JAMES, G., D. WITTEN, T. HASTIE, AND R. TIBSHIRANI (2013): *An Introduction to Statistical Learning*. Springer, New York, first edn.

KUHN, M. (2020): *caret: Classification and Regression Training*.R package version 6.0-86.

LESMEISTER, C. (2015): *Mastering machine learning with R: Master machine learning techniques with R to deliver insights for complex projects*, Community experience distilled. Packt Publishing, Birmingham and Mumbai.

SRIVASTAVA, N., G. HINTON, A. KRIZHEVSKY, I. SUTSKEVER, AND R. SALAKHUTDINOV (2014): "Dropout: A Simple Way to Prevent Neural Networks from Overfitting," *Journal of Machine Learning Research*, 15(56), 1929–1958.

VENABLES, W. N., AND B. D. RIPLEY (2002): *Modern Applied Statistics with S*. Springer, New York, fourth edn., ISBN 0-387-95457-0.