

# Feature engineering is your ticket to survival

@FilipVitek, Director Data Science

### Who the hell is Filip Vitek?







#### **TECH STACK** & kafka amazon REDSHIFT SQL Server AVRO Amazon S3 Parquet Spark Scala python

### 2 reasons why Feature engineering is king

kaggle Think of last time you were designing the Machine learning model. In which of the following steps did you rely on some of the ready-made packages? (e.g. SciKit Learn, ...) 5 2 1 3 XGBoost Model Model Variable winning many of the **Feature** Sampling fitting comparison selection transform. & Cross competitions on kaggle [DecTree, validation Algorithms are becoming \*\*\*\* commodities, you barely can beat others by "better algorithm" Rate of machine replacement quality compared to human

MORE

INFO

4



*"What alternative to AI do we, humans, have?"* http://mocnedata.sk/en/**what-alternative-do-we-have-to-ai**/



#### Most common FE approaches in business (coming soon)







### Taylor Polynomial ... as inspiration for feature creating

$$egin{aligned} T_n f(x;a) &= \sum_{k=0}^n rac{f^{(k)}(a)}{k!} (x-a)^k \ &= &f(a) + rac{f'(a)}{1!} (x-a) + rac{f''(a)}{2!} (x-a)^2 + \ldots + rac{f^{(n)}(a)}{n!} (x-a)^n \end{aligned}$$

Same principle	Example: Autonomous car
0] the variable itself	What is our actual speed?
1] <b>relative change</b> (in time)	Are we already breaking or do we accelerate?
2] changes tendencies	With full push to brakes how strong negative acceleration can we still achieve?

CAUTION! It is second partial derivative so it needs not to be dFx \* dFx, but can be dFx\*dFy as well.

6







#### **Common FE approaches in business** (and their pitfalls)





7



### Data underdogs ... and their impact



![](_page_7_Picture_3.jpeg)

![](_page_7_Picture_4.jpeg)

![](_page_7_Picture_6.jpeg)

### **Real examples of Unconventional Feature Generation**

![](_page_8_Picture_1.jpeg)

#### How old are you, Bernard?

- Nothing like "National ID" for German insurance companies
  = they have no clue about age of customer
- Important for setting proper communication (web vs. call vs. paper letter)
- First name + Region predicting 92% accurately the decade when the customer was born

[cut-off of approx. 20% Individuals]

9

#### Estimating age of car

- Modelling risk profile of the cars
- For many cars the age of the car was unknown in data
- Market relatively stable in past
- Duration of the contract with our company & age of owner

[even if a miss, still positively correlated]

![](_page_8_Picture_13.jpeg)

![](_page_8_Picture_14.jpeg)

Detecting **commercial** customer

- Quite a few small companies without license
- Too small to detect via IP address range
- Using standard desktop OS versions
- Pattern of use strong within working hours, weak outside

[nightmare of time zones from UTC]

#### Zodiac, are you kidding me?

- Probability to have car accident
- As Joker card for model
- Strong objection from Data Scientist "*This is not serious* work, we protest."
- Ended up as the Second strongest parameter in model.
- Later confirmed in 4 other countries in same issue

[I have a hypothesis why it works]

![](_page_8_Picture_28.jpeg)

![](_page_8_Picture_30.jpeg)

Summary of older EN blogs of mine http://mocnedata.sk/en/what-to-read-in-english-here/

![](_page_8_Picture_32.jpeg)

## Thanks for your attention and I am ready NOW to answer YOUR QUESTIONS

![](_page_9_Picture_1.jpeg)

www·TheMightyData·com

If you have Q's later ...

Mgr. Filip Vítek Data Science Director TeamViewer, Berlin

filip.vitek@teamviewer.com

![](_page_9_Picture_6.jpeg)

https://sk.linkedin.com/in/vitekfilip

@FilipVitek

(a)