# Machine Learning in Deutsche Telekom



AI

SHARED
SERVICES

# DTSE in Europe

### GERMANY

**Cologne**
1007

**Darmstadt**
388

**Hamburg**
146

**Leipzig**
361

**Nuremberg**
130
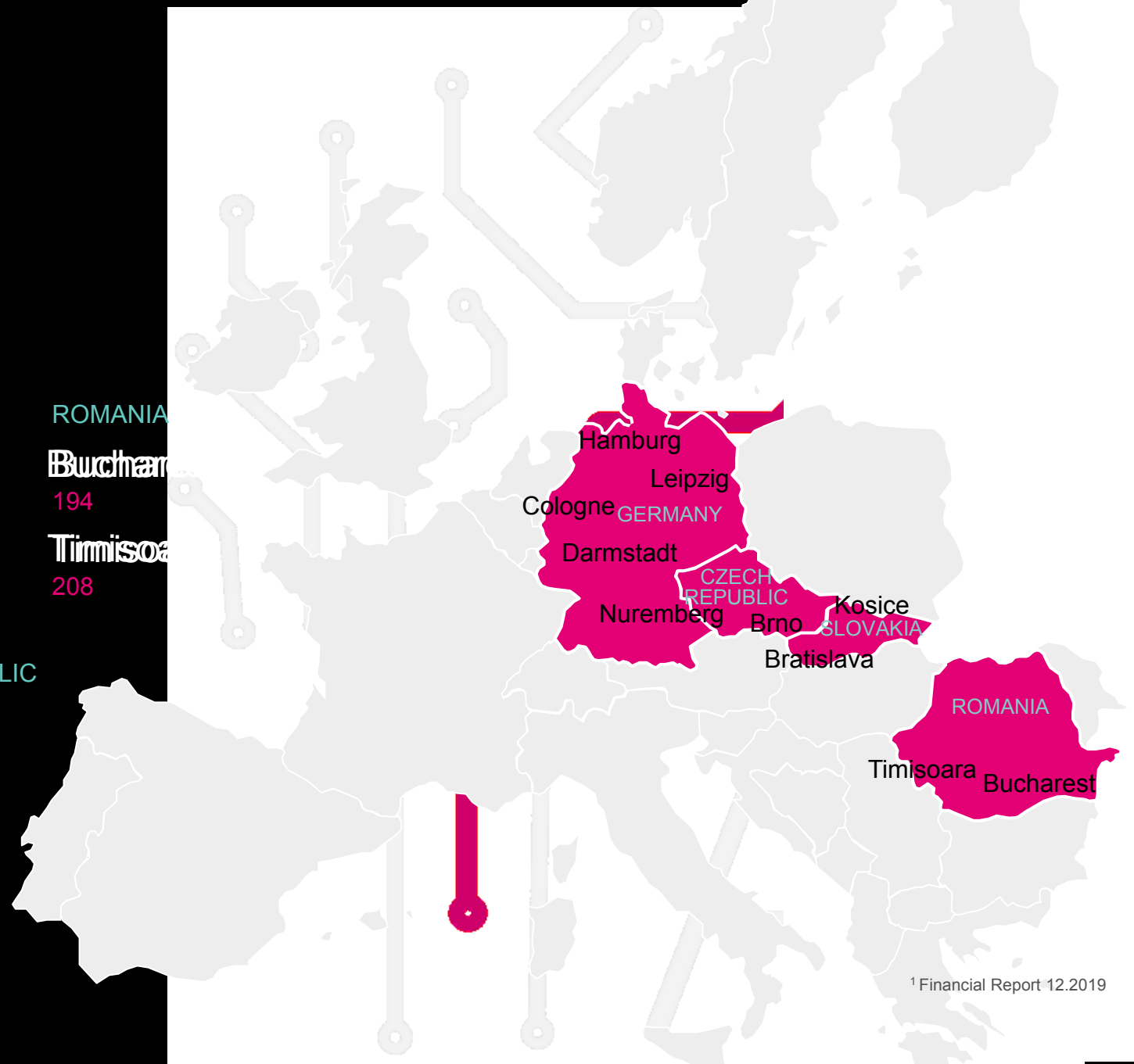
### SLOVAKIA

**Bratislava**
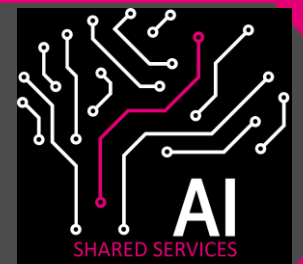617

**Kosice**
237

### CZECH REPUBLIC

**Brno**
133

### ROMANIA

**Bucharest**
194

**Timisoara**
208



Hamburg
Leipzig
Cologne GERMANY
Darmstadt
CZECH REPUBLIC
Nuremberg  Brno
Kosice
SLOVAKIA
Bratislava
ROMANIA
Timisoara  Bucharest

[1] Financial Report 12.2019

**T···**  LIFE IS FOR SHARING.

# AI Shared Services: your partner for the leading digital telco

## We embrace AI and agile working to enable data-driven decision making

**> 40**
INTERNATIONAL TECH TALENTS

**> 120**
AI COMMUNITY MEMBERS

## Our recipe for great AI projects

AGILE WORKING & RAPID DEPLOYMENT

EUROPEAN CROSS-FUNCTIONAL TEAMS

SMALL INVESTMENT & END-TO-END SERVICE

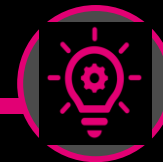BEST TECHNOLOGY & CONTINOUS IMPROVEMENT

Follow us on LinkedIn

TEXT PROCESSING

PREDICTIVE ANALYTICS

ESG ACCELERATION

AND MUCH MORE...

