

# FJFI & PRAXE

JAKUB STECH, PETRA KOSTAKOVA

7.3.2018

# O CEM TO BUDE

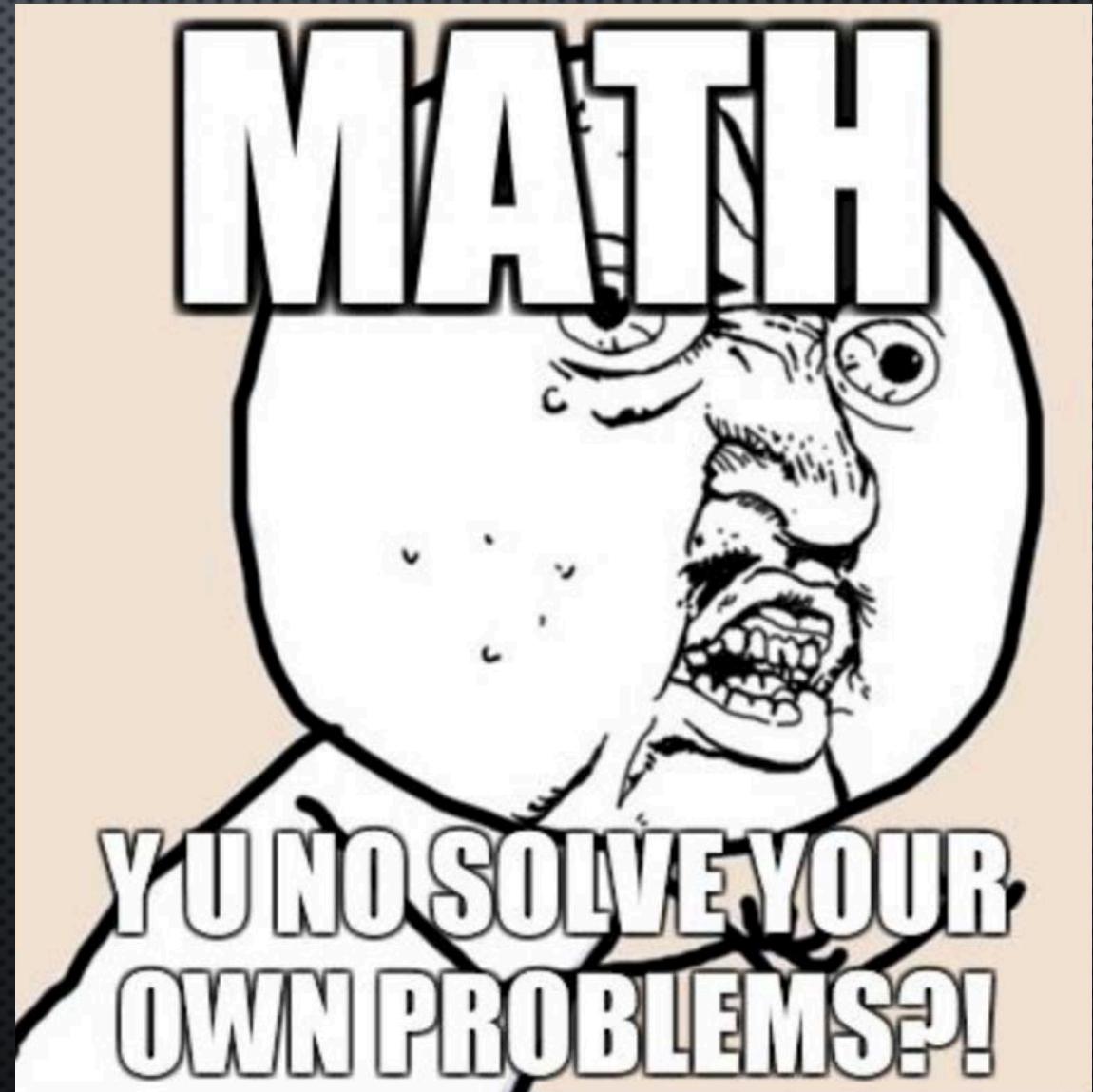
JAKUB

- PREDSTAVENI
- PRACE DATA ANALYTIKA+
- USE CASES

PETRA

- RANDOM FORESTS
- USE CASE

**VAROVANI**  
(PRED MOJI CASTI)



# KDO JSME

## JAKUB STECH

- ING NA FJFI AMSM 2016
- PHD 2016 - ???
- AS UTIA
- DATASENTICS 2016

# KDO JSME

## PETRA KOSTAKOVA

- ING NA FJFI MM 2007
- RIP PHD
- DATASENTICS 2016

**Petra Kocábová**

rok: 2007

zaměření: Matematické modelování

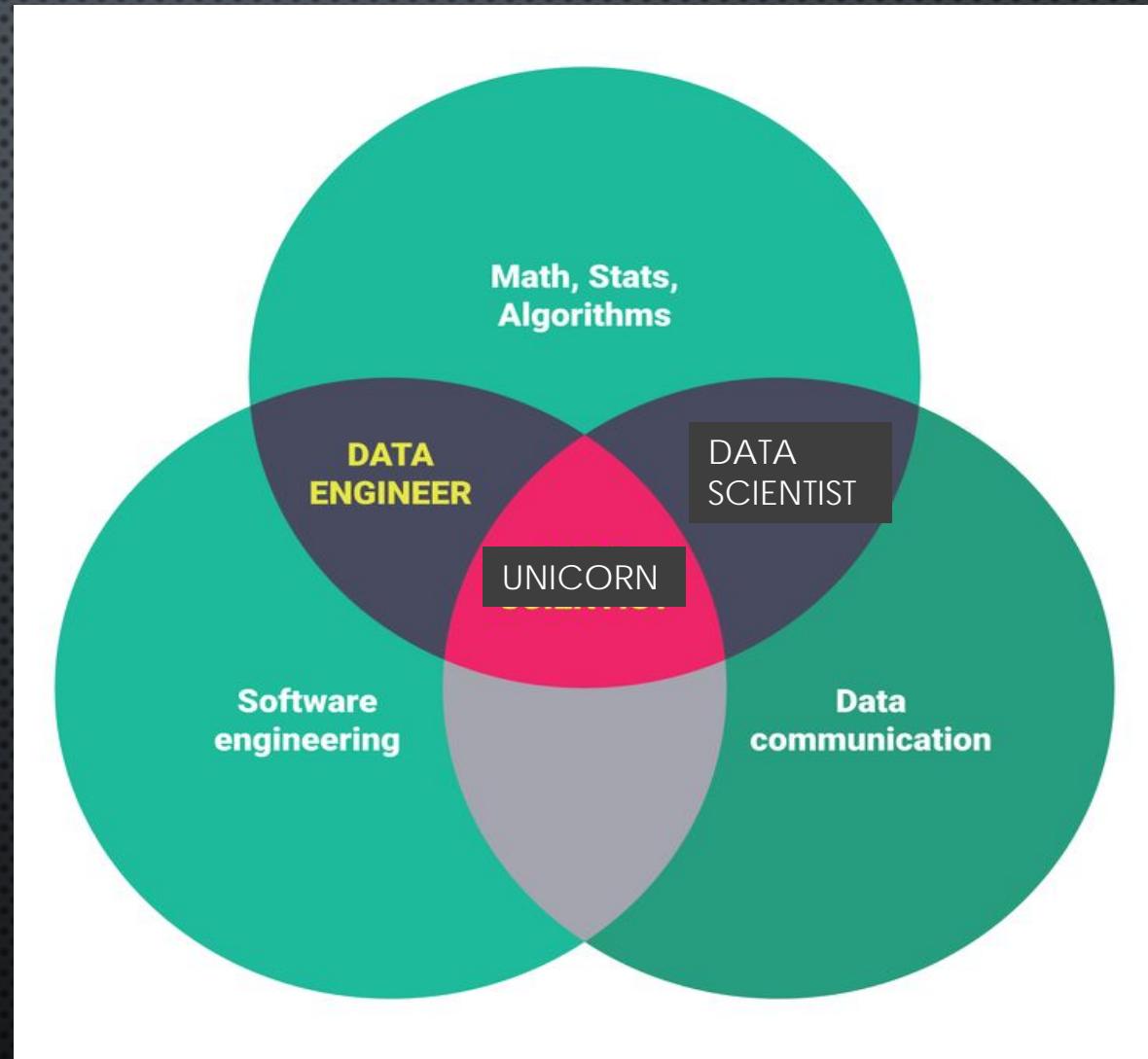
práce: Kvantová mechanika na násobně souvislých  
varietách Quantum mechanics on multiply  
connected manifolds

školitel: Prof. Ing P. Šťovíček, DrSc.

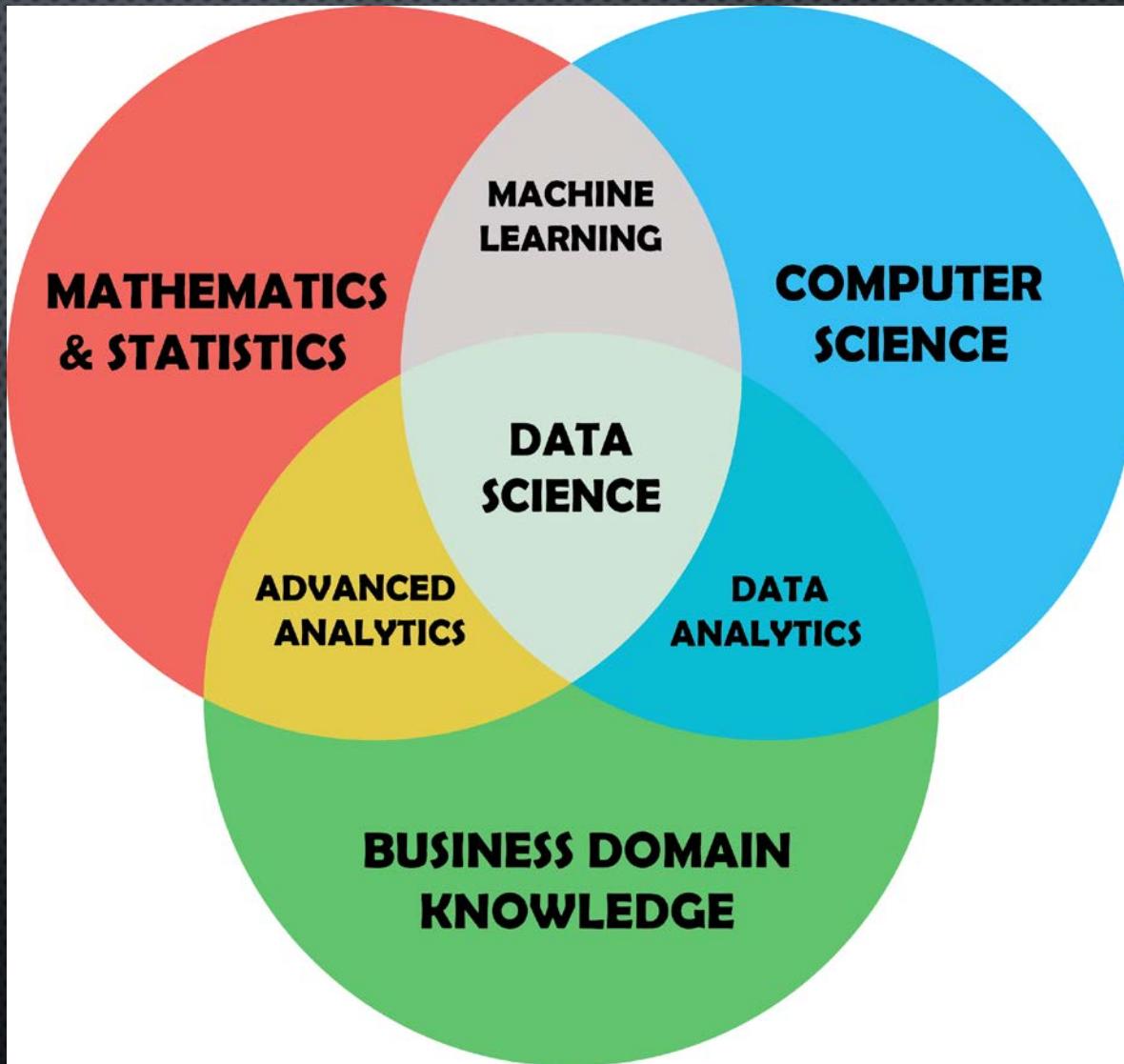
# **DATASENTICS**

- **MALY AGILNI BUTIK NA DATOVOU ANALYZU, ML, TECH (50% FJFI)**
- **DATA ENGINEERS vs DATA SCIENTISTS**
- **BIG DATA, LSTM 120L DEEP**
- **INSURANCE COMPANIES, ECOMMERCE, FMCG, HR, ...**

# TEAM

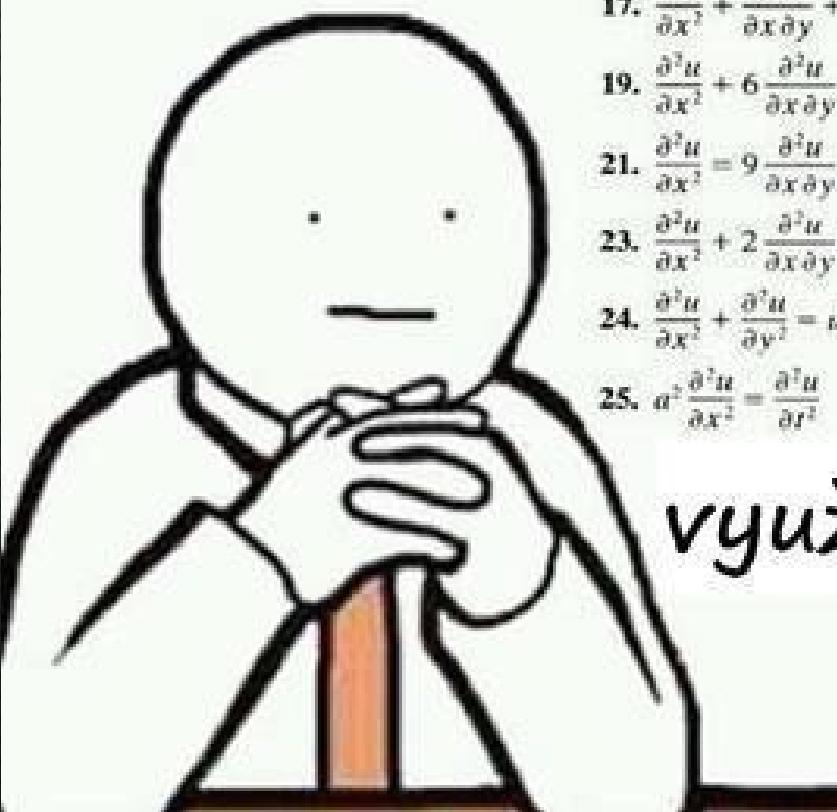


# DATA SCIENTIST



...FJFI?

