

How To Do Cool Feature Engineering In Python

@FilipVitek, Director Data Science

Who the hell is Filip Vitek?









TECH STACK & kafka amazon REDSHIFT SQL Server **AVR** Amazon S3 Parquet 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, ...) XGBoost Your ticket to survival winning many of the 5 competitions on kaggle Model Model Feature Sampling Feature fitting creation comparison & Cross transform. & selection validation Algorithms are becoming commodities, you barely can beat Rate of machine replacement quality compared to human others by "better algorithm"

4



$$\bigcirc$$

... but you need to be true to your self on WHO YOU ARE

Analytics as Single continent



Analytics Archipelago



V - I - B - A

... the MBTI of the analytics. Find out which type of Data Scientist YOU are:

http://mocnedata.sk/en/VIBA-type-of-analyst/

What do I expect ... of COOL Feature Engineering



B Sorting out the hopeless cases



Prune to have more efficient model training & operation



Signal, if I missed anything relevant





Python is ... the Microsoft Excel[™] of our era



It Became standard. There are options, but why to bother to even try. (Think Lotus 1-2-3)

Everybody claims knowledge of it, but knowledge of most people is very shallow. (Frustrating to test Data Science hires for basic Python ... and see them struggle with RegEx)

Microsoft[®] Visual Basic for Applications Tool is really powerful, but you need to possess certain skills beyond elementary use.

To rely hone the power of it, you need to know more than the default options/libraries. (... which most people don't)

=1*(0.5-0.4-0.1)

VLOOKUP blinds (back search, MIX, Case Sensitive)

Z-score glitch

"Don't CHALLENGE or REVIEW, just CONSUME."

(Have you ever checked what Excel calculates? / Have you challenged any Sci-kit learn routines?)



SciKit Learn ... Our U-bahn of the Machine Learning (?!)

Supervised Learning

(GLM, LinDiscAnal, KernelRidge, SVM, StochGradDescent, NearestN, NaiveBayes, DecisionTrees, Feature selection, Ensemble methods, MulticlassAlgor, Isotonic Regress., Prob. Calibration, NeuralNetworks, ...)

Unsupervised Learning

(GaussianMixture, ManifoldLearning, Clustering, Biclustering, MatrixFactorization, Covariance estimation, OutlierDetect, DensityEstimate, NeuralNet, ...)

Model selection

(CrossValid, HyperParemeters, Model evaluation, Model persistence, Validation curves.)

✓ Dataset Transformations

(Pipelines, Feature extraction, Preprocessing, Impute, DimensionReduction, Projections, KernelApprox, PairWiseMetrics, TargetTransform)



(Toy/Real datasets, Generated datasets, Loading, Incremental learning, PredicitonThrouput, Parallelism)



Feature selection tools

- Low Variance Removal
- Univariate feature selection (Select K-Best)
- Recursive Feature Elimination (only backwards)
- SelectFrom Model (Tree)
- Including into Pipeline
- Principal Component Analysis
- Independent Component Analysis



How does SciKit Learn ... meet our expectations?



Sorting out the hopeless cases



Prune to have more efficient model training & operation



Signal, if I missed anything relevant

How to compensate for that in Python space ...



Build your OWN FEATURE engine

- Calculate variable statistics [see also transformations slide, ...]
- 2 Generate obvious suspects [aggregations, time windows, ...]
- 3 Indicate missing info categories [compare to dictionary, Expl. score ...]
- 4 Hard criteria knock-out [Variance, NonNulls, distinct X, ...]
- 5 Binning & Categorical decomposition [Forced binning if > N]
- 6 Univariate correlation & Log -P [Simple tree is enough]
 - Bivariate relations [Cut off for Categorical dummies by Support]
 - **Decision on ranking** of parameters [Simple, Stage based, ...]



7



Thanks for Your attention and I am ready to answer YOUR QUESTIONS

Feel free

to contact

me

in

Mgr. Filip Vítek

TeamViewer, Berlin

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

+ 49 1525 309 8505

Data Science Director

filip.vitek@teamviewer.com

@FilipVitek



www.TheMightyData.com

Join the community



http://mocnedata.sk/en/lets-keep-in-touch/

BACK UP SLIDES

Is Principal Component Analysis your Friend or Foe?



PCA

- Pure Machine to Machine interface
- Data-space visualization required
- Overcomes mutual correlations of features without even explicitly checking for them



- Feature selection procedure
 [even in SciKit Learn]. Reduction ≠ selection
- Humans using the result of predictions
- Had to do oversampling in process of the model preparation
- Neural network one of the rival models
- Non linear effects of the variables

Real examples of Unconventional Feature Generation

Unconventional approach to generating features



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 point for approx. 25% Individuals]

Fee increase tolerance

- Fee increase sensitivity for retail bank
- In search for metric that would tell: How "lazy" user is?
- Limited space, banking feels very un-emotional
- Lowest amount ever withdrawn from the ATM

[worked surprisingly well, due to large coverage]





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 Scientists: "*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]







Data underdogs ... and their impact

Who will win the car race to nearest lights?



Data fields that are "just identifiers"

Contact & Transactional data

No obvious relations as champion challengers (Joker cards)

Unusual aspect of usage

"Ryanair-like" data test



(1333333333 /xxxx 1333<mark>5</mark>3333 /xxxx)



Variable Transformations ... simply & shortly

Variance X_i, Y Kurtosis, Skewness UNI Pearson correlation 5th /95th percentile

