Dr. Mike Ashcroft, ESDA, Sweden

  - Root Cause Analysis

  - FinTech

Using statistics with modern IT

Features = structured file (tidy) taken
        from unstructured data, like
        pictures, Web pages, Audio, Video, Text

Select Features ⟹ Transform Features
        ⇓                      ⇓
Extract Features ⟶      Select Features

Supervised : target variable known a d i the data

Unsupervised : target variable not known and i the data

X = Input Features, Y = Target Variable

## Regression

Target variable is a real number.

Residuals = errors of a given model
= how far off are the real values

$$MSE = Mean\ Square\ Errors$$

$$= \frac{1}{n} \sum_{i=1}^{n} (\bar{y}_1 - g_1)^2$$

n = cases = rows in dataset

↳ average of the errors

## Ordinary Least Square (OLS) Regression

lm = Regression line that minimizes
the sum of the square residuals

Console:

? lm
?? "search phrase"

No guaranty is that the mathematics is right.
Python has a significant advantage in deep learning.

In Classischer Statistik würde man manuell
versuchen, daß die Werte zu einer linearen Regression
passe.

In Data Science, der Rechner macht das.

Eine Möglichkeit ist

## POISSON Regression
## (Generalized linear Models)

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + b_3 x_3 + \dots$$

$$x_0 = 1$$

$$\hat{y} = b_0 x_0 + b_2 x_0 + b_3 x_0 + \dots$$

$$\hat{y} = \sum_{x-1}^{n} b_n x_n$$

$$\log(\hat{y}) = \sum_{x=1}^{n} b_n x_n$$

$$\Downarrow$$

$$f(\hat{y}) = \sum_{x=1}^{n} b_n \cdot x_n$$

just some sort of transformation function

The transformation is called „link function".

↓↓

an aid to predict the values of y

to get the real values parameter

type = „response"

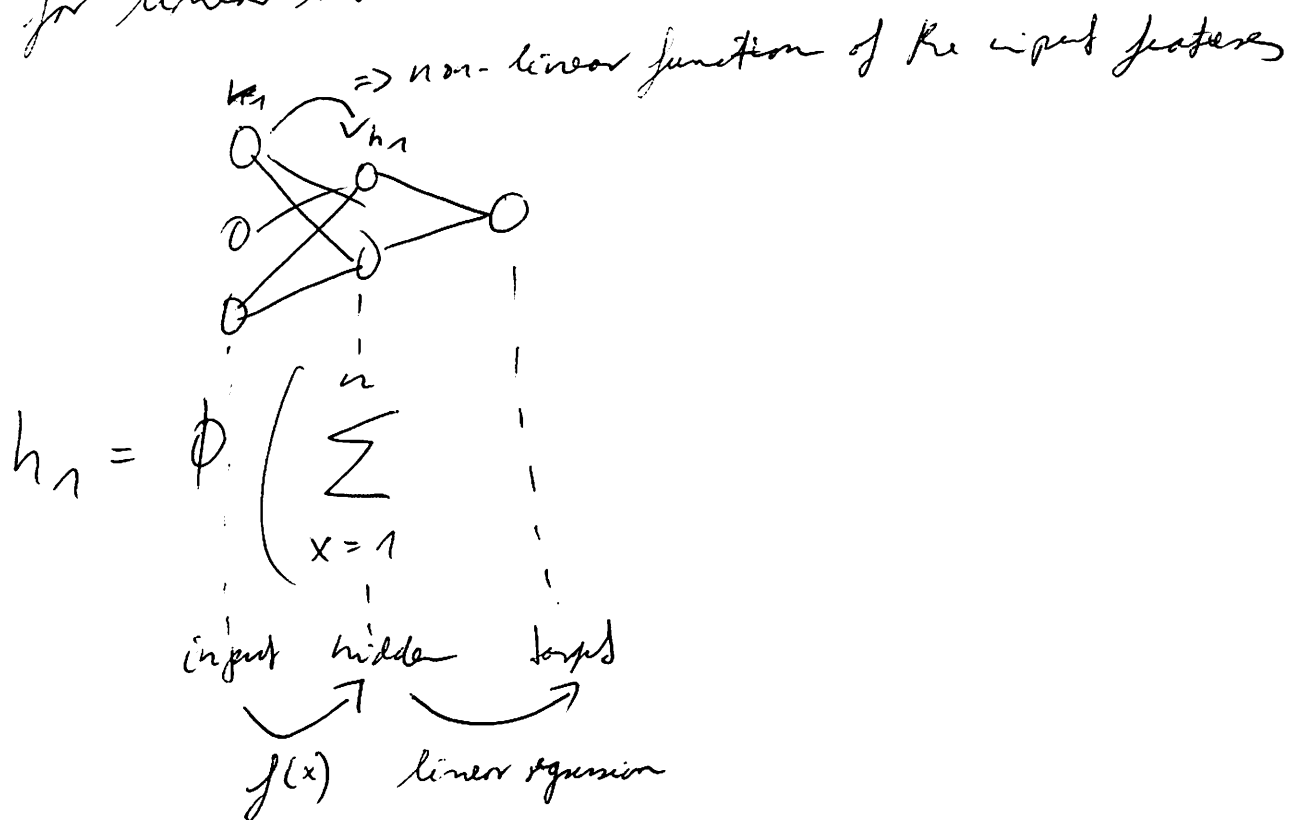must be used.