**Jakub Kondek, M.Sc.**

[j.kondek@telekom.com](mailto:j.kondek@telekom.com)
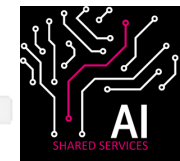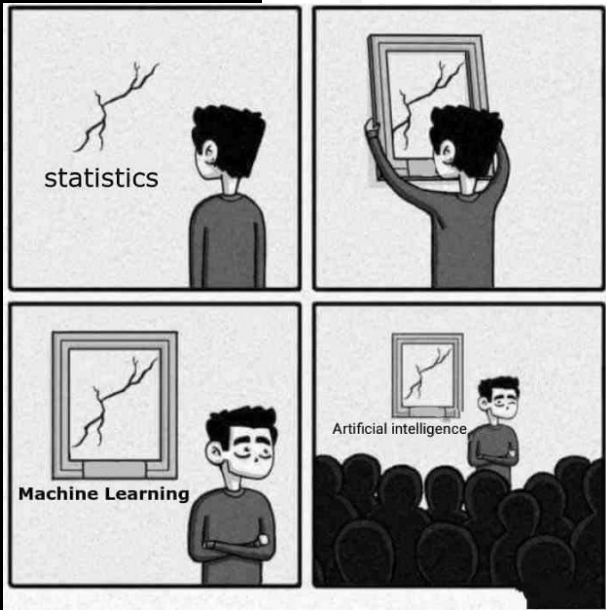
Senior Data Scientist

Turista

Blues enthusiast

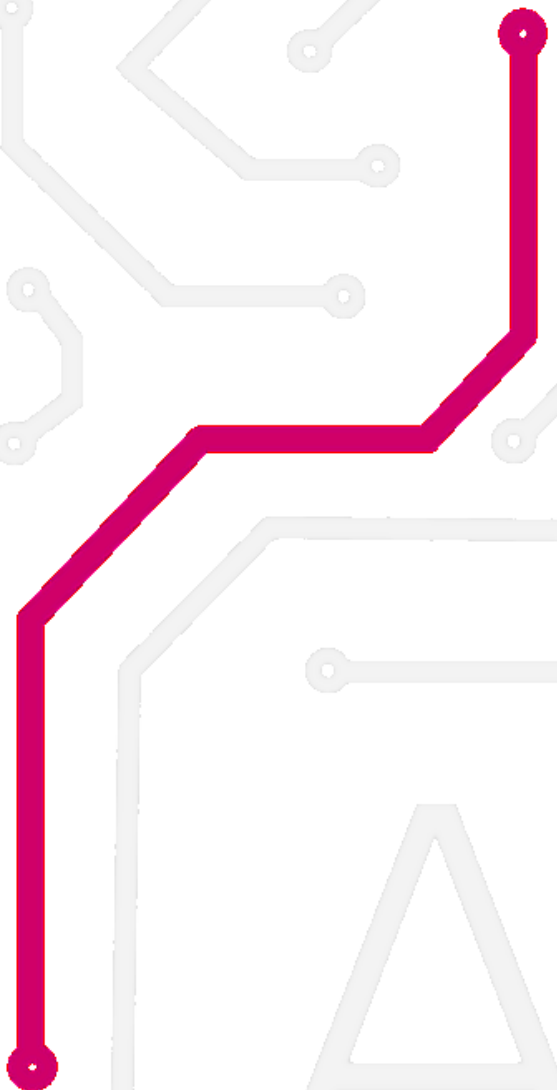[stepan.vondracek@telekom.com](mailto:stepan.vondracek@telekom.com)

LIFE IS FOR SHARING.

AI
SHARED SERVICES

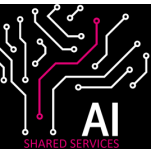https://www.instagr[...]or[...]s/

# DAILY Challenges of Data Scientist

- A customers do not understand what ML does（can do）

- Communication of technological topics

- Unclear requirements *("we would like to implem... some AI...")*

- Stated problems are much simpler

- System integration
  - Is data which was used for training available... predicting
  - Do I have access/rights to data