Stále čakám na ten deň,  
kedy konečne toto



$$17. \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial x \partial y} + \frac{\partial^2 u}{\partial y^2} = 0$$

$$19. \frac{\partial^2 u}{\partial x^2} + 6 \frac{\partial^2 u}{\partial x \partial y} + 9 \frac{\partial^2 u}{\partial y^2} = 0$$

$$21. \frac{\partial^2 u}{\partial x^2} = 9 \frac{\partial^2 u}{\partial x \partial y}$$

$$23. \frac{\partial^2 u}{\partial x^2} + 2 \frac{\partial^2 u}{\partial x \partial y} + \frac{\partial^2 u}{\partial y^2} + \frac{\partial u}{\partial x} - 6 \frac{\partial u}{\partial y} = 0$$

$$24. \frac{\partial^2 u}{\partial x^2} + \frac{\partial^2 u}{\partial y^2} = u$$

$$25. a^2 \frac{\partial^2 u}{\partial x^2} = \frac{\partial^2 u}{\partial t^2}$$

$$18. 3 \frac{\partial^2 u}{\partial x^2} + 5 \frac{\partial^2 u}{\partial x \partial y} + \frac{\partial^2 u}{\partial y^2} = 0$$

$$20. \frac{\partial^2 u}{\partial x^2} - \frac{\partial^2 u}{\partial x \partial y} - 3 \frac{\partial^2 u}{\partial y^2} = 0$$

$$22. \frac{\partial^2 u}{\partial x \partial y} - \frac{\partial^2 u}{\partial y^2} + 2 \frac{\partial u}{\partial x} = 0$$

$$26. k \frac{\partial^2 u}{\partial x^2} = \frac{\partial u}{\partial t}, \quad k > 0$$

využijem v reálnom  
živote...

**CLV**  
**ARPU**  
**AOV**  
**CPC**  
**CPA**  
**YOY**  
**RFM**  
**RMSE**  
**ROI...**



# **AMSM**

- ROZ (CLUSTERING, PCA)
- TIN, ZLIM (GLM)
- SKE, MEX
- DATS (SQL)
- DRO
- X NEUR, SSI, PRAKТИКЕ UKAZKY

**PLEASE**

**TELL ME MORE**

memegenerator.net

# **CESTA PROJEKTU**

- 1 BUSINESS ANALYZA
- 2 SEZNANI A LOAD DAT
- 3 VIZUALIZACE
- 4 EXPLORACE
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE

# **CESTA PROJEKTU – BUSINESS ARCHITECT**

- 1 BUSINESS ANALYZA – POCHOPENI POTREB ZAKAZNIKA, 80:20, PoC
- 2 SEHNANI A LOAD DAT
- 3 VIZUALIZACE
- 4 EXPLORACE
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE

# CESTA PROJEKTU – DATA ENGINEER

- 1 BUSINESS ANALYZA
- 2 SEHNANI A LOAD DAT (DB, KBC, CSV)
- 3 VIZUALIZACE
- 4 EXPLORACE
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE - PERSONALIZACE DLE VÝSTUPU, POTREB, SYSTEM, PENEZ



## Overview

### DataSentics



Overview

Extractors

Transformations

Writers

Orchestrations

Storage

Jobs

Applications

**This project will expire in 18 days**

Please [contact support](#) for project plan upgrade.



#### Storage

**113.00 MB**

3,633,705 Rows

#### Access

**10** Users

3 API Tokens

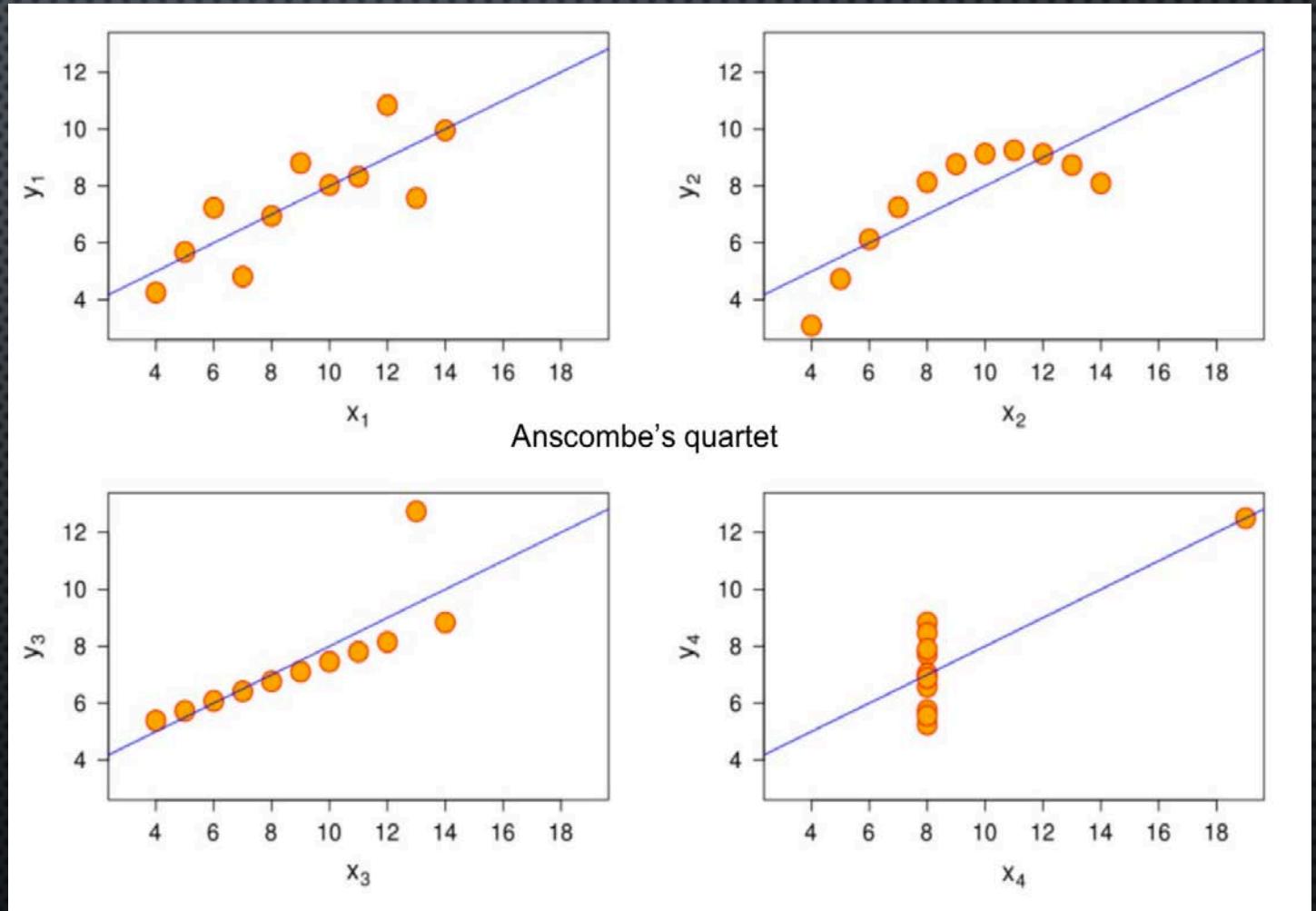
# **CESTA PROJEKTU – DATA SCIENTIST**

- 1 BUSINESS ANALYZA
- 2 SEHNANI A LOAD DAT
- 3 VIZUALIZACE (QLIK, TABLEAU)
- 4 EXPLORACE (R, PYTHON)
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE

# 3 VIZUALIZACE

... EH, NESTACI CISLA?

# 3 VIZUALIZACE



## 3 VIZUALIZACE

All the summary statistics you'd think to compute are close to identical:

- The average  $x$  value is 9 for each dataset
- The average  $y$  value is 7.50 for each dataset
- The variance for  $x$  is 11 and the variance for  $y$  is 4.12
- The correlation between  $x$  and  $y$  is 0.816 for each dataset
- A linear regression (line of best fit) for each dataset follows the equation  $y = 0.5x + 3$