## Polynomial Regression

transform the input features

| X | y |
|---|---|
| 3 | 7 |
| 2 | 6 |
| 4 | 11 |
| 6 | 6 |

=>

| x | $x^2$ | $x^3$ | y |
|---|-------|-------|---|
| 3 | 9 | 27 | 7 |
| 2 | 4 | 16 | 6 |
| 4 | 16 | 64 | 11 |
| 6 | 36 | | |

projection from 2 Dim to 3 Dim
linear model and transform it
back again to 2 Dim

## Neural Networks

are used to get the functions for the transformations for linear models

$\Rightarrow$ non-linear function of the input features

$$h_1 = \phi \left( \sum_{x=1}^{n} \right)$$

input    hidden    target

$f(x)$    linear regression

weights          coefficients
=                =
parameters of the     coefficients of the
transformation       linear regression


weights are taken ~~by the~~ from the data !

## Error Distributions

$$\hat{y} = \underbrace{f(x)}_{\text{point estimate}} + \underbrace{N(0.6)}_{\substack{\text{Error Distribution} \\ \text{Function}}}$$

point estimate
from NNET()
or LinReg,
Poisson Reg,
Poly Reg

### Example

$$x = \text{wind} = 20$$
$$f(x) = \text{Point Estimate of Ozone} = 7$$

$$\hat{P}(y) = f(x) + N(0.6)$$

Because we know that the point estimates are not correct. Is there a way to give a range of possible correct values? Yes: Error Distribution around the estimate.

The error distribution is just the distribution of the residuals. We take the standard deviation (SD) of the residuals.

## How to control complexity? Use Regularization

- made it difficult for parameters to take all values
- just bind the to a range of values
- this will prevent us from the need to throw away variables / features

$$\hat{y} = f(X, \beta)$$

$\beta$ = minimize        + Sum of the        $\boxed{L1}$
the MSE              absolute values
( Mean Squared Error)
                    $\underset{\downarrow}{}$
                    regularisation function
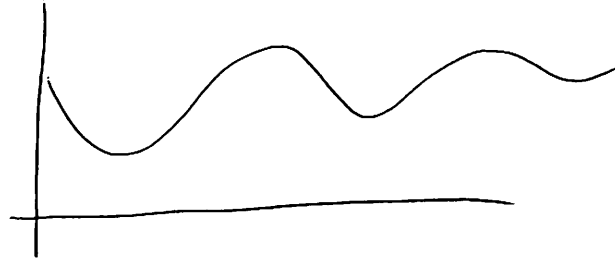
Do: penalise the parameter values

$L1$ = continuous feature selection

$L1$ = some parameters go to zero and stick,
thus they can be thrown away

$L2$ = some parameters go to zero but rise again,
thus all variables needs to be kept

Complexiness = Unregelmäßigkeit der Kurve

BSP:



## NNeb

decay = l (lambda) = L1/L2 regularisation für

Neural Networks

How do we find the right model?

