

Stories: data mining v businessové praxi

Aneb výzvy při automatizaci datové dělničiny

Viktor Brada & Hynek Walner

MFF UK, 28.2.2018



Avast, Lagardere, Kiwi.com, Alza, Ahold, ČSAS, Sporty.cz, Snowboards.cz, Vivantis, Rockaway, Leo Express, Koloniál, Prodeti.cz, BigBrands, Bibloo, Rozbaleno.cz, Le Premier, Liftago, Lidská Síla, Adastra, Bonavita, COOP, Nielsen, Agrozet, Seznam, Škoda Auto, Vodafone UK, Sklizeno, Volkswagen, Daher, Freshlabels, Heureka, MF ČR, ...

viktor@stories.bi, hynek@stories.bi



PostgreSQL



...

Surrey	England	SU957435	495783	143522	51.18291	-0.63098	GU7 2	Wa
Somerset	England	ST707226	370749	122688	51.00283	-2.41825	BA8 0	So
Worcestershire	England	SO744675	374477	267522	52.30522	-2.37574	WR6 6	Ma
Essex	England	TM006190	600637	219093	51.83440	0.91066	CO5 7	Co
Worcestershire	England	SO995534	399538	253477	52.17955	-2.00817	WR10 2	W
Essex	England	TL575115	557500	211500	51.78000	0.28172	CM5 0	Ep
Essex	England	TL571112	557170	211283	51.77815	0.27685	CM5 0	Ep
Devon	England	ST140106	314090	110654	50.88896	-1.92659	B98 8	Ea
Powys	Wales	SO053712	305386	271238	52.33104	-1.52539	S17 3	Po
Gloucestershire	England	SO863164	386373	216413	51.84615	-2.89607	HR2 0	Gl
South Yorkshire	England	SK327822	432719	382237	53.33603	-1.58043	CV8 1	Sh
Worcestershire	England	SP051681	405102	268179	52.31170	-1.58497	CV8 1	Re
South Yorkshire	England	SK317807	431711	380770	53.32291	-3.00118	EX13 5	Sh
Herefordshire	England	SO385304	338538	230489	51.96946	-1.58043	CV8 1	Co
Warwickshire	England	SP286718	428680	271865	52.34411	-1.58043	CV8 1	Wa
Essex	England	TL957435	500650	174500	51.88336	0.89730	CO2 7	Co
Warwickshire	England	SP283719	428376	271950	52.34490	-1.58497	CV8 1	Wa
Devon	England	TL295975	329500	97500	50.77278	-3.00118	EX13 5	Ea
Shropshire	England	SP350666	350666	333195	52.89395	-1.58043	CV8 1	Sh
Staffordshire	England	SJ979577	397900	357773	53.11713	-2.03283	ST13 8	Sta
Greater Manchester	England	SJ095065	399599	306500	53.146514	-2.15963	M18 8	Ma
City of Edinburgh	Scotland	N1270744	327066	674436	55.95738	-3.16970	EH7 5	Cit
Staffordshire	England	SJ907486	390708	348651	53.03506	-2.14002	ST2 8	Cit
Surrey	England	TO045675	504500	167500	51.39693	0.49929	KT16 8	Su

You lose on noise.

Every modern business has a wealth of information, but 95 % of it is useless clutter. And that's a burden. A cost. A cognitive, technological, and financial cost.

Stories.

-1.92659	B98 8
-1.52539	S17 3
-2.89607	HR2 0
-1.58043	CV8 1
-1.58497	CV8 1
-3.00118	EX13 5

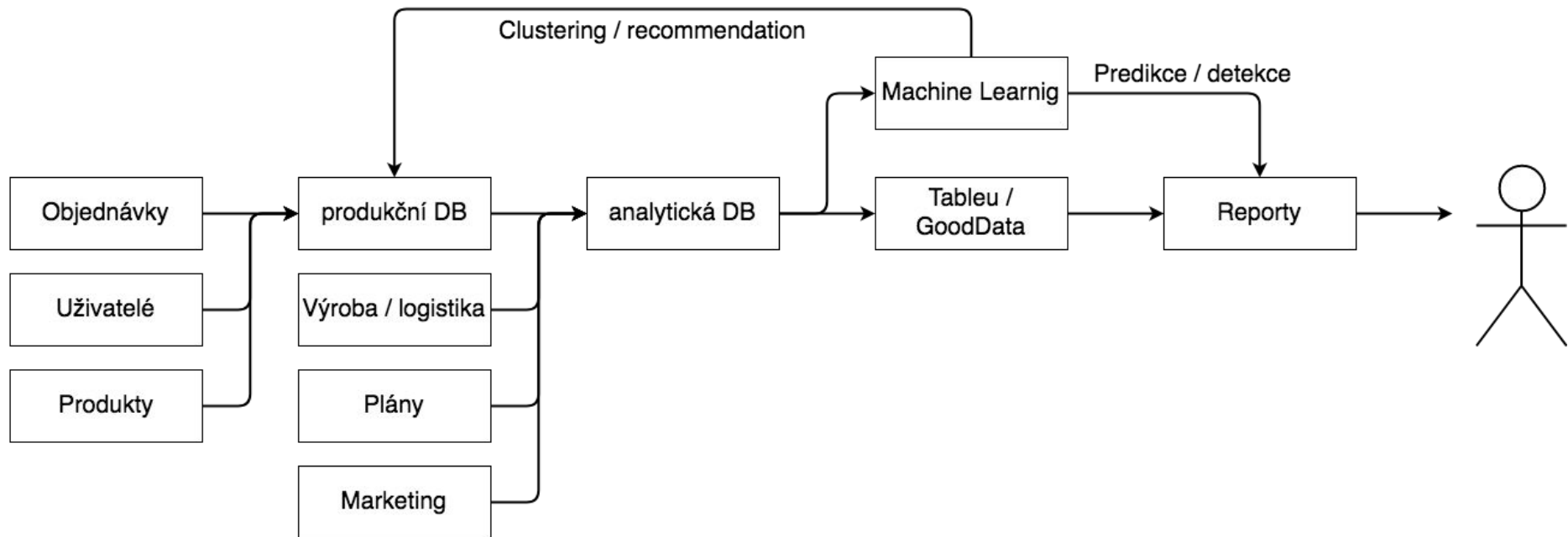
Co je Business Intelligence?

Co je Business Intelligence?

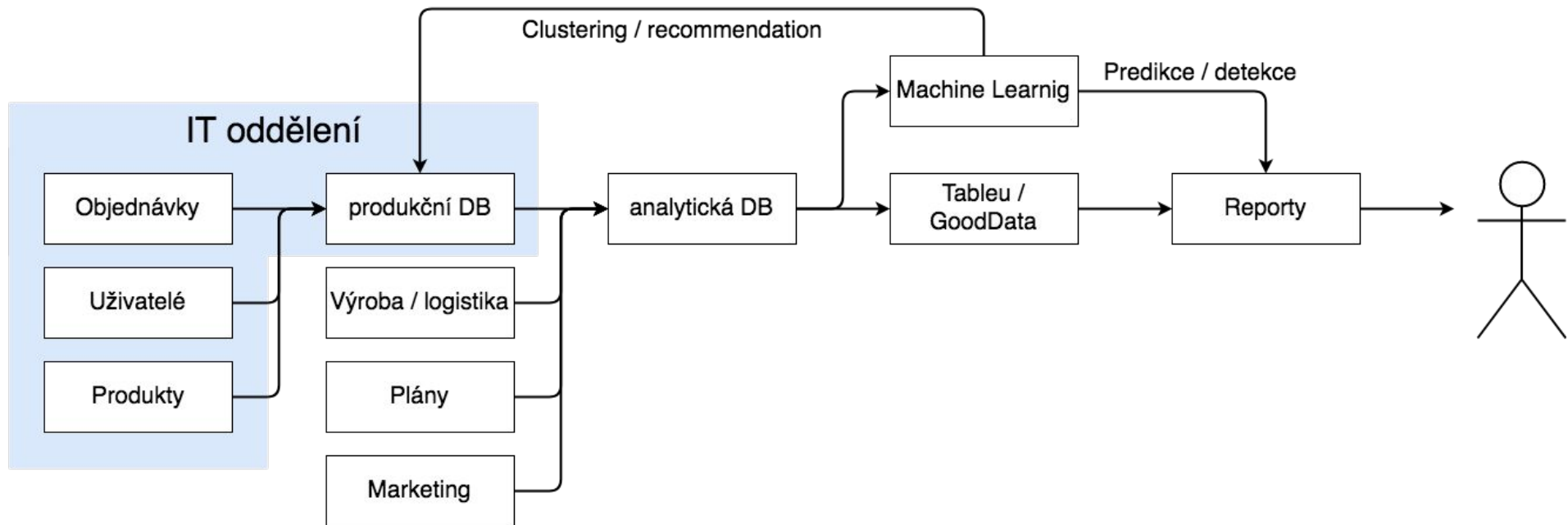
[...] applications, infrastructure and tools [...] that enable access to and analysis of information to improve [...] decisions and performance.

<https://www.gartner.com/it-glossary/business-intelligence-bi/>

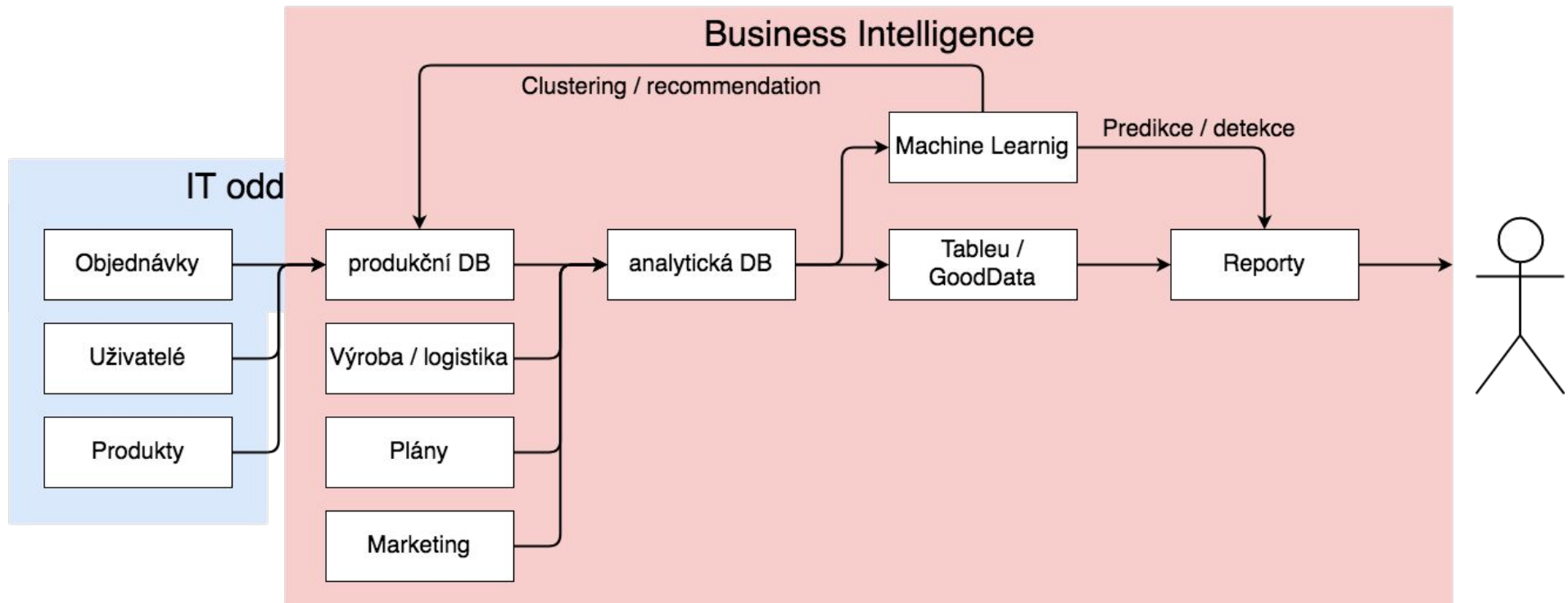
Co je BI?



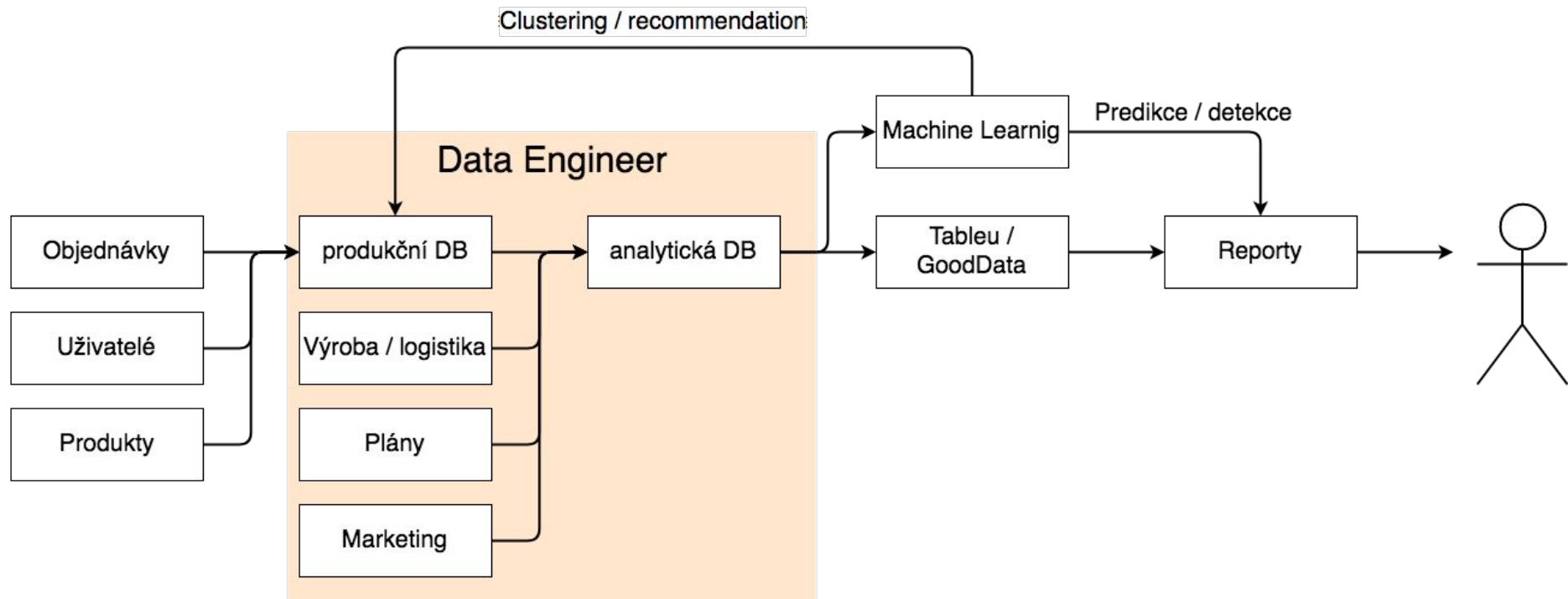
Co je BI?