# 3 VIZUALIZACE

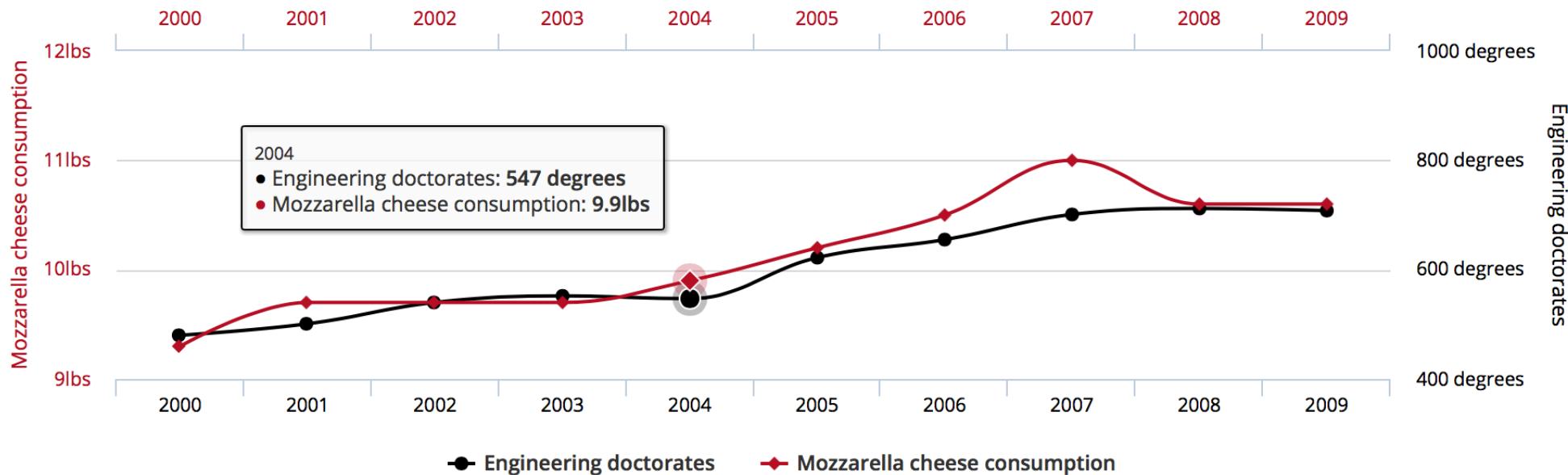
... POZOR NA BUSINESS KONTEXT

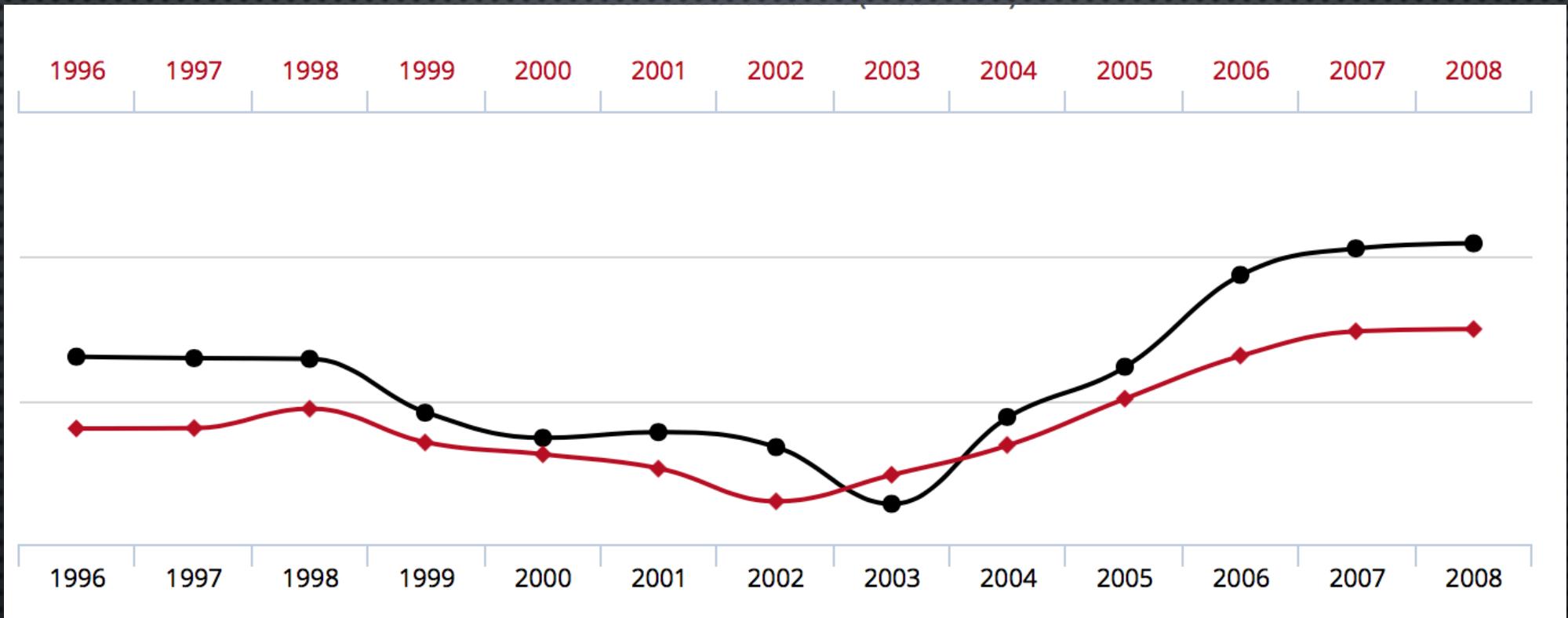


# Per capita consumption of mozzarella cheese correlates with Civil engineering doctorates awarded

☰

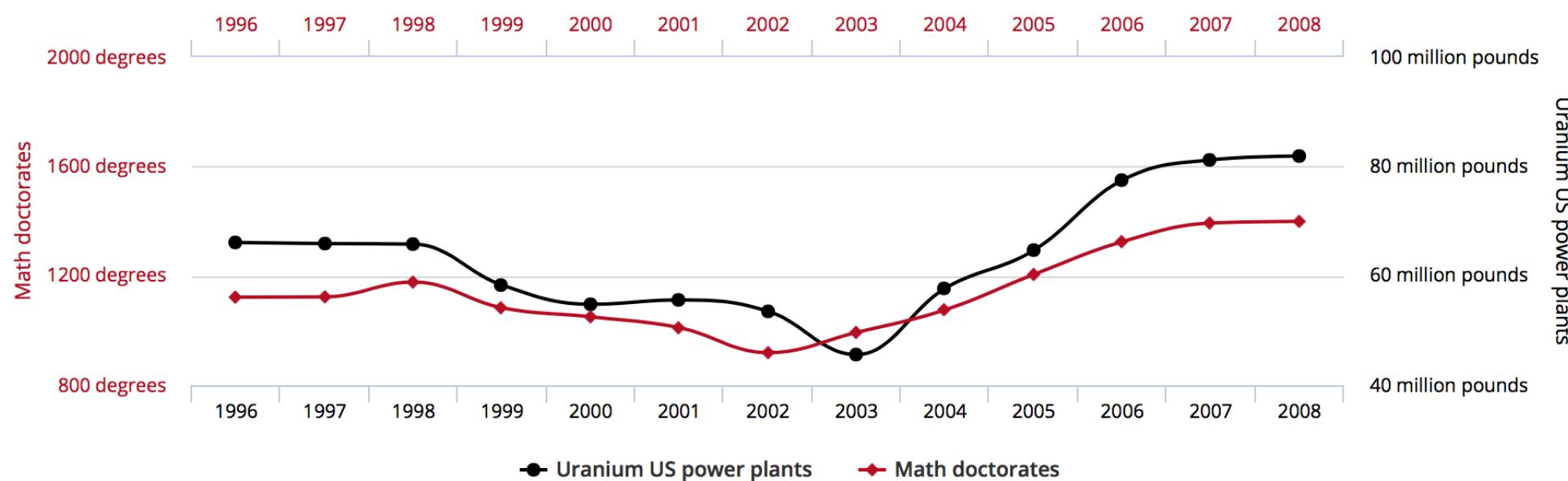
Correlation: 95.86% ( $r=0.958648$ )





# Math doctorates awarded correlates with Uranium stored at US nuclear power plants

Correlation: 95.23% ( $r=0.952257$ )



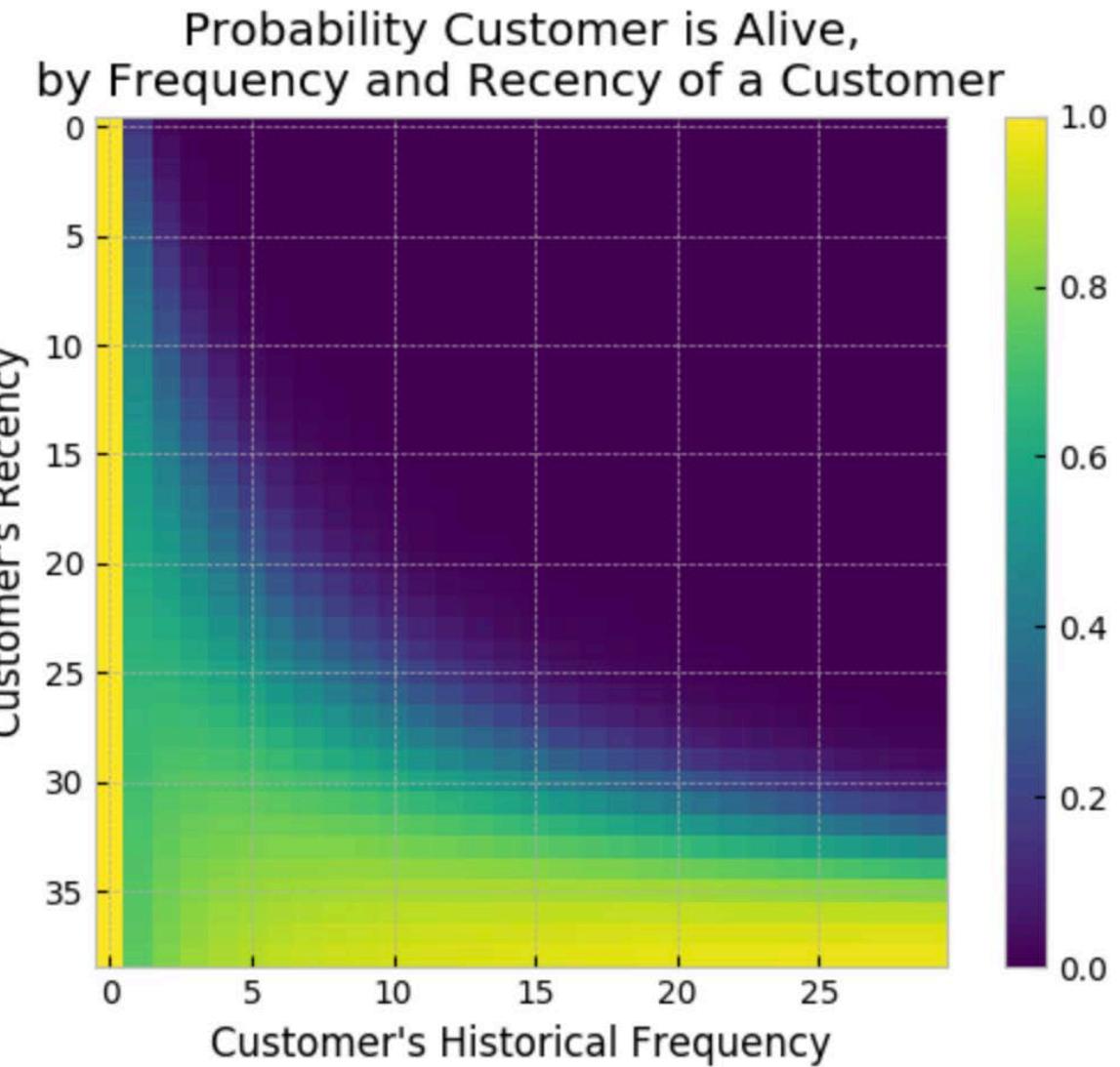
## 4 EXPLORACE

- CUSTOMER JOURNEY (WEB, NAKUPY)
- SKU, REVENUE MEZI N A N+1 NAKUPEM
- ACTIVE USER BASE

# **5 DATOVÁ ANALYZA, MACHINE LEARNING**

- PROBABILISTIC: CLV, CHURN RATES
- REGRESSION: (G)LM, (BAGGED) RANDOM FORESTS
- CLUSTERING: K-MEANS, SVM, RF
- NLP: TF-IDF, WORD2VEC

# CHURN RATES



# **CESTA PROJEKTU – DATA ENGINEER**

- 1 BUSINESS ANALYZA
- 2 SEHNANI A LOAD DAT (DB, KBC, CSV)
- 3 VIZUALIZACE
- 4 EXPLORACE
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE - PERSONALIZACE DLE VÝSTUPU, POTREB, SYSTEM, PENEZ

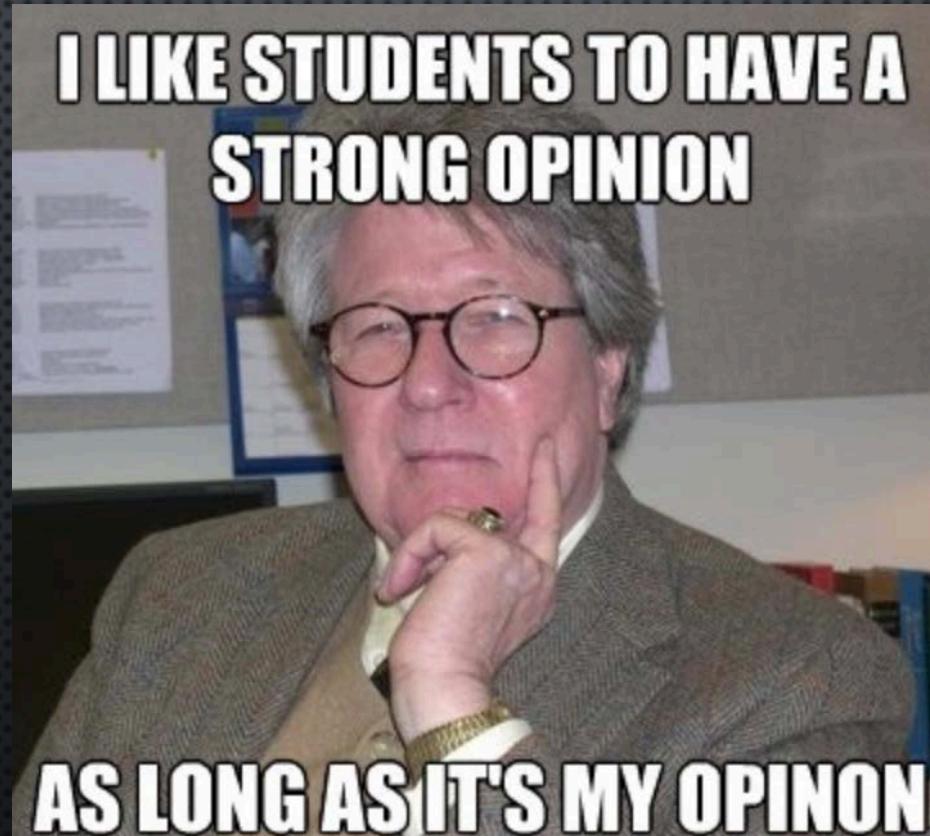
**DAFUQ YOU'RE TALKING ABOUT**



**SHOW ME ██████████ EXAMPLES!**

memegenerator.net

# 1 DYNAMICKY DIGITALNI MARKETING



# 1 DYNAMICKY DIGITALNI MARKETING

VYUZITI VICE DATOVYCH ZDROJU  
DORUCIT PRESONALIZOVANY OBSAH

- NIZSI CHURN RATE
- VYSSI CTR, NIZSI CPA => VYSSI REVENUE

AUTOMATIZACE

# CO NEDELAME



# CO DELAME



The screenshot shows the Keboola platform interface. On the left is a sidebar with the following navigation items:

- Overview
- Extractors
- Transformations
- Writers
- Orchestrations
- Storage
- Jobs

At the top center is the title "Transformations" with a refresh icon. To the right are two buttons: "Sandbox" and "+ New Bucket". The main area displays a list of transformations:

Search		
Data_studio_customer support	No description	
DataSentics analytika	analyticke treansformace od DataSentics	
1 clv	R	clv
1 current_clv	R	actual monthly consuption for FB audiences
1 mkt_30	R	1st and 2nd buy in last 30 days
1 order_table	R	No description

# JAK TECOU DATA



PrestaShop



Facebook  
Custom Audience  
And What It Means For Your Brand

happy  
marketer



# CO DELAME



The screenshot shows the Keboola platform interface. On the left is a sidebar with the following navigation items:

- Overview
- Extractors
- Transformations
- Writers
- Orchestrations
- Storage
- Jobs

At the top center is the title "Transformations" with a refresh icon. To the right are two buttons: "Sandbox" and "+ New Bucket".  
The main area displays a list of transformations:

Data_studio_customer support	No description			
DataSentics analytika	analyticke treansformace od DataSentics			
1 clv	R	clv		
1 current_clv	R	actual monthly consupption for FB audiences		
1 mkt_30	R	1st and 2nd buy in last 30 days		
1 order_table	R	No description		

# CO DELAME

Okruhy uživatelů						
	Název	Typ	Velikost	Dostupnost	Datum	
<input type="checkbox"/>	nad10	Vlastní okruh uživatelů Seznam zákazníků	Méně než 1000	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 14:09	
<input type="checkbox"/>	do10	Vlastní okruh uživatelů Seznam zákazníků	Méně než 1000	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 14:09	
<input type="checkbox"/>	do4	Vlastní okruh uživatelů Seznam zákazníků	Méně než 1000	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 14:08	
<input type="checkbox"/>	do2	Vlastní okruh uživatelů Seznam zákazníků	Méně než 1000	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 14:04	
<input type="checkbox"/>	do1	Vlastní okruh uživatelů Seznam zákazníků	Méně než 1000	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 14:03	
<input type="checkbox"/>	current_clv	Vlastní okruh uživatelů Seznam zákazníků	4 400	<span>Připraveno</span> Poslední aktualizace: 19.2.2018	19.2.2018 11:05	



rohlik.cz

31 říjen v 9:38 ·

...

Oblíbené pesto Barilla v akci za 69,90 Kč! Skvělé jídlo s ním vykouzlíte za pár minut 🍅

CENOVKA  
DNE



**69,90**  
~~99,90~~  
**-30%**

Objevte naší cenovku dne!

Nakupujte zde >>>

ROHLIK.CZ



Super



Komentář



Sdílet



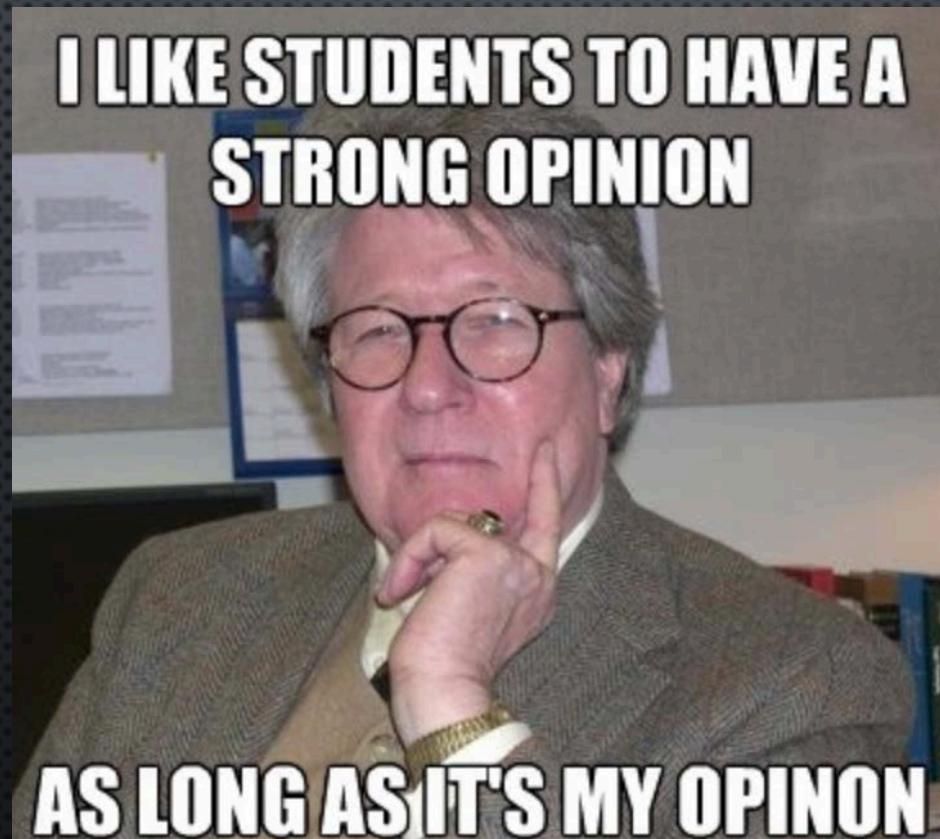
▼



Vy a 8 dalších

Hlavní komentáře ▼

2 TAXI



## **2 SENSITIVITA ZAKAZNIKU**

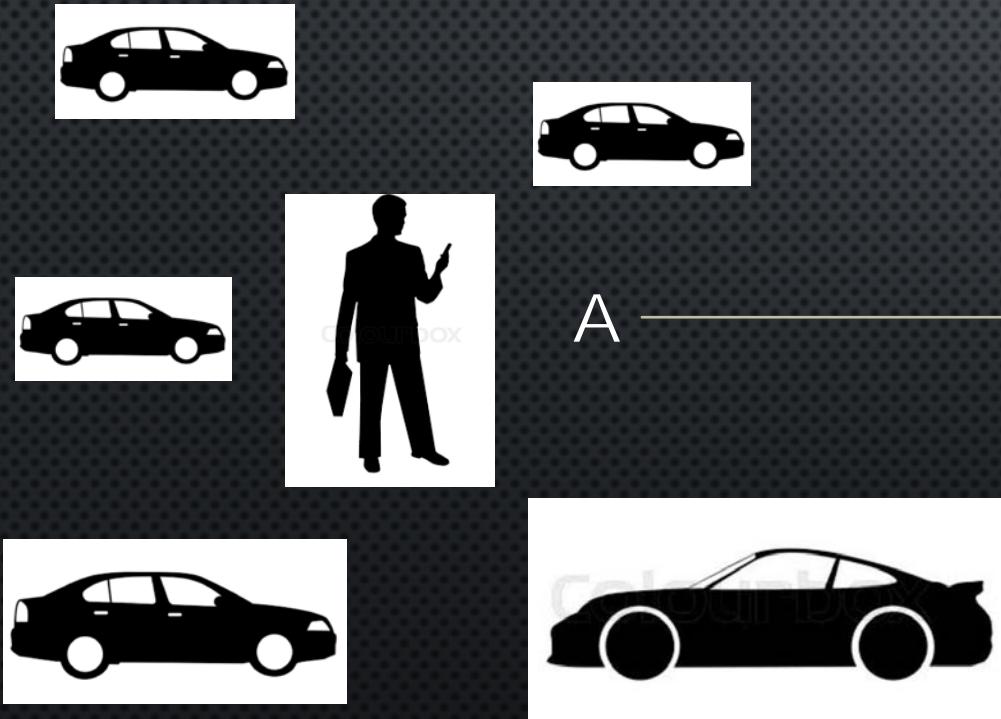
### **EXTRAKCE PREFERENCI ZAKAZNIKU**

- NIZSI CHURN RATE**
- ZVYSENI OBSAZENOSTI**
- ZVYSENI REVENUE**

**CC POJISTOVNA, TAXISLUZBA**

**GLM, (RF)**

# TAXI PROBLEM



Cena  
Prijezd  
Delka jizdy  
Komfort auta  
Hodnocení řidiče

# TAXI MODEL

- LOGISTICKA REGRESE
- $P(JIZDA = ACCEPTED) = A * CENA + B * PRIJEZD + C * DELKA JIZDY + D * KOMFORT AUTA$
- BUSSINESS JIZDY = KLADNE “A”
- NEBUSSINESS = ZAPORNE “B” KAM AZ ZDRAZIT VIZ THRESH

OTAZKY...?

# TREES, RANDOM FOREST, GRADIENT BOOSTED TREES

APLIKACE V POJIŠŤOVNÁCH

# BUSINESS MOTIVACE

- JAKÉ OTÁZKY ŘEŠÍ POJIŠŤOVNY A JAKÉ DATA MAJÍ?

# BUSINESS MOTIVACE

- JAKÉ OTÁZKY ŘEŠÍ POJIŠTOVNY A JAKÁ DATA MAJÍ?

## DATA

- DATA Klientská – nastavení smlouvy, počet smluv, klientská historie, demografické údaje, interakce se zákazníkem – CRM databáze o cca 400 sloupcích
- Obohacení o statistiky, průzkumy, banko produkty

# BUSINESS MOTIVACE

- JAKÉ OTÁZKY ŘEŠÍ POJIŠŤOVNY A JAKÁ DATA MAJÍ?

## OTÁZKY

- REPORTING 😊
- CROSS SELL, UP SELL, SPRÁVNÝ ČAS NA OSLOVENÍ Klienta
- OPTIMÁLNÍ NASTAVENÍ PRODUKTU
- PREDIKCE RIZIKOVÝCH SKUPIN ATD.

# MACHINE LEARNING

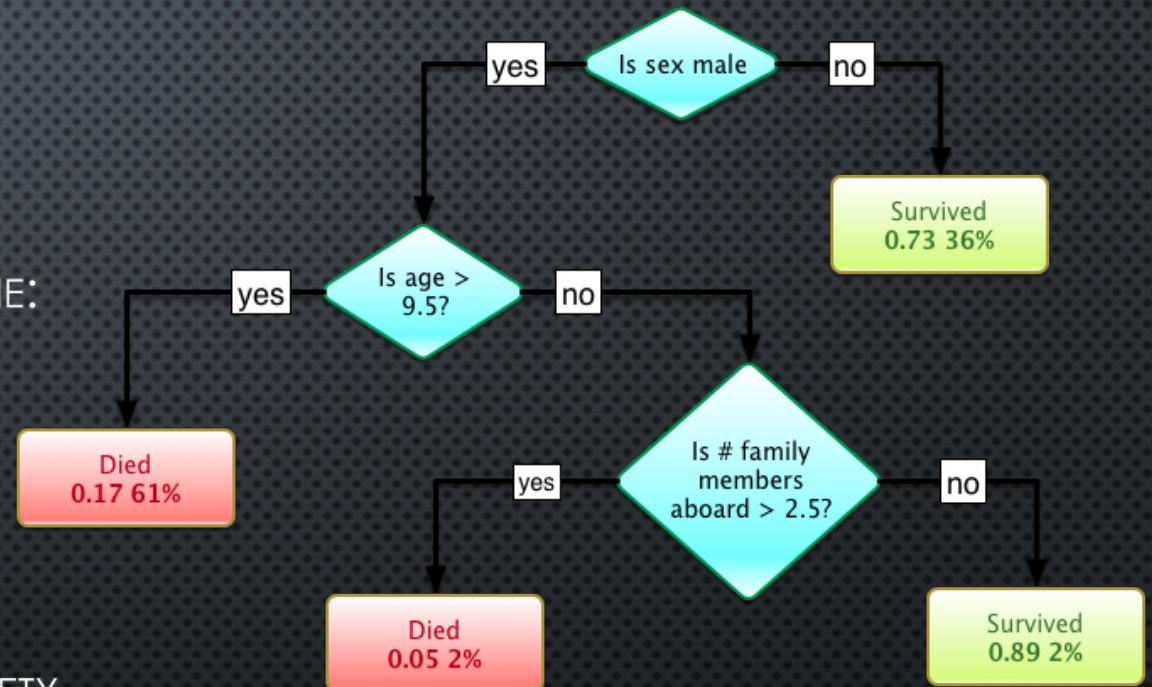
- OBECNÁ KONCEPCE:
- $F(x) = y$
- BIAS – VARIANCE TRADEOFF (UNDERFITTING, OVERFITTING)

# STROMY

IKUZNE ALGORITMUS KRU KUZDELENI DU LISTU –  
NEJBĚŽNĚJŠÍ MAXIMÁLNÍ NÁRŮST INFORMAČNÍ ENTROPIE:

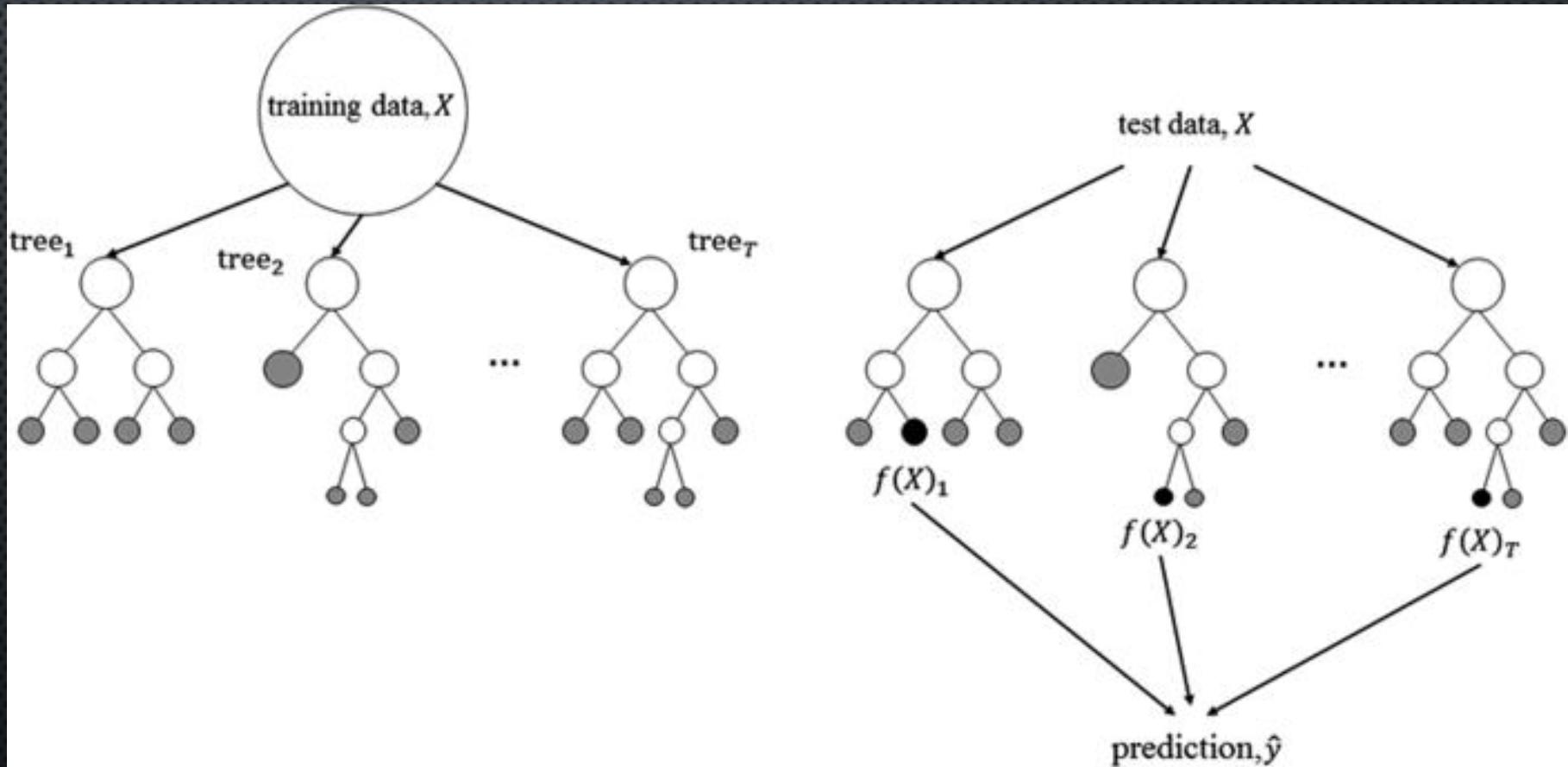
$$H(T) = I_E(P_1, P_2, \dots, P_j) = \sum_{i=1}^j P_i \log_2 p_i$$

- VÝHODY: LZE DO NĚJ VLOŽIT KATEGORICKÉ I NUMERICKÉ HODNOTY, SNADNÁ INTERPRETOVATELNOST, FUNGUJE PRO VELKÉ DATASETY, DATA NENÍ TŘEBA MOC PŘEDUPRAVOVAT
- NEVÝHODY: NON-ROBUST, OVERTFITTING, NEMUSÍ



# RANDOM FOREST

NOVÉ STROMY VZNIKAJÍ  
SAMPOVÁNÍM  
DATASETU/PROMĚNNÝCH



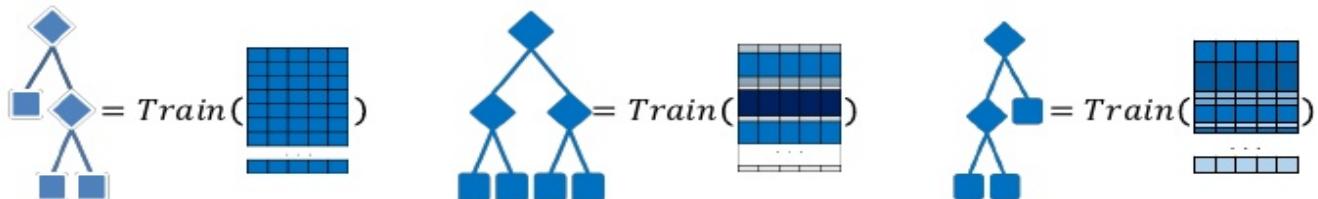
# GRADIENT BOOSTED TREES

- JEDEN Z NEJLEPŠÍCH ALGORITMŮ PRO KLASIFIKACI – KAGGLE
- OPEN SOURCE IMPLEMENTACE PRO R A PYTHON (A DALŠÍ)
- PARALELNÍ -> RYCHLÉ
- UMÍ POSTIHNUŤ I KOMPLEXNÍ PROBLÉMY
- DOCELA SLUŠNĚ INTERPRETOVATELNÉ

## Boosting: Iterative Tree Construction

*"Best off-the-shelf classifier in the world" – Breiman*

- Reweight examples for each subsequent tree to focus on **errors**



- Numerically: gradient descent in function space

– Each subsequent tree approximates a step in  $-\frac{\partial L}{\partial f}$  direction

– Recompute **target labels**

$$y^{(m)} = - \left[ \frac{\partial L(y, f(x))}{\partial f(x)} \right]_{f(x)=f^{(m-1)}(x)}$$

– Logistic loss:  $L(y, f(x)) = \log(1 + \exp(-yf(x)))$

$$y^{(m)} = \frac{y}{1 + \exp(yf(x))}$$

– Squared loss:  $L(y, f(x)) = \frac{1}{2}(y - f(x))^2$

$$y^{(m)} = y - f(x)$$

# VYHODNOCENÍ MODELŮ

- ROC KŘIVKA

# CO TO ZNAMENÁ V PRAXI

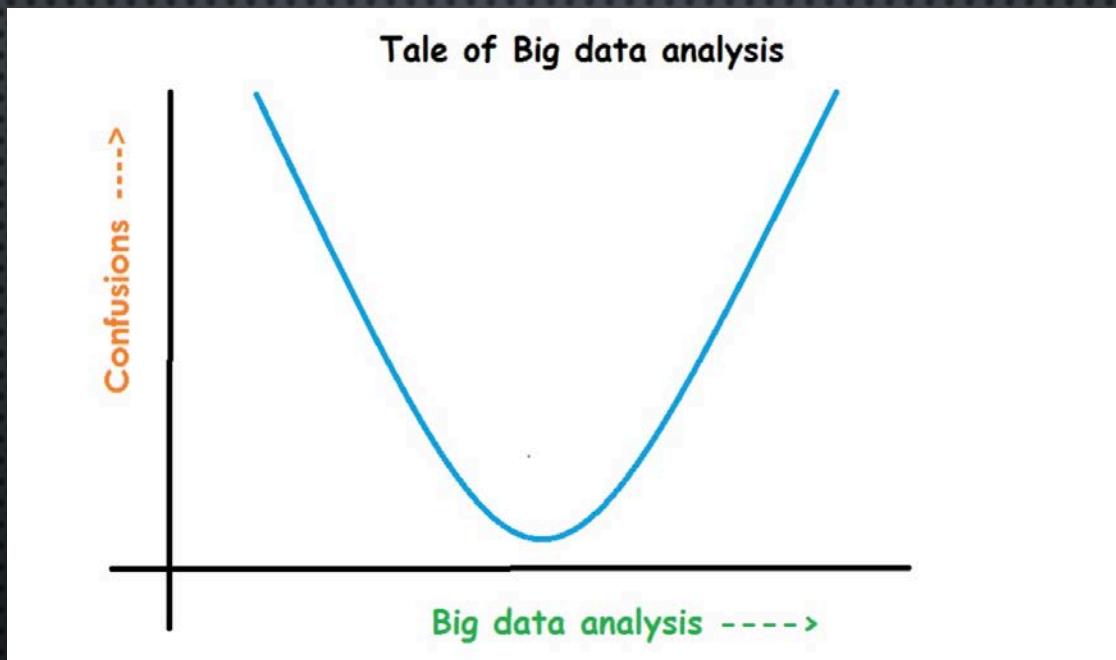
- MÍT K DISPOZICI DATA V ANALYTICKÉ DATABÁZI
- MÍT NAPOČÍTANOU L2 VRSTVU CRM
- 1000 – 2000 ŘÁDKŮ V RKU (TAKOVÝ DROBNÝ ZÁKLAD)
- DŮRAZ NA INTERPRETACI
- AŽ SEDMINÁSOBNĚ VĚTŠÍ ÚSPĚŠNOST PŘI OSLOVENÍ, NESPAMUJEME VŠECHNY

OTAZKY...?

# ZAJIMAVÉ ZDROJE

- DATA CAMP
- [HTTPS://WWW.R-BLOGGERS.COM/](https://www.r-bloggers.com/)
- MEETUP PRAHA
- DATA ANALYSTS, DATA ENGINEERS & DATA SCIENTISTS - CZECH&SLOVAK GROUP
- JAKUB.STECH@DATASENTICS.COM

# DEKUJEME ZA POZORNOST



# KDE JSME SKONCILI

- FJFI DOBRY PREDPOKLAD
- DRO
- DATA SCIENCE WORKFLOW

# **CESTA PROJEKTU – DATA SCIENTIST**

- 1 BUSINESS ANALYZA
- 2 SEHNANI A LOAD DAT
- 3 VIZUALIZACE (QLIK, TABLEAU)
- 4 EXPLORACE (R, PYTHON)
- 5 DATOVÁ ANALYZA, MACHINE LEARNING
- 6 PRODUKCIJNALIZACE



rohlik.cz

31 říjen v 9:38 ·

...

Oblíbené pesto Barilla v akci za 69,90 Kč! Skvělé jídlo s ním vykouzlíte za pár minut 🍅

CENOVKA  
DNE



**69,90**  
~~99,90~~  
**-30%**

Objevte naší cenovku dne!

Nakupujte zde >>>

ROHLIK.CZ



Super



Komentář



Sdílet



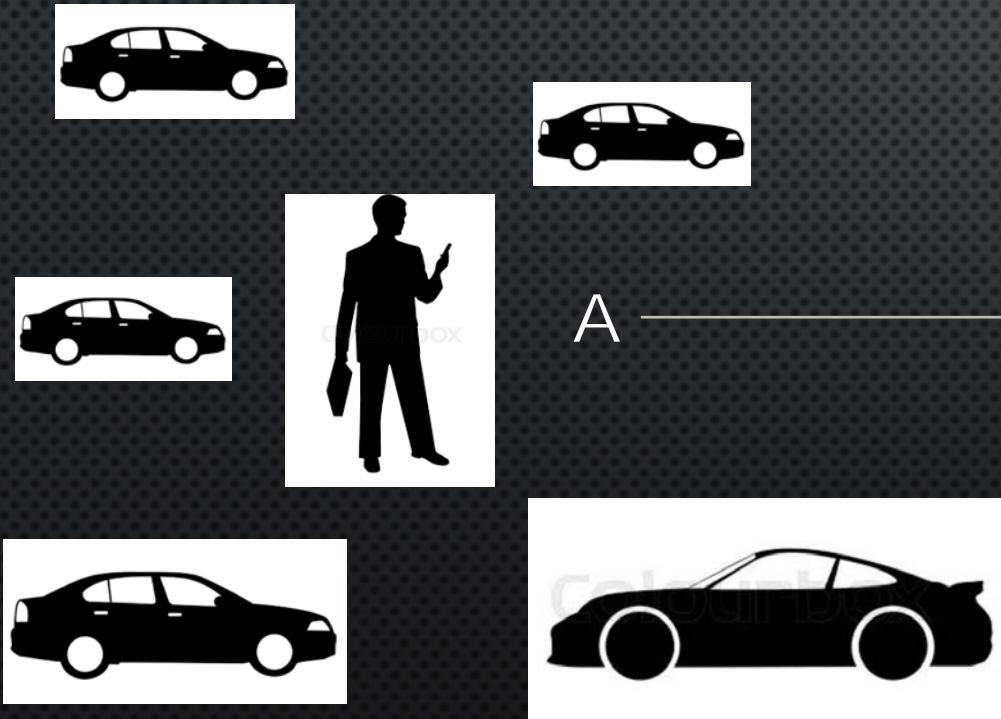
▼



Vy a 8 dalších

Hlavní komentáře ▼

# TAXI PROBLEM



Cena  
Prijezd  
Delka jizdy  
Komfort auta  
Hodnoceni ridice

# KDE JSME SKONCILI

...KOHO TO ZAUJALO?



# CO DNES?

- GLM v POJISTOVNE
- PROBABILISTIC CHURN RATE + CLV, DATASENTICS CLV v ECOMMERCE

# **1 CALLCENTRUM POJISTOVNA**

## **SITUACE**

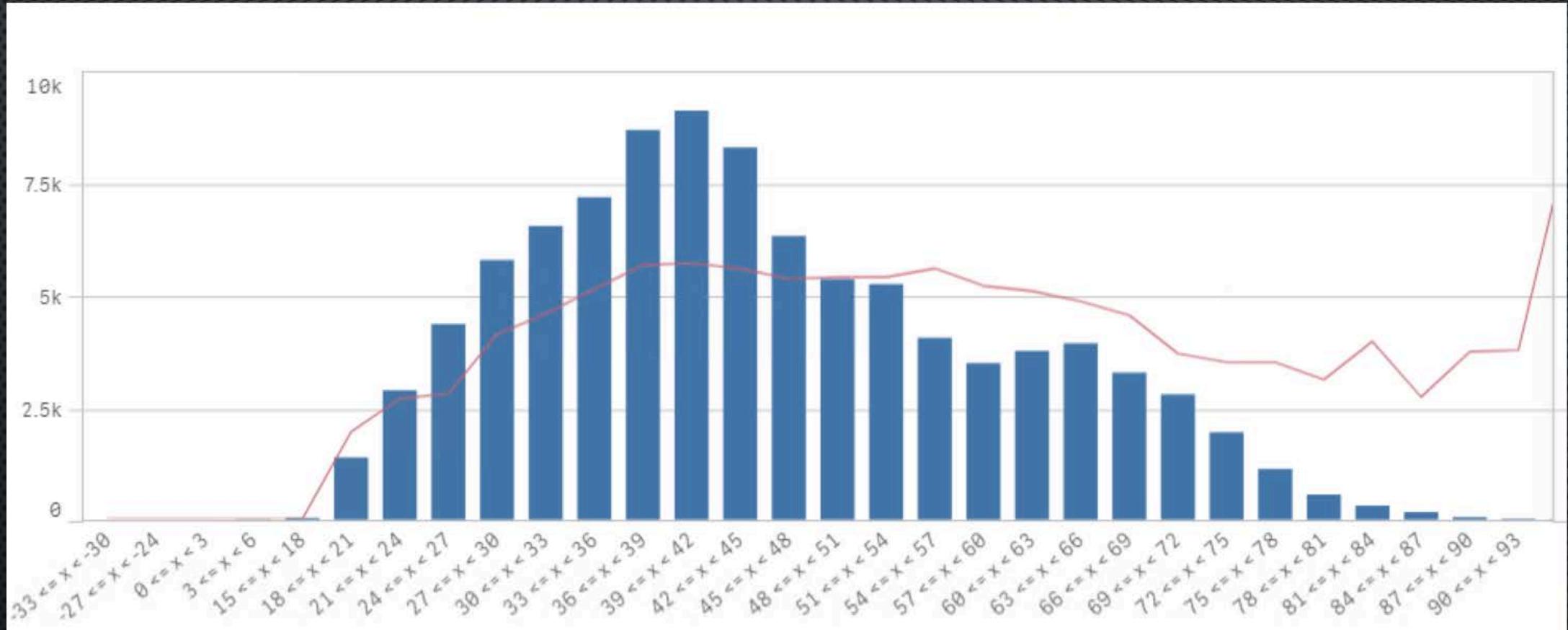
- KALKULACE NA WEBU -> PRVNI NABIDKA**
- KDYZ NESJEDNA -> CALLCENTRUM -> DRUHA NABIDKA**

# **1 CALLCENTRUM POJISTOVNA**

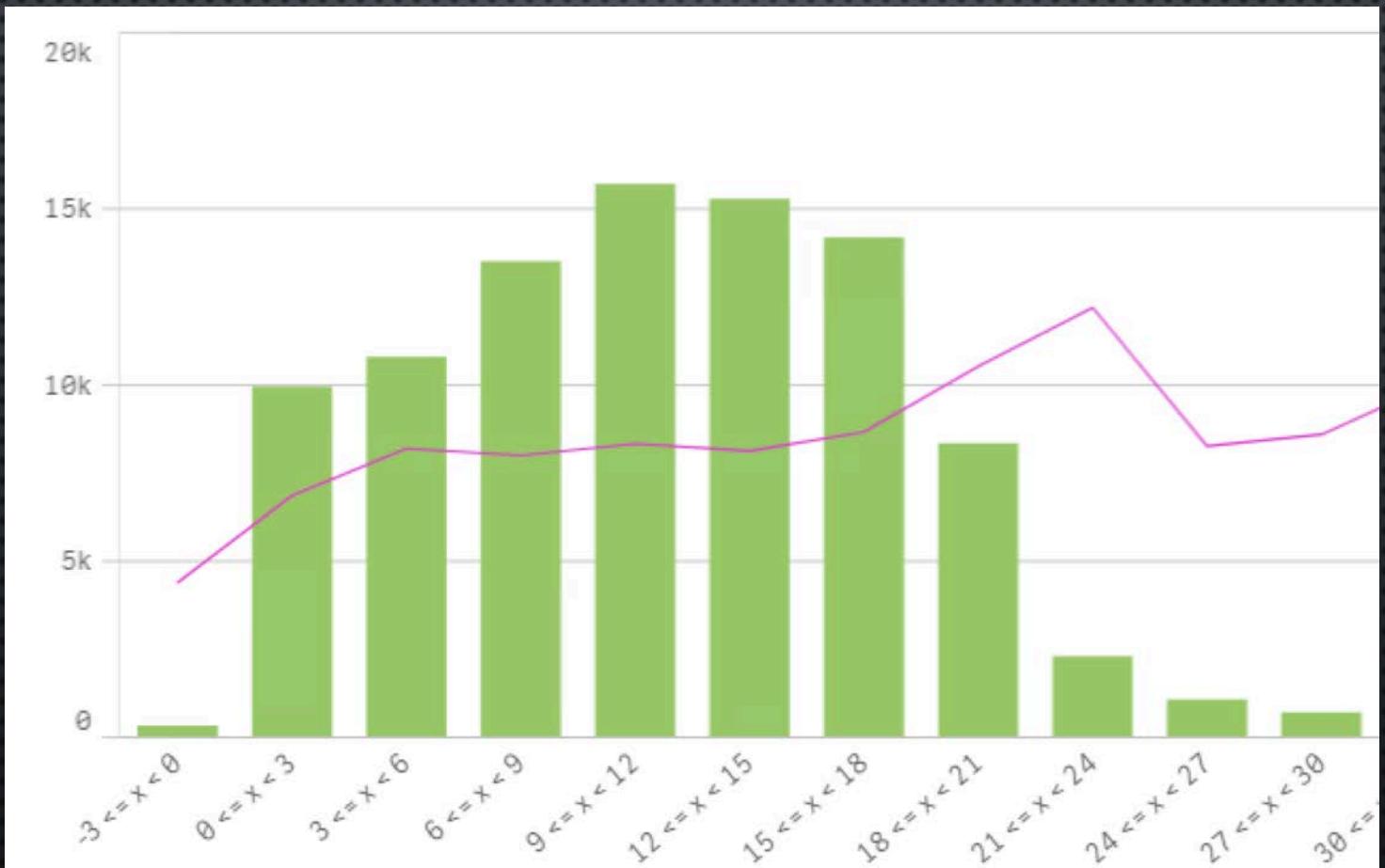
## **BUSINESS MOTIVACE**

- UPRAVA/PERSONIFIKACE NABIDKY NA WEBU/CALLCENTRU**
- PRIORITIZACE NAVOLAVANI**
- UPSELL HAV, PPC, ...**

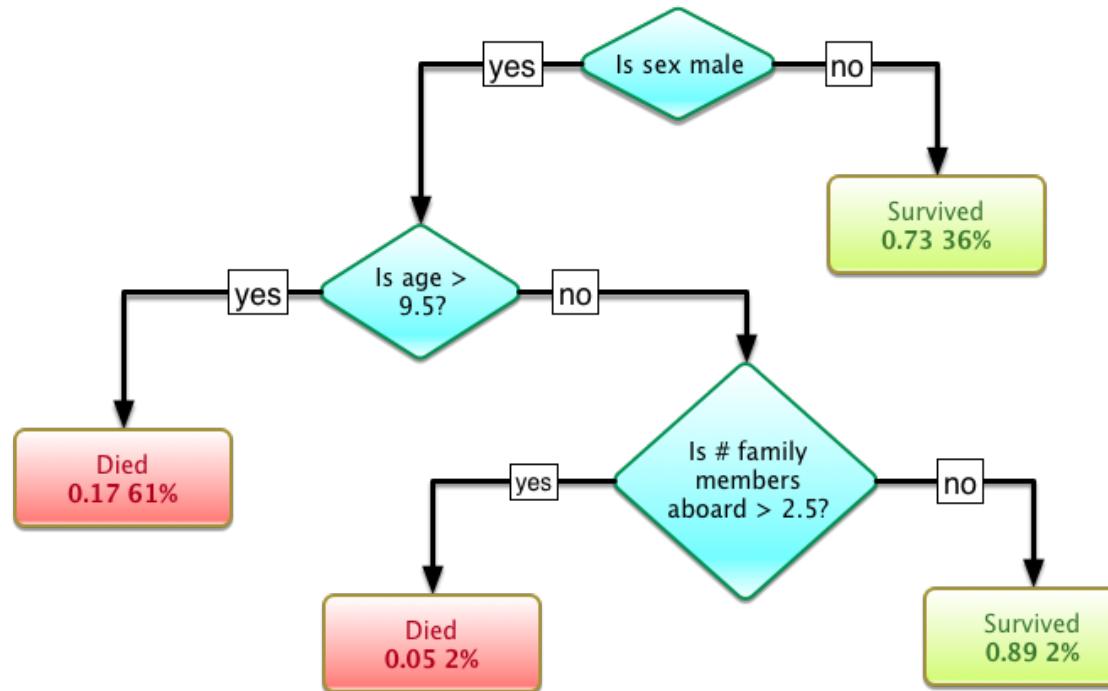
# VEK



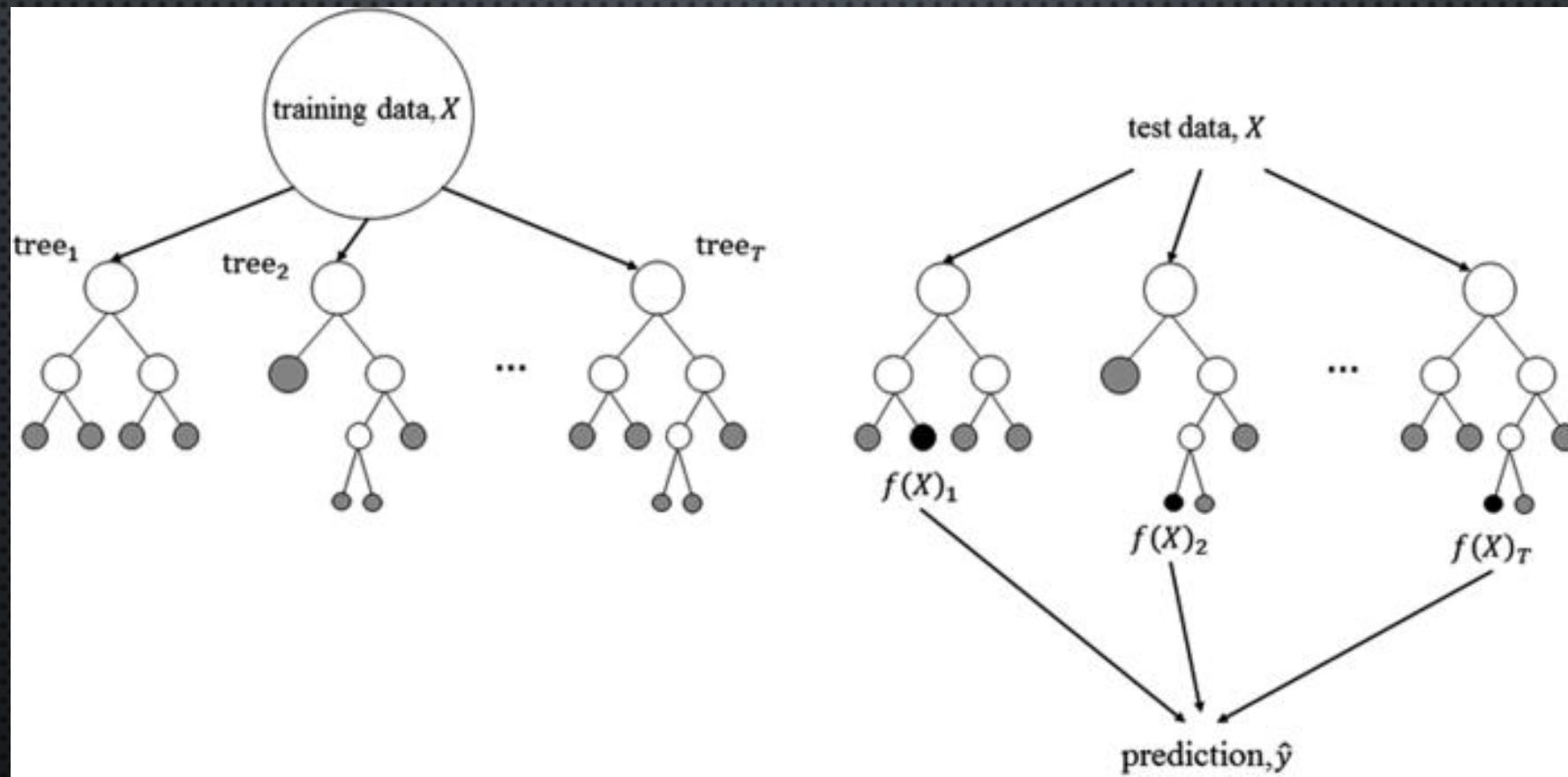
# STARÍ AUTA



# STROMY



# RANDOM FOREST

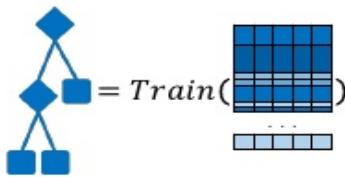
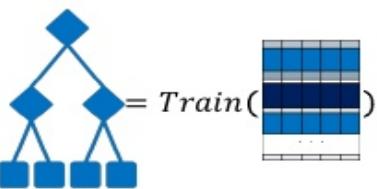
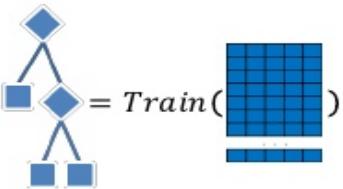


# GRADIENT BOOSTED TREES

## Boosting: Iterative Tree Construction

*“Best off-the-shelf classifier in the world” – Breiman*

- Reweight examples for each subsequent tree to focus on **errors**



- Numerically: gradient descent in function space

- Each subsequent tree approximates a step in  $-\frac{\partial L}{\partial f}$  direction

- Recompute **target labels**

$$y^{(m)} = - \left[ \frac{\partial L(y, f(x))}{\partial f(x)} \right]_{f(x) = f^{(m-1)}(x)}$$

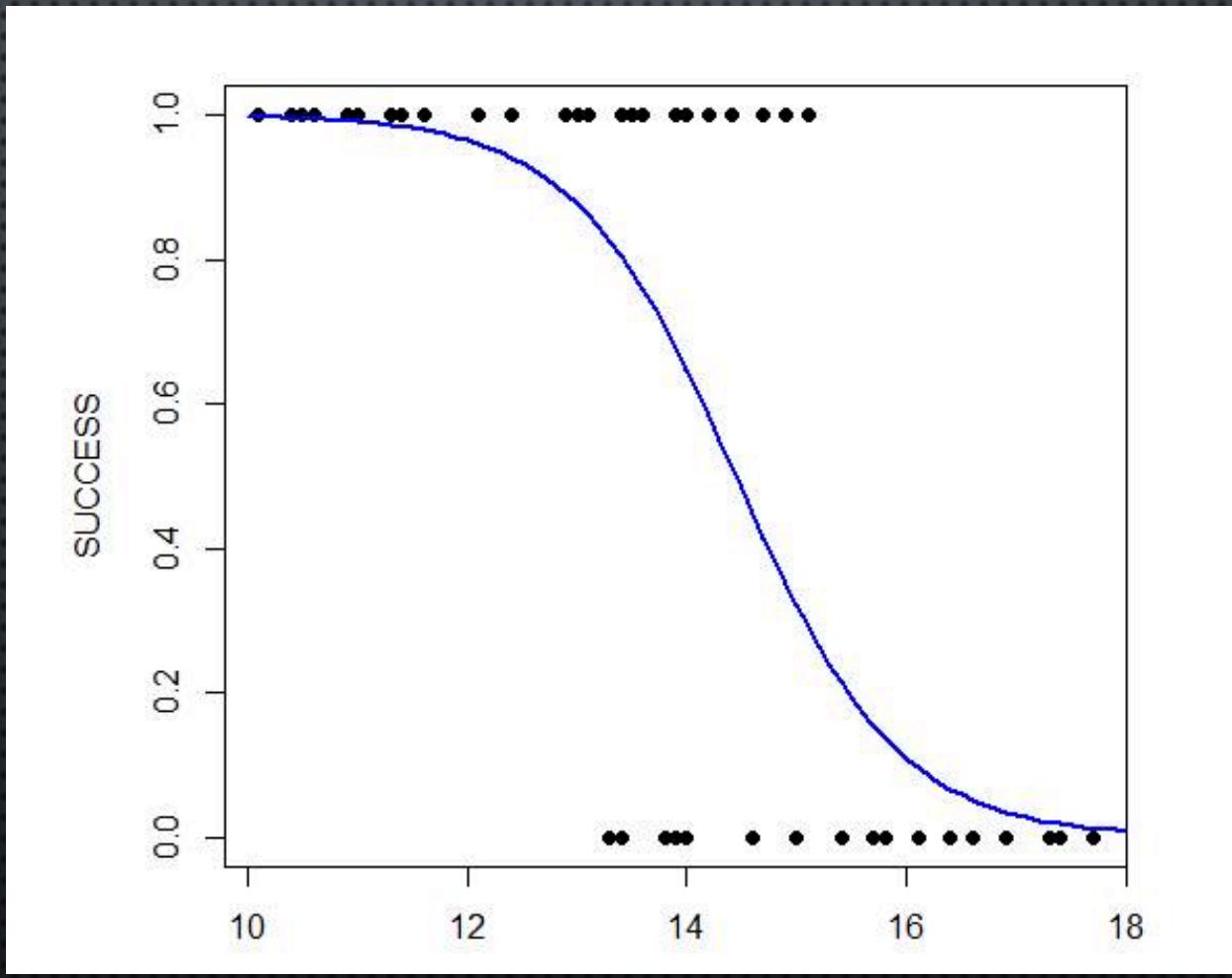
- Logistic loss:  $L(y, f(x)) = \log(1 + \exp(-yf(x)))$

$$y^{(m)} = \frac{y}{1 + \exp(yf(x))}$$

- Squared loss:  $L(y, f(x)) = \frac{1}{2}(y - f(x))^2$

$$y^{(m)} = y - f(x)$$

# GLM



# GLM

$$y = e^{(b_0 + b_1 * x)} / (1 + e^{(b_0 + b_1 * x)})$$

$$\text{logit}(p) = b_0 + b_1 X_1 + b_2 X_2 + b_3 X_3 + \dots + b_k X_k$$

$$\text{logit}(p) = \ln\left(\frac{p}{1-p}\right)$$

# GLM

- $P/(1-P) = ODDS$

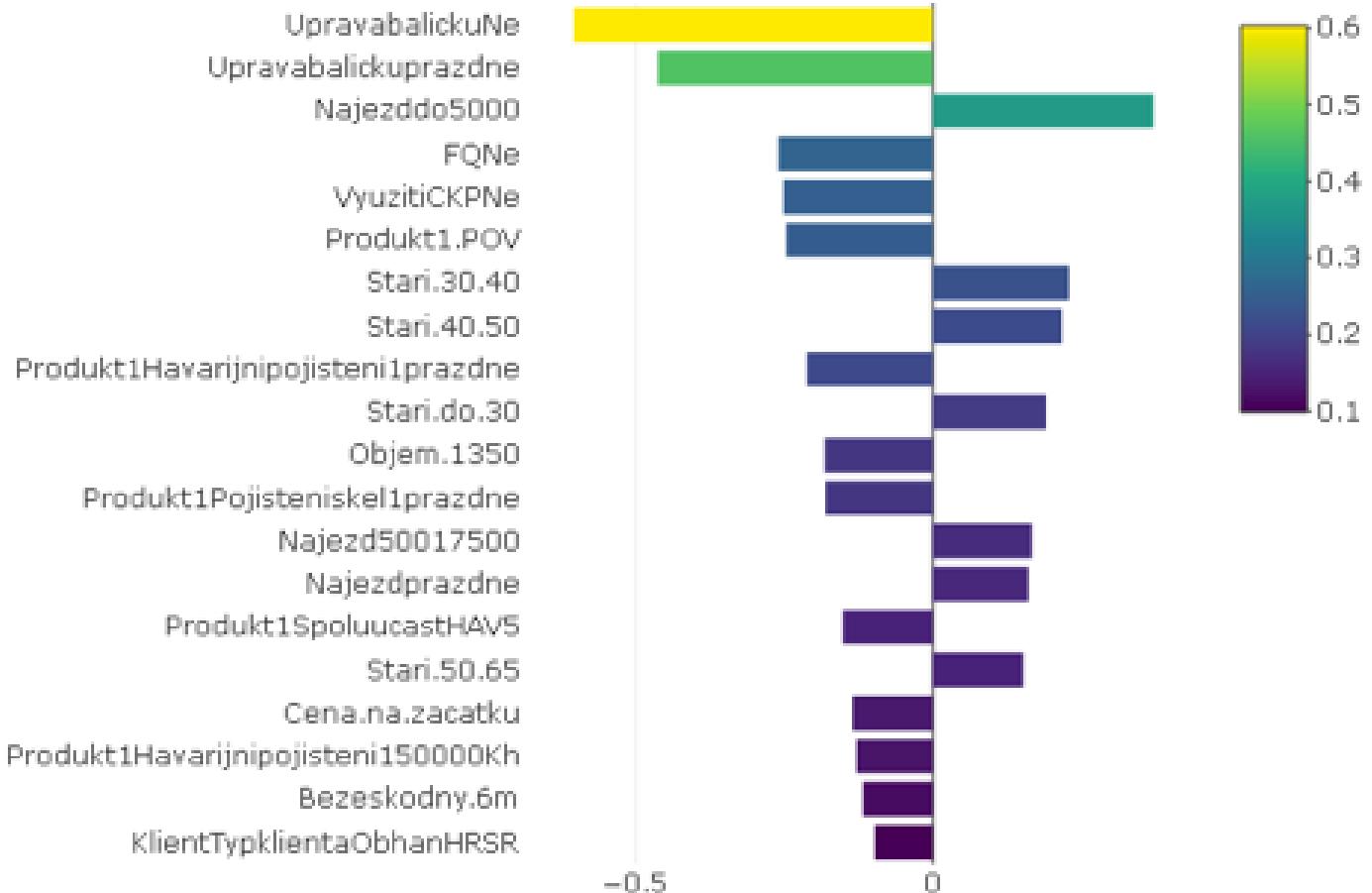
$$odds = e^{(b_0 + b_1 * x)}$$

# **1 CALLCENTRUM POJISTOVNA**

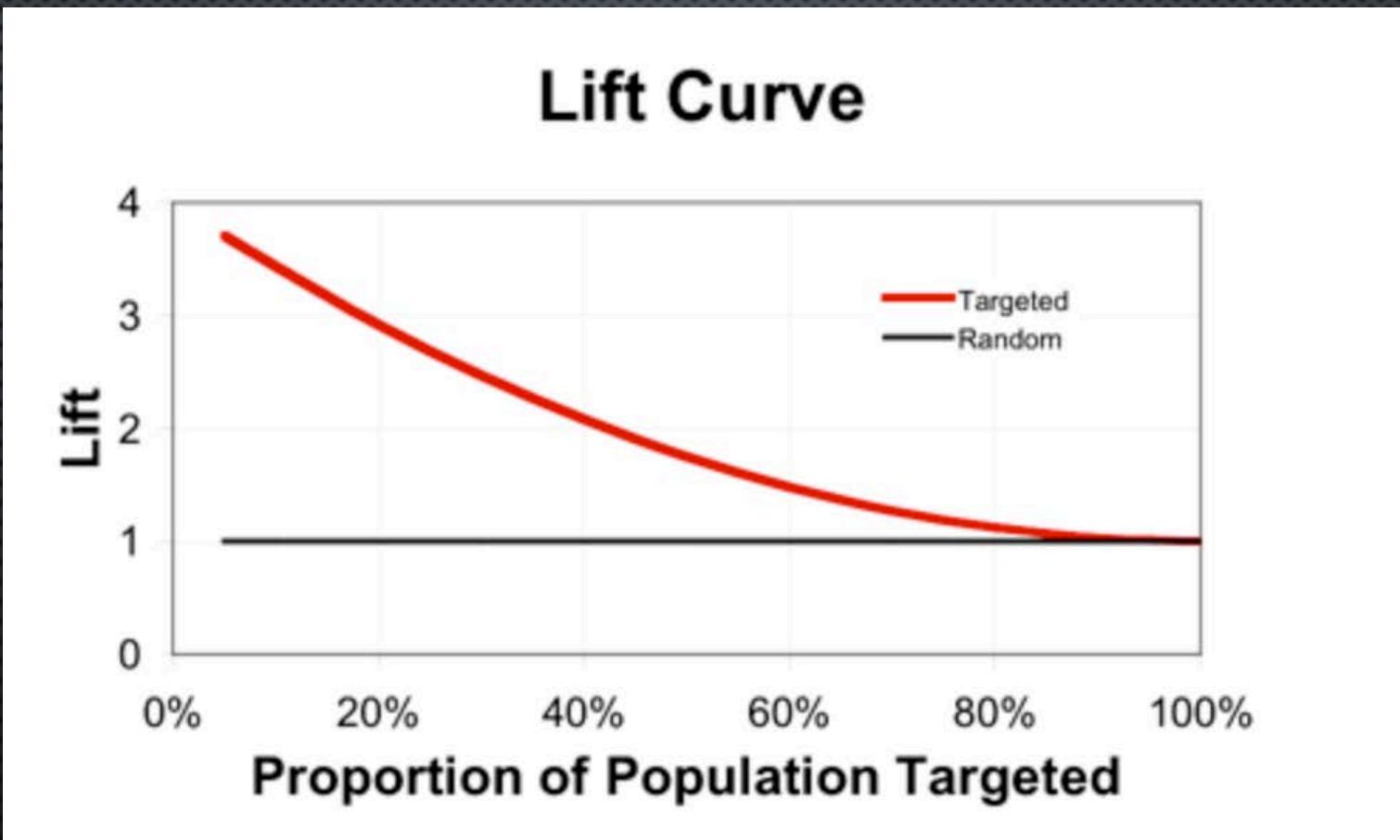
**GLM v CC**

- LOGISTICKA REGRESE**
- TARGET = SJEDNANO**
- CENA, VEK,**

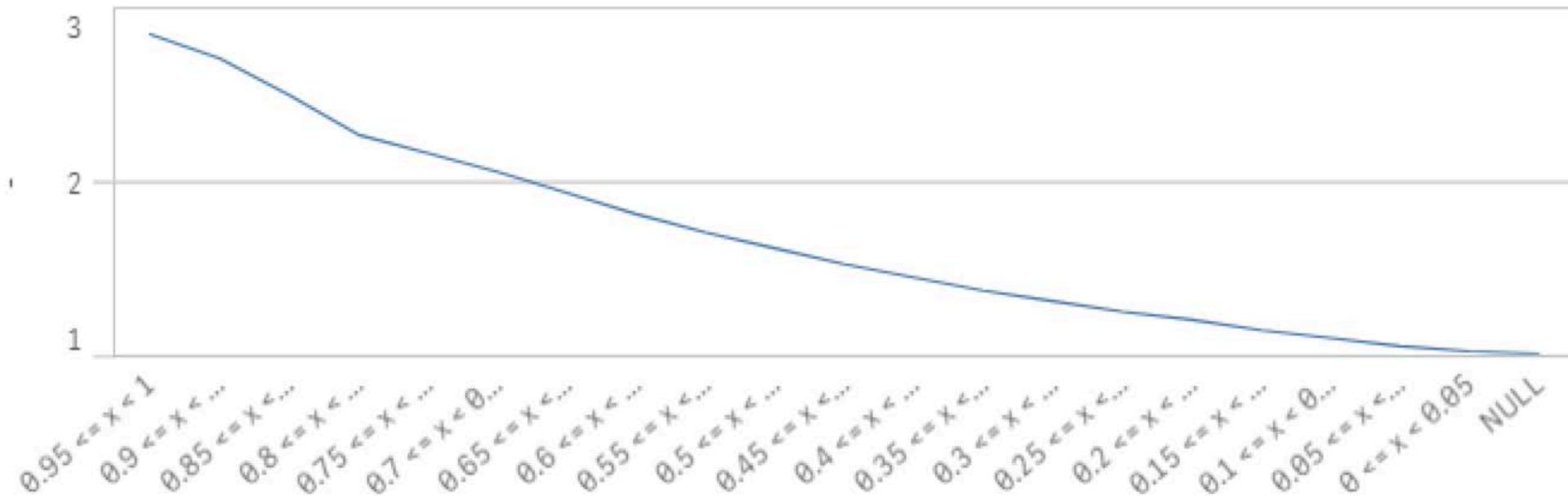
### Most important variables



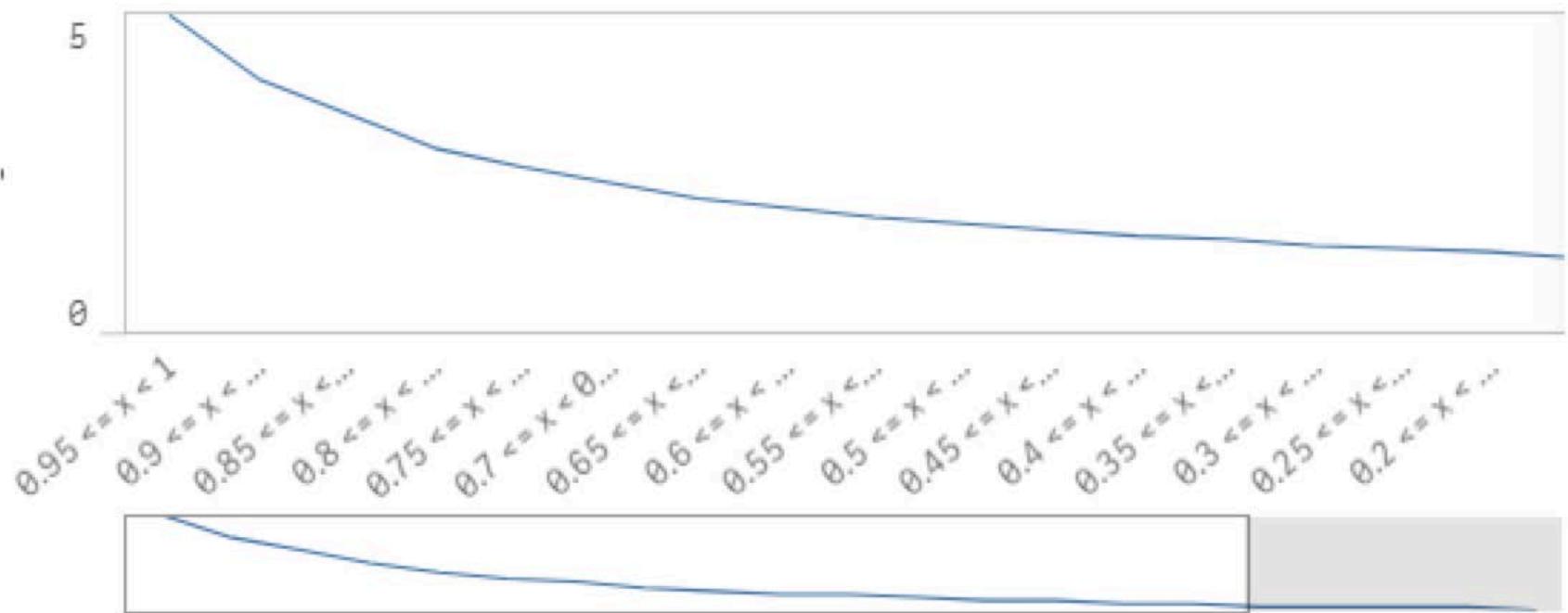
# LIFT KRIVKA



## Lift curve - probability of new contract



## Lift curve - probability of new contract \* yearly payment



## **2 MANA**

### **SITUACE**

- **EXPLORACE**
- **CASHFLOW**
- **OPTIMALIZACE MKT SPENDU**

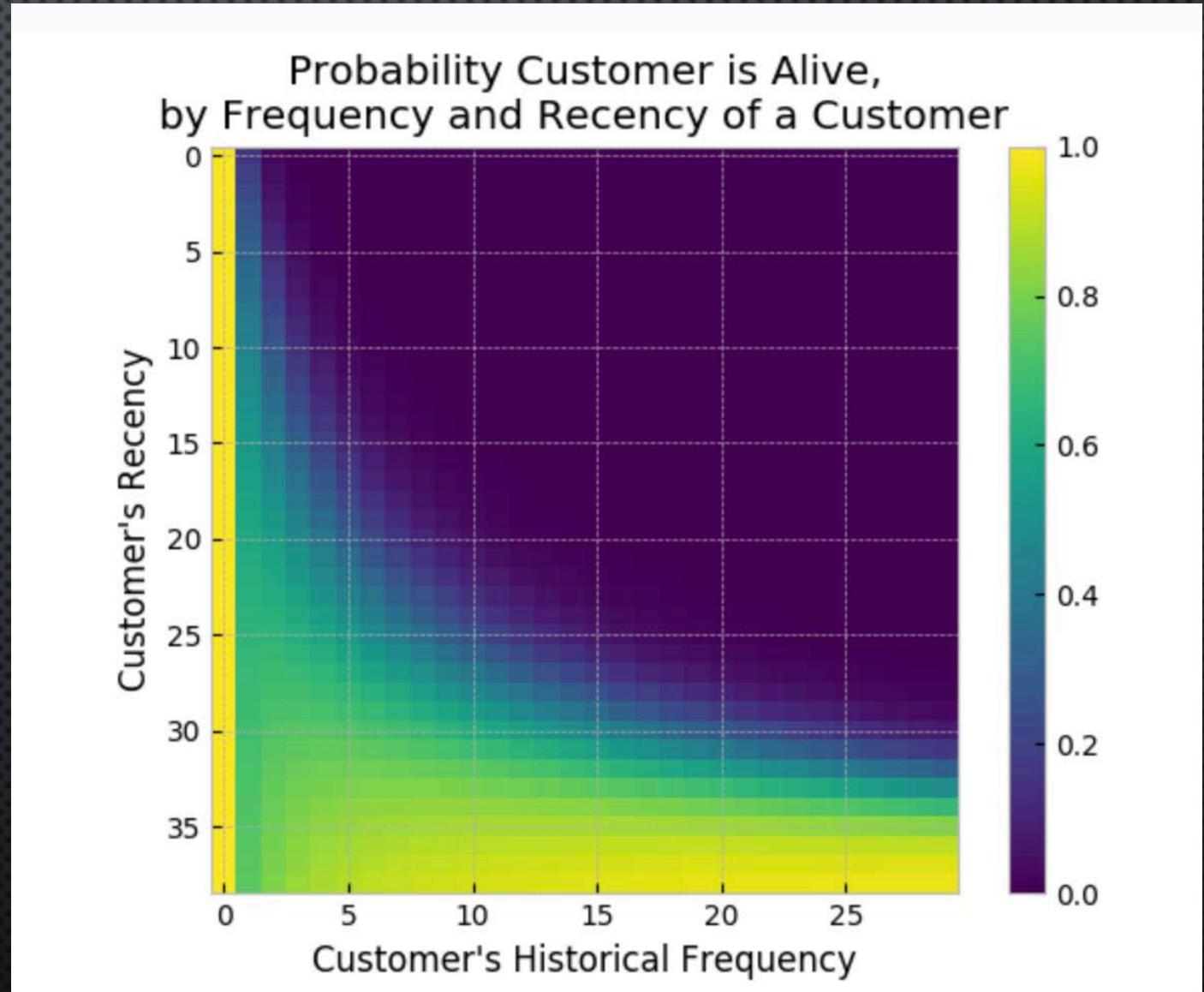
# CLV/CHURN RATES – PROBABILISTIC APP

- NUMBER OF TRANSACTIONS - POISSON PROCESS  $\lambda$ .
- TIME BETWEEN TRANSACTION - EXPONENTIAL  $\lambda$ ,

## THE BASIC PARETO/NBD MODEL

- “COIN” CUSTOMER CHURNS - PARETO
- “DICE” ORDERS - NEGATIVE BINOMIAL

# CHURN RATES



# **DATASENTICS CLV**

- SPOTREBA
- PRECHODOVE MATICE SPOTREBY
- MRTVY NENI MRTVY JEN MALO SPOTREBUJE

# DATASENTICS CLV

	do 1	do 2	do 4	do <10	do >10
z 1	94%	1%	1.1%	1.4%	2.5%
z 2	1.5%	95%	1.5%	0.9%	1.1%
z 4	1.9%	6.5%	79.3%	8.6%	3.7%
z <10	1.1%	3.1%	11.8%	70.4%	13.7%
z >10	3.6%	2.5%	4.3%	19.4%	70.3%

# **DATASENTICS CLV**

- **NAMODELOVAT X NOVÝCH ZAKAZNIKU**
- **STREDNÍ HODNOTA = CLV**

# **SROVNANI OBOU PRISTUPU**

- POTREBA HODNE ZAZNAMU PRO PROBABILISTIC APP
- DEFINICE AKTIVNIHO/MRTVEHO ZAKAZNIKA – 3 NAKUPY
- PROBLEM PREDIKCE AKTUALNI SPOTREBY

OTAZKY...?

# DEKUJI ZA POZORNOST