Build different models with different regularisation!
See how low do the models perform on the
evaluation ~~in test test~~ test data.

Wenn der Suchraum gefunden wurde dann
Grid Search angewendet werden.

Lambda is not a parameter of the model itself
But of the function that call it. It is therefore a
hyperparameter. Iterations will produce different
models if lambda is adjusted.

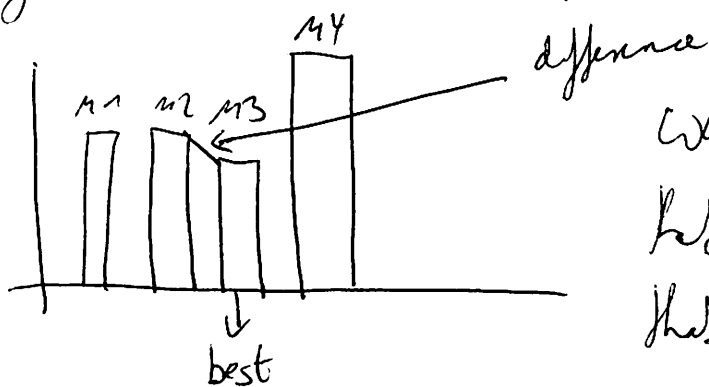max it => early stopping in iterations for
neural networks

If we have more data, the best models
become the more complex models.

## How to split your data?

Is there a objective way to split the data?

① Use learning jouls

② Significance test for evaluating model
performance.



Is M3 the best model?

What is the probability
that M3 is better assuming
that in the long run all
models perform best.

# ML Bootcamp

Today: Feature Selection    } Pre-Processing
Feature Transformation

Target: the find minimum of feature with
maximum information

= smallest number of most informative variables

Ways to do it:

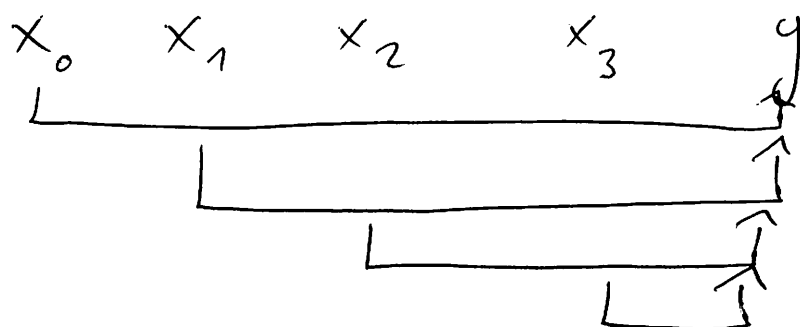① Business Expert knowledge

② Be careous: Experts can be wrong!

② Pairwise statistical analysis  ⟵————————

③ Model validation, in extreme 1 model for
every subset of features
— theoretical ideal
— computational difficult, thus ————

## Pairwise Statistical Analysis

$$X_0 \quad X_1 \quad X_2 \quad X_3 \quad y$$

What is the information the $x_n$ gives to explain $y$?

| Case | feature | | target | |
|------|---------|---|--------|---|
| ① | Num | | Num | → 1. Approach |
| ② | Num | | Cat | |
| ③ | Cat | | Cat | |
| ④ | Cat | | Num | |

2. Approach
3. Approach

(Num - Num)

Case ① : <u>simple correlation</u> = Pearson's corr coeff

$x \to y$
$x^2 \to y$  — for not linear correlation,
$x^3 \to y$  ~~not dimensional~~ and take the maximum score

Use corr.test() in R !

Case ③ :        Mutual Information ( Cat - Caf )

Entropy = measure of randomness



Interval of 3

- measure the information that x and y share
- if we know ① is weight pdf of all people and
  ② is weight of chinese people, we can reduce the
  value range dramatically.

Case ② + Case ③ : ( Cat-Num / Num-Caf )

Conditional Distribution Divergence ( CDD )

— the less the distributions overlap, the better if
  we estimate what the value is

— if normal distribution : mean, sd, variance

— Method : Bhattacharyya distance

## Excercise : „Feature Selection . $R$"

### Duncan dataset

X:
Type (cat)
Income (num)
Education (num)

Y: Prestige (num)

① first only numerical

② bin income input feature

③ bin target variable

## Feature Transformation

changing the coordinate system

### ① Center and Scaling Variables

it's important because the algorithms pay attention to the range of values which might be by accident wider or narrower.
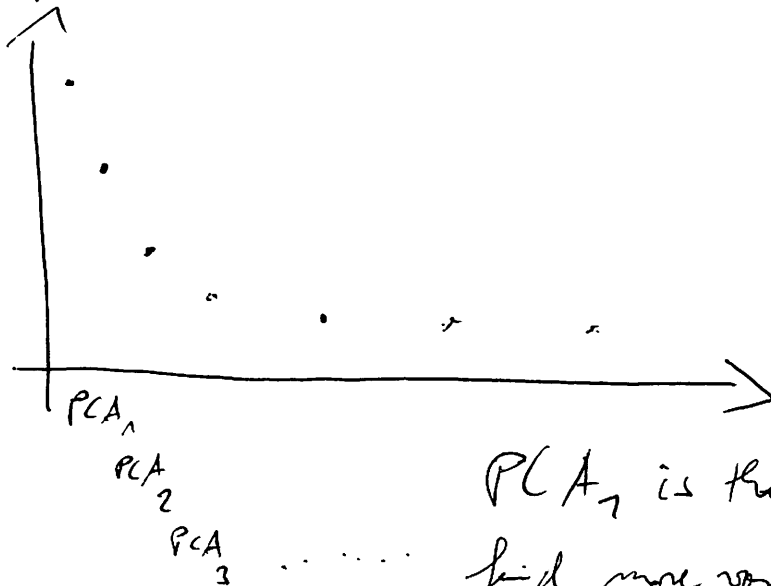
SSD: — mean is „0".
↓
R:: scale()

— the value is how many SD the real value is away from the mean.

— the data cloud is the same as with original values

— all values are around „0"

② PCA = Principle Component Analysis

- Scale and center all values before using PCA !!!  (V)
- gives us the maximum variance of the variables

- R :: prcomp()

- Assumption: the more spread out the variable values are, the more information there is

- Throw away all but the first few variables of the PCA.

- PCAs contain the information how many variance
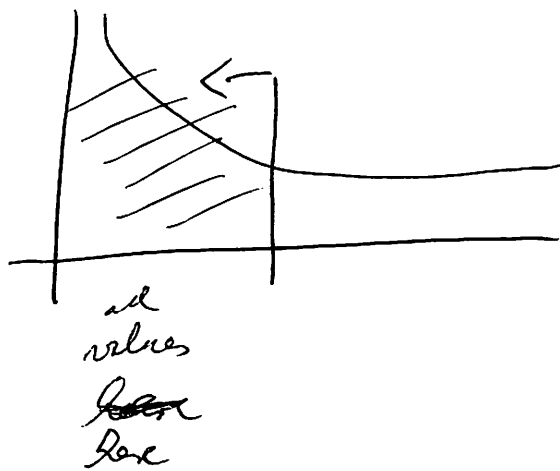  = information each component has



PCA₁
  PCA₂
    PCA₃  . . . . .

$PCA_1$ is the direction where to find more variance.

- Most of the time the first 10 PCA's should be chose cause there is a quick drop out of components.

How to decide how many features to
keep?

Rule of thumb: all value left of a angle
of 45° of PCA curve



all
values
~~keep~~
here

pr comp ( scale = TRUE )
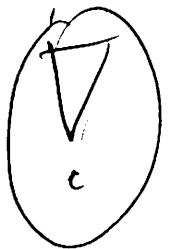


scale should always be TRUE
center = TRUE by default

predict can transform all new data
according to the transformed data:

predict( myPCA, newData )

newData is transformed the same way as
myPCA. Instead of myPCA any other
transformation can be used.