- Can I send prediction somewhere reasonably

  \*\*\*

  Possible challenge for statistics/econo... graduate -> mostly IT terminology

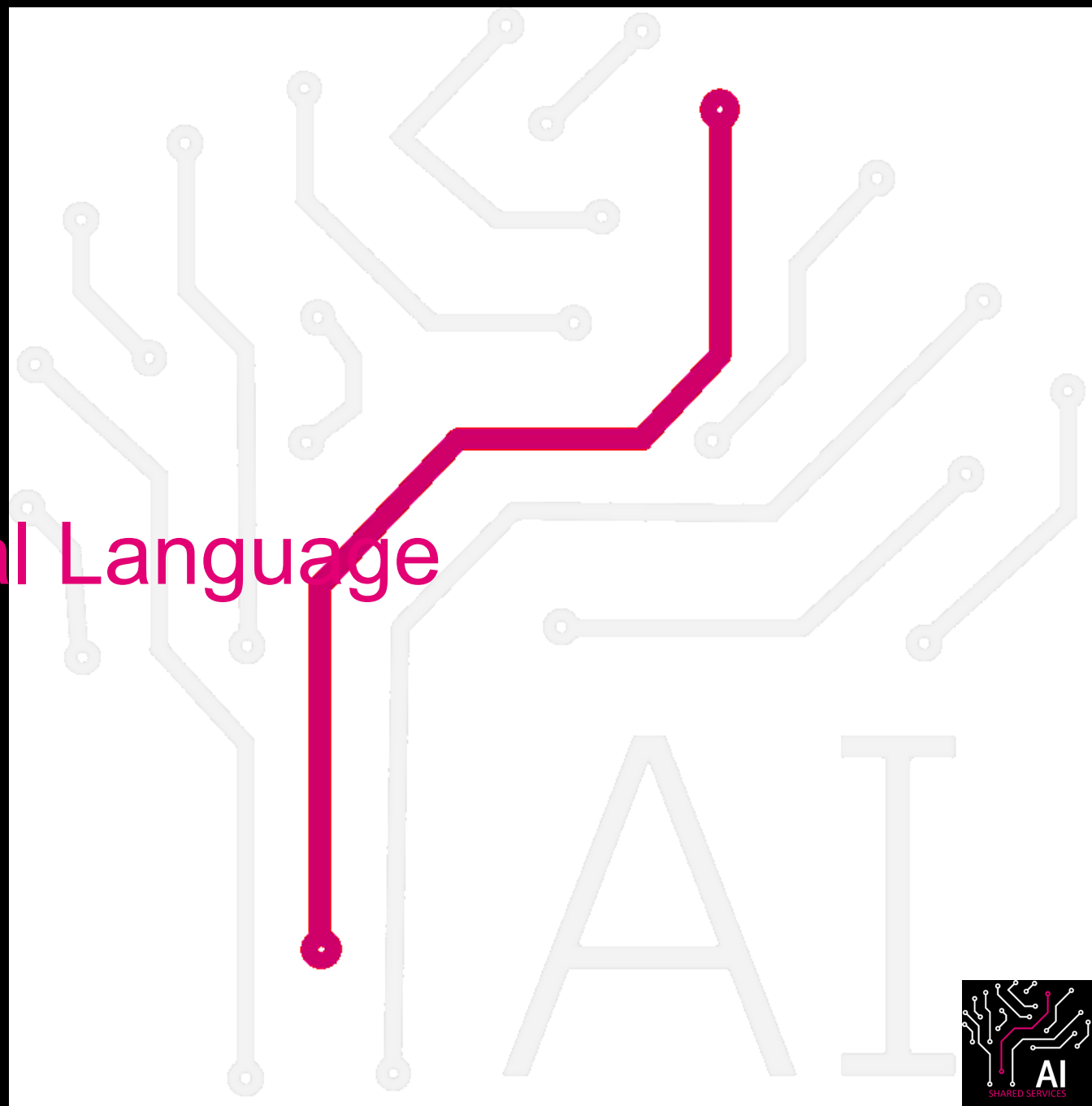# What We Use in AIS

- Python，SQL

- Python IDEs （Pycharm，Visual Code）

- Dedicated ML server + Cloud (Openshift

- containerization (Docker, Kubernetes)

- REST API for communication between th

- Webservice Frameworks (fastapi)

**T** · · ·    LIFE IS FOR SHARING.

# Part I. - Natural Language Processing

# NATURAL Language Processing (NLP)

- A field of AI dealing with interaction between computers and humans using the natural language

- Started in 1950s（Turing test）
- Considered to be a difficult problem in computer science
- Research in NLP is still going on
- 3 main NLP subfields:
  - Speech Recognition
  - Natural Language Understanding
  - Natural Language Generation

LIFE IS FOR SHARING.

# NLP Model Types (eras)

- Rule-based

- Statistical — the "traditional approach"

- Deep learning — the state of the art, "modern" appro

LIFE IS FOR SHARING.

AI
SHARED SERVICES

# Rule-Based Modeling

- A hand-crafted system of grammar rules based on lin
  - regular expressions, context-free grammars
  - often requires a skilled expert — a linguist
- Useful when we don't have enough data
- Very good interpretability
- Poor generalization and maintenance

  - E.g. rule-based chatbots

# Statistical Modeling

- Traditional ML models – training and testing data

- Requires moderate amount of data with annotations

- Heavy on feature engineering
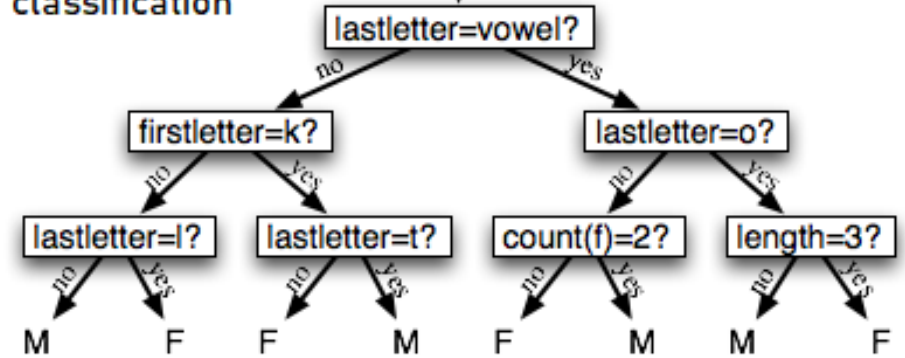  - word frequency, number of characters, edit distance capitalized, plural etc.

- Linear classifiers, Decision trees etc.


- Language Model
  - a probability distribution over sequences of words
  - can be also used for Language Generation

- N-gram Language Model
  - N-gram frequencies pre-counted on training corpus
  - P("closer"|"getting a step") > P("coffee"|"getting



Name gender classification

LIFE IS FOR SHARING.

# Deep Learning

- Feature engineering is generally skipped
  - raw data as an input（word embeddings）
  - network learns important features itself
- Large training corpus
- Good generalization
  - transfer learning – reusing models trained on different tasks
- Poor interpretability
- Sequence models
  - Recurrent Neural Networks（RNN，LSTM，GRU）
  - Temporal Convolutional Networks（TCN）
- Tasks
  - Classification，Regression，Sequence-to-sequence

# Selected "methods" of working with text

# BAG-of-words (BOW)

- A simplifying representation, disregards word order

- Text is represented as a bag（multiset）of its words

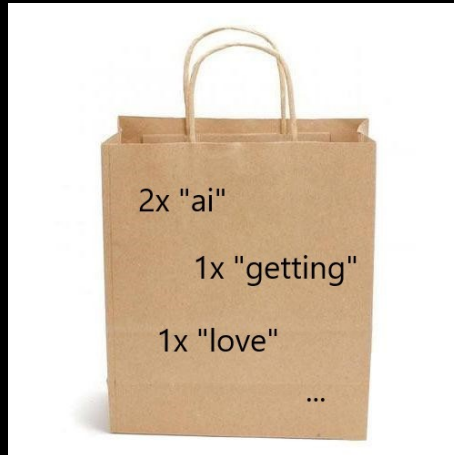  - multiplicity － number of occurrences of each word

2x "ai"

1x "getting"

1x "love"

…

```python
from sklearn.feature_extraction.text import CountVectorizer

vectorizer = CountVectorizer(tokenizer=nltk.word_tokenize, vocabulary=dictionary)

bow = vectorizer.transform([sentence])

list(zip(dictionary.keys(), bow.toarray()[0])) #word counts

[('.', 2),
 ('a', 1),
 ('ai', 2),
 ('am', 1),
 ('closer', 1),
 ('getting', 1),
 ('i', 2),
 ('journey', 1),
 ('love', 1),
 ('my', 1),
 ('on', 1),
 ('step', 1)]

bow.toarray() #bag of words representation

array([[2, 1, 2, 1, 1, 1, 2, 1, 1, 1, 1, 1]])
```
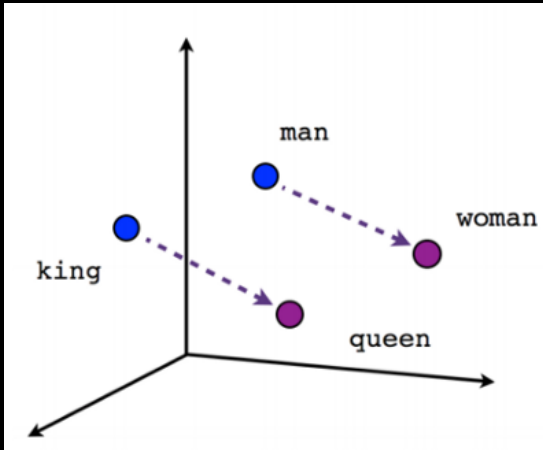
LIFE IS FOR SHARING.

AI
SHARED SERVICES

# WORD EMBEDDINGS

- Mapping of words into vectors of real numbers

- Words are closer to each other, if they occur in similar context

- King − Man + Woman = Queen

- Most popular embedding models are Word2Vec (Tomas Mikolov), FastText, or BERT



```
import fasttext

model = fasttext.load_model("./fastText/cc.en.300.bin")

cosine_distance(model.get_word_vector("journey"), model.get_word_vector("voyage"))

0.42419618368148804

cosine_distance(model.get_word_vector("journey"), model.get_word_vector("coffee"))

0.9299457967281342
```

# LLMs @AI Shared Services – Product Athena

AI
SHARED SERVICES

# Athena Overview

- Semantic search over various data types

- Enables fast orientation in vast amount of data

- Features:

  - Language independent

  - Summarization

  - Similarity matching

  - Text generation

  - Connection to web data



**Options**
Show debug info

Which database you want to use?
riskdata1

**File Upload:**

Drag and drop files here
Limit 200MB per file • JPG, JPEG, PNG, GIF, TIFF, PDF, DOCX, XLSX, PPTX, TXT, CSV, DOC, RTF

Browse files

Answers in audio (only english)

**Athena - semantic search**

Which sports teams does Telekom sponsor?                40/100

Options
Max. number of documents from retriever         Answer mode
10
1                                    10         Expert

🔍 Search

Telekom sponsors several sports teams, including the German record champion FC Bayern Munich, the German Football Association (DFB), Hamburger SV, Borussia Mönchengladbach, and 1. FC Köln [Document 2]. They also have partnerships with the National Men's Football Team of Germany, the 3. Liga, and the FLYERALARM Frauen-Bundesliga [Document 7]. Additionally, Telekom is involved in the eSports scene and sponsors SK Gaming [Document 8].

👍 👎  ✏ Specify

Part II. - Time Series Forecasting in DTSE

# Jak odhadnete počet lidí, kteří přijdou na naši přednášku?

- Počet lidí na semináři minulý týden

- Počet lidí ve druhém týdnu loňské

- Zohlednění počtu studentů v ročn

- Geniální anotace, která určitě přilá

# Predikce časových řad

1. PRAVIDELNĚ uspořádané hodnoty

2. Predpoklad, že (některé) aspekty z
   pokračovat v budoucnosti a ovlivní

3. Možnost zahrnout další vysvětlující

- V Deutsche Telekom hlavně měsíčn

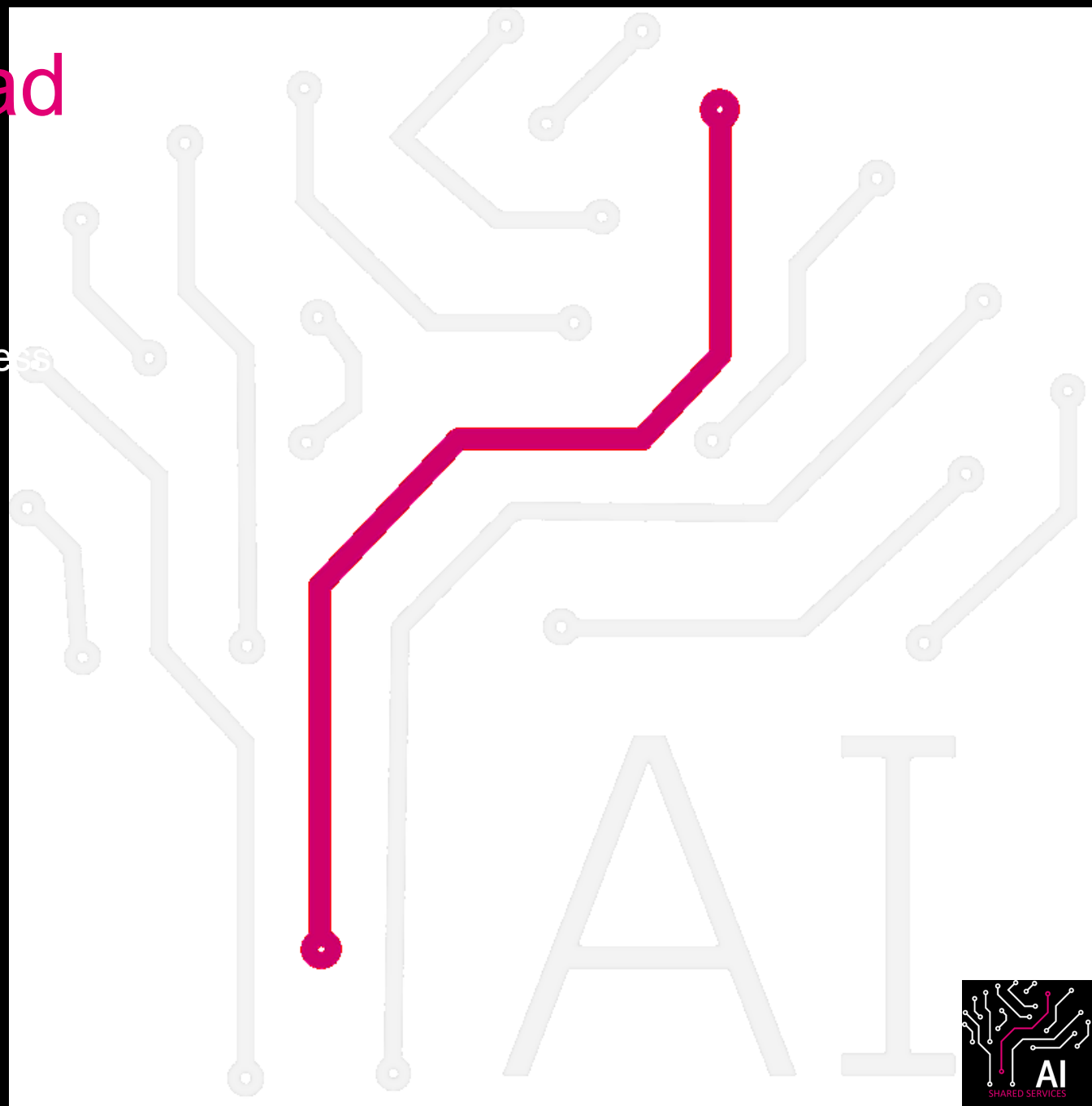| Date | Position | Value |
|------|----------|-------|
| 2022-01-01 00:00:00 | Revenues | 104 |
| 2022-02-01 00:00:00 | Revenues | 101 |
| 2022-03-01 00:00:00 | Revenues | 103 |
| 2022-04-01 00:00:00 | Revenues | 108 |
| 2022-05-01 00:00:00 | Revenues | 112 |
| 2022-06-01 00:00:00 | Revenues | 115 |
| 2022-07-01 00:00:00 | Revenues | 129 |
| 2022-08-01 00:00:00 | Revenues | 134 |
| 2022-09-01 00:00:00 | Revenues | 127 |
| 2022-10-01 00:00:00 | Revenues | 116 |
| 2022-11-01 00:00:00 | Revenues | 111 |
| 2022-12-01 00:00:00 | Revenues | ? |
| 2023-01-01 00:00:00 | Revenues | ? |
| 2023-02-01 00:00:00 | Revenues | ? |

# Proč?

- Zlepšená schopnost rozhodování
- Zlepšená schopnost plánování
- Konkurenční výhoda
- Řízení rizika

# Predikce časových řad

- Konvenční postup:
  1. Definice problému
  2. Sběr dát a získávání insights od business
  3. EDA – vizualizace, popisné statistiky
  4. Preprocessing
  5. Výběr modelů a jejich hyperparametrů
  6. Cross-validace
  7. Výběr finálního modelu (ensembling?)
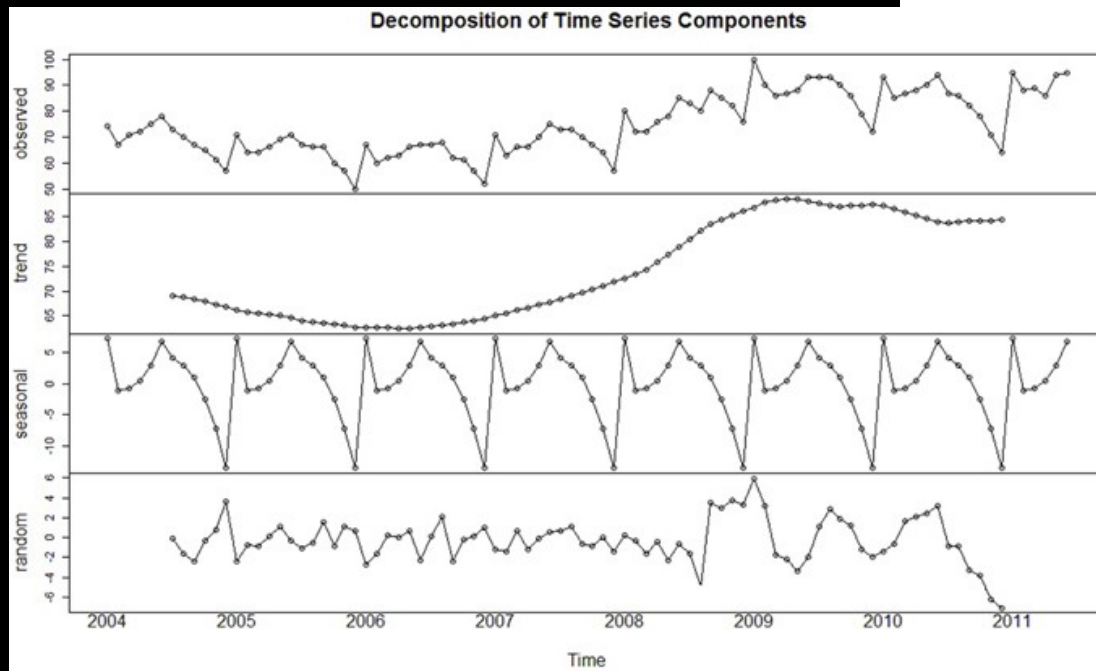  8. Predikce
  9. Postprocessing
  10. Delivery