# ML Bootcamp

Hastie, et al.:

○ Elements of statistical learning
  supervised + unsupervised

○ Neural Networks and learning machines

○ Pattern recognition          unsupervised

(LiT)

supervised

---

Excercise: Feature Transformation

○ pca. $R^c$

---

## Methods Attributes

type of input variable ( nom, cat, ord, num)

type of target variable    ( dtv. )

Assumptions            Subject to overfitting

basic shape of data

concept (theory)  , implementations names(in R,
calculation procedure      Python, Julia, Scala, Java),
outcome(s)          gif. algorithms or links to
                    them in other books

interpretation of values

target of the question / use case

white box / black box method

at least interpretable, could be filled also with business knowledge

black celling

## Classification

on categorial target variable

MSE = Misclassification Error

1. Technique: QDA Quadratic Discrimination Analysis

                    LDA
2. Technique: Linear Discriminant Analysis

LDA + QDA work well on small data

3. Technique: Logistic Regression

linear classifier : where

$P(Y=1)$       $P(Y=0)$

only for use cases where y can be
only „0" or „1".

Neural Networks: sequence of linear regression

Neural : sequence of logistic regression
Networks for binomial target variables

GLM : link function, binomial & instead
of, poisson for linear regression

Neural Networks for Classification

To Do

- Build a landscape of methods

- Gaps a matrix of methods
  using the attributes on page 8 of 2017-08-17

## SVM    Support Vector Machines

It is easier to draw a ~~line to build a~~ linear separal
groups if there are more dimensions.

If the data is not linear separable we introduce
a cost factor „c" and optimise the data from
there taking into account the cost „c".

A linear kernel delivers #a separation on
the original data.

Bisher habe wir „hyperparameters" behandelt habe,
von, indem wir Sequenze von mögliche werte definiert
habe und die Modelle mit alle ~~profit~~ definierte
werte zu sehen und die Performance zu teste.

De die SVM's ~~hypra~~ hyperparameter benötige,
kann man auf vorberechnete SVM-Modelle z.B. von
Google zurückgreifen.

It is essential to select the hyperparameters very carefully.

---

Everything that's going on is using Euclidian distance.

---

In SVM is choosing the Kernel is very important as well as choosing values for the hyperparameters of the chosen kernel. It is uncommon to change the distance function from Euclidian to something else although it would be possible.

Exkurs:

## Techniques für Time Series Date

simple

① ARIMA

② HIDDEN MARKOV MODELS (HMMS)

③ Dynamic Baysian Network → white box technique

④ Recurrent Neural Nets → black box techniques

complex

If you want to be sure to work right on real
data use ③ or ④. If you need to put in
business expert knowledge ~~into~~ to it or need to
interpret the results the only model is ③ Dynamic
Baysian Networks.

## Confusion Matrices, Precision, Recall

used for to evaluate how well a model
is doing as an alternative with MSE

Mean Square Error.

Confusion Matrix, Balanced ~~Accuracy~~ Accuracy + ROC
Curves

"Classification. R

## Decision Trees

- Decision Trees will not be used on real-world data.

- Cause they perform too bad.

- There are more useful methods like

  a) bagging

  b) boosting

- These other methods build on „Decision Trees" though.

(+) for trees we do not need to do any feature selection

## Tree bagging

- We reduce the variance and not ~~effecting~~ effecting the bias.

- If we to reduce the correlation between trees we are pushing down the variance in all trees.

$$\downarrow$$

## Random Forest    (Forest = bunch of decision trees)

- Trees only on a small random sample of variables.

- This way we get very different trees.

---

- Random Forest perform pretty well.

- By using more trees, Random Forest does not overfit.

Evaluation:

Tests — Out put -bag error estimate
to & test data

↳ we do not need to split
between training and & evaluation
data

(+) • only one parameter to tune, which can be done
using a function
- simple to do and understand
↓
tune RF()

- much more simple than Neural Networks or SVM,
although NN and SVM might perform better,
but fine tuning takes a long long time.

- can be massively ~~parallisted~~ parallised

(−) • you need <u>feature selection</u> cause it looks
only at a small subset of feature. So if not
done a lot of trees will be build of data with no information,
thus producing a lot of ~~noise~~ noise.

## Boosting with AdaBoost

- increasing the number of trees will lead to overfitting (V 6)

- Can not be parallelized, cause it works sequentially

- a good package for boosting is "gdm" package