**T** • • •  LIFE IS FOR SHARING.

AI
SHARED SERVICES

# Komponenty časových řad



Revenues [babickine_halusky ™]

# Komponenty časových radov



Decomposition of Time Series Components

**Trend**
**Seasonality**
**Residuals**
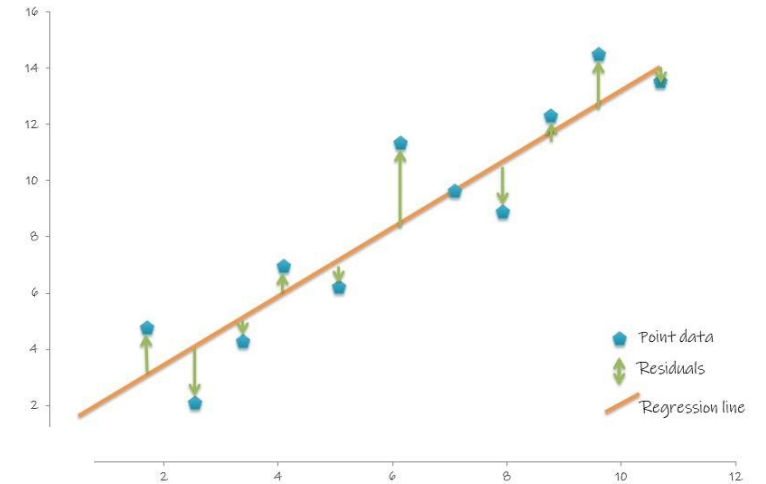
**Monthly Sales in U.S. Beer, Wine and Liquor Stores 1992-2017**
in million US$/month (NAICS 4453), dotted line: mov 12-month avg
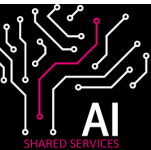Source: Bureau of Census
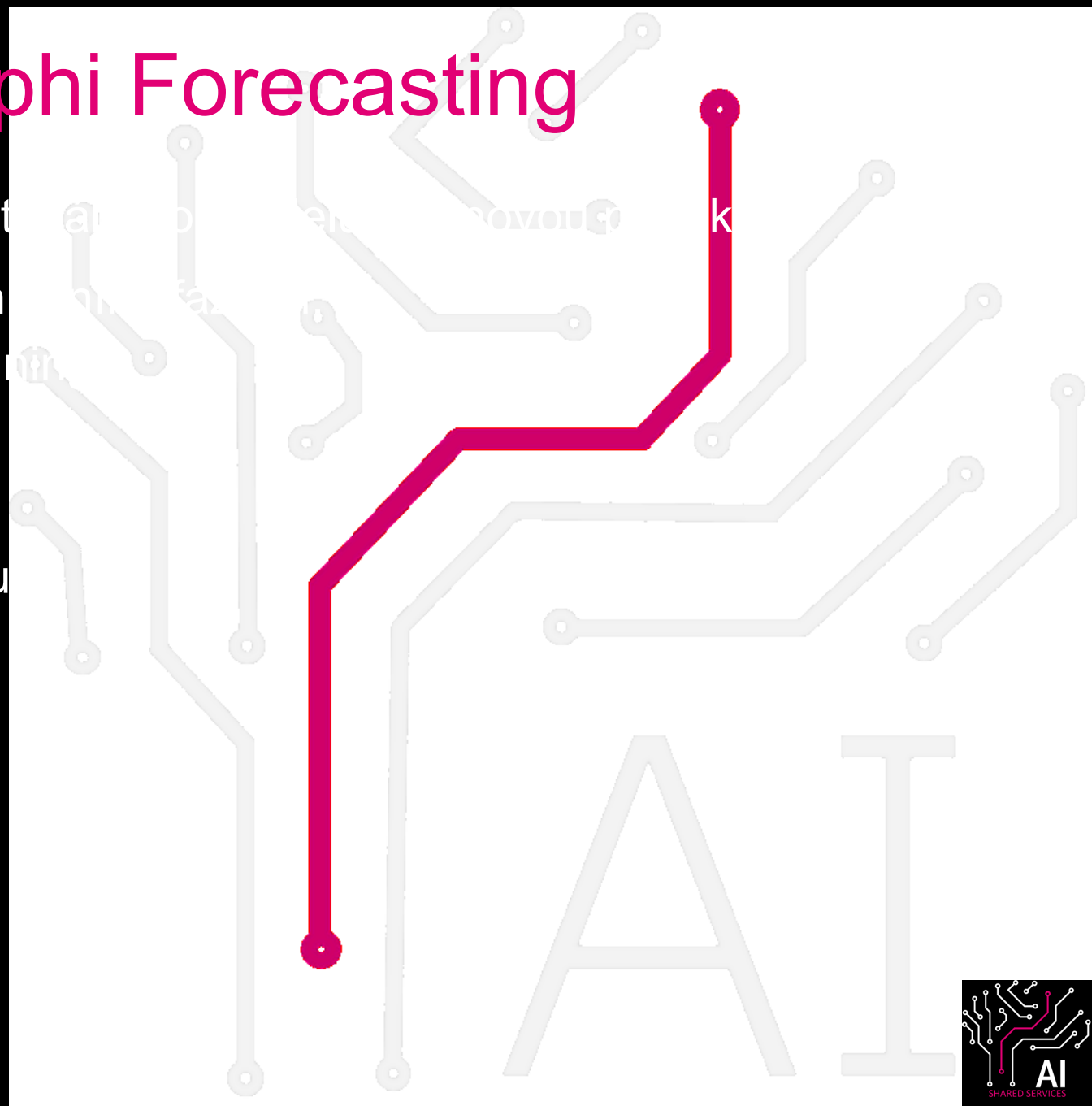
# Chyba predikce

$$e_t = y_t - \hat{y}_t$$



Point data
Residuals
Regression line

- Common error metrics:
  - MAE (Mean Absolute Error)          $MAE = mean(|e_t|)$
  - MSE (Mean Squared Error)           $MSE = mean(e_t^2)$
  - RMSE (Root Mean Squared Error)     $RMSE = \sqrt{mean(e_t^2)}$
  - MAPE (Mean Percentage Error)       $MAPE = mean(100 \cdot e_t / y_t)$

# Auto ML řešení - Delphi Forecasting

- Každý měsíc probíhá nový výběr opt                    k

- Delphi algoritmus prechádza třemi h
  - Model training ('hyperparameter tur
  - Model selection / evaluation
  - Prediction

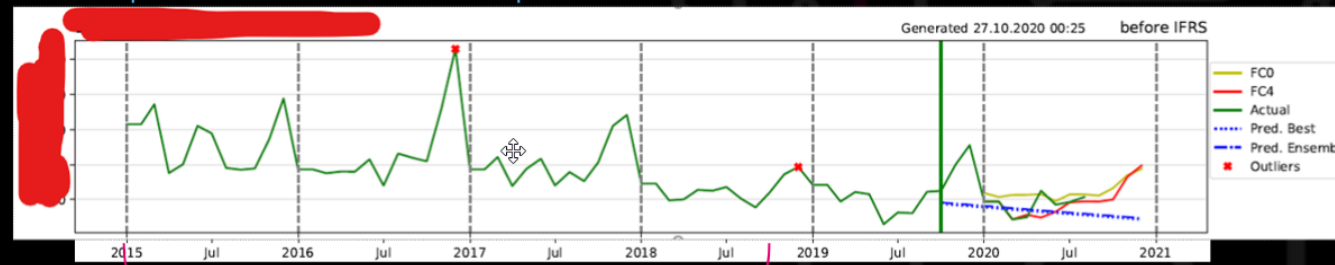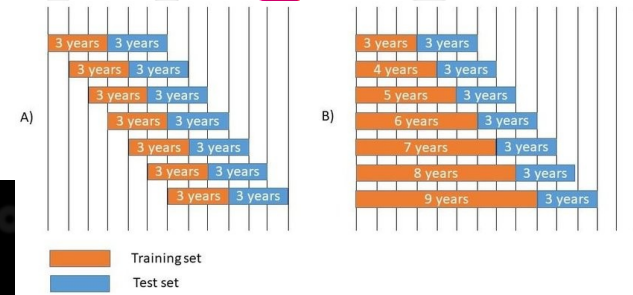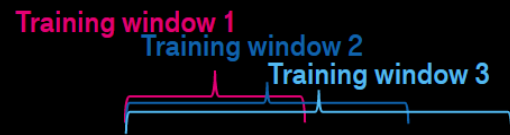- Separátně pro každou časovou řadu

AI
SHARED SERVICES

# Používané technologie

# Tréninování Modelu

- Každý model má hyperparametry, které s... brát vysvělující proměnné?)

- K rychlejšímu výběru parametrů využívá... hyperparamery, které se nezdají, že by m...

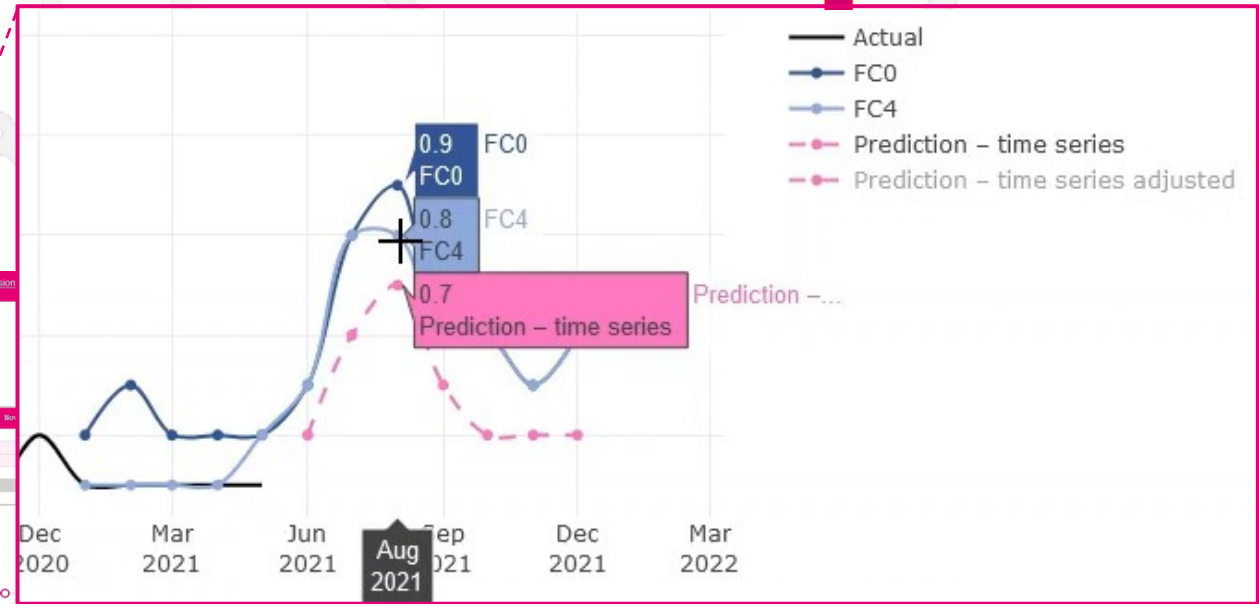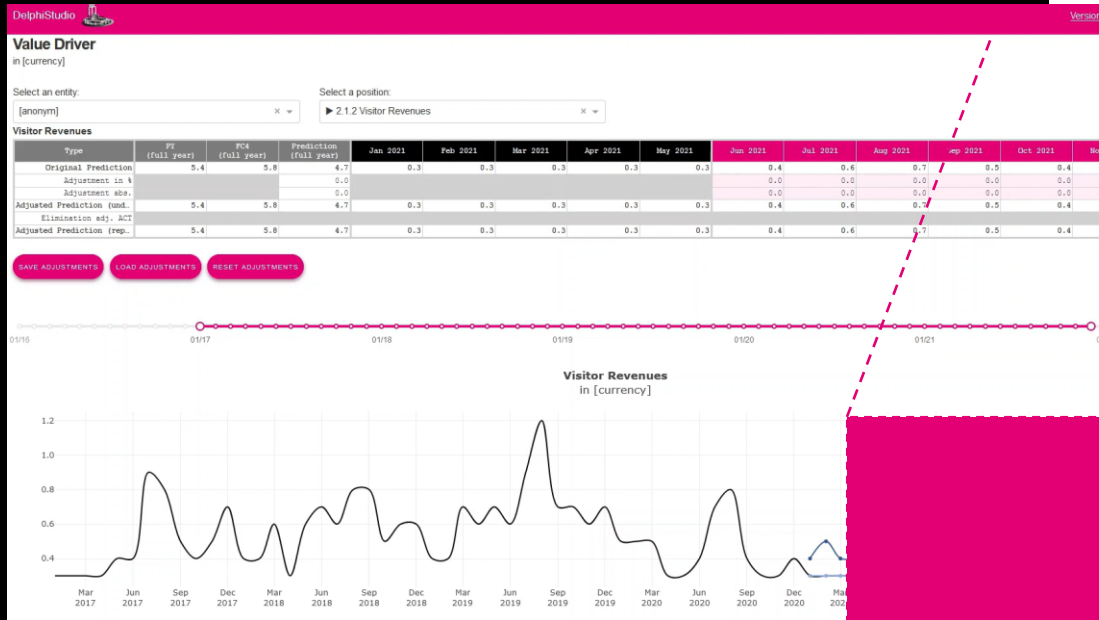- Model se fituje na základě zvolené komb... a následně ... error metriky

# DelPy TDG Live demo

# Delphi Studio tour (1/3): allows to visualize and compare historical, forecasted, and adjusted data

# Delphi Studio tour (2/3): adjust forecast to arrive at ML-assisted expert projection



| Type | Aug 2021 | |
|---|---|---|
| Original Prediction | 0.7 | |
| Adjustment in % | 0.0 | |
| Adjustment abs. | 0.2 | |
| Adjusted Prediction (und... | 0.9 | |
| Elimination adj. ACT | | |
| Adjusted Prediction (rep... | 0.9 | |

LIFE IS FOR SHARING.

# Delphi Studio tour (3/3): switch between various P&L positions

# ČÁST III. - Názory. Chcete někdo názory?

# Specializace je vždy relativní



## Data Scientist
also known as Data Managers, statisticians.

A data scientist will be able to take data science projects from end to end. They can help store large amounts of data, create predictive modelling processes and present the findings.

*Skills:* Mathematics, Programming, Communication

*Will use programmes such as:*
SQL, Python, R

## Data Engineers
also known as database administrators and data architects.

They are versatile generalists who use computer science to help process large datasets. They typically focus on coding, cleaning up data sets, and implementing requests that come from data scientists.

*Skills:* Programming, Mathematics, Big data

*Will use programmes such as:*
Hadoop, NoSQL, and Python

## Data Analysts
also known as business Analysts.

They typically help people from across the company understand specific queries with charts.

*Skills:* Statistics, Communication, Business knowledge

*Will use programmes such as:*
Excel, Tableau, SQL

**T** · · ·   **LIFE IS FOR SHARING.**

AI
SHARED SERVICES

# Naučte se psát kód koncepčně a tvořit úhledné projekty

- Čistý kód

- Vhodné pojmenovávání proměnných

- Modularita

- Pouze relevantní části

- Snadněji se tak buduje portfolio projektů

- Role engineeringu/architektury ML platform postup... ...AutoML frameworks a množství dotupných předko...

# Neztraťte se v záplavě "cool data science"
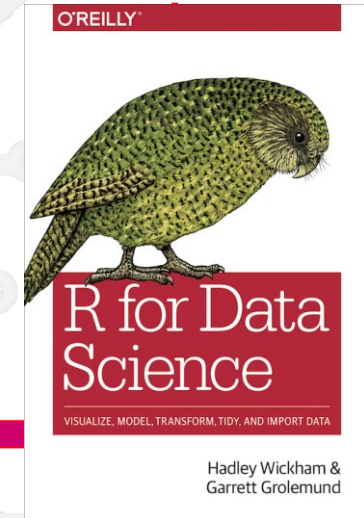
- Najděte balanc mezi užitečnými všeobecnými zna...
- Don't skip the basics
- Nejtěžší je získat svou první práci v oboru
- Neztraťe se v záplavě kurzů a materiálů, opravdu

# Zdroje ke studiu

- [Science as Amateur Software Development](#)

- [Value in Data Science Beyond Models in Production | RStudio](#)

- [Forecasting – Principles and Practice](#)

- [Hands on ML](#)

- [Andrew Ng Courses](#)

- [R for DS + tidyverse tutoriály](#)

- LinkedIn

# Závěr

- Data jsou a budou relevantní obor

- Nezapomínejte na základy - matematika, statistika

- Buďte trpěliví (při studiu, hledání práce, při práci

- Soft skills jsou v IT někdy trochu opomíjená dove

# OTÁZKY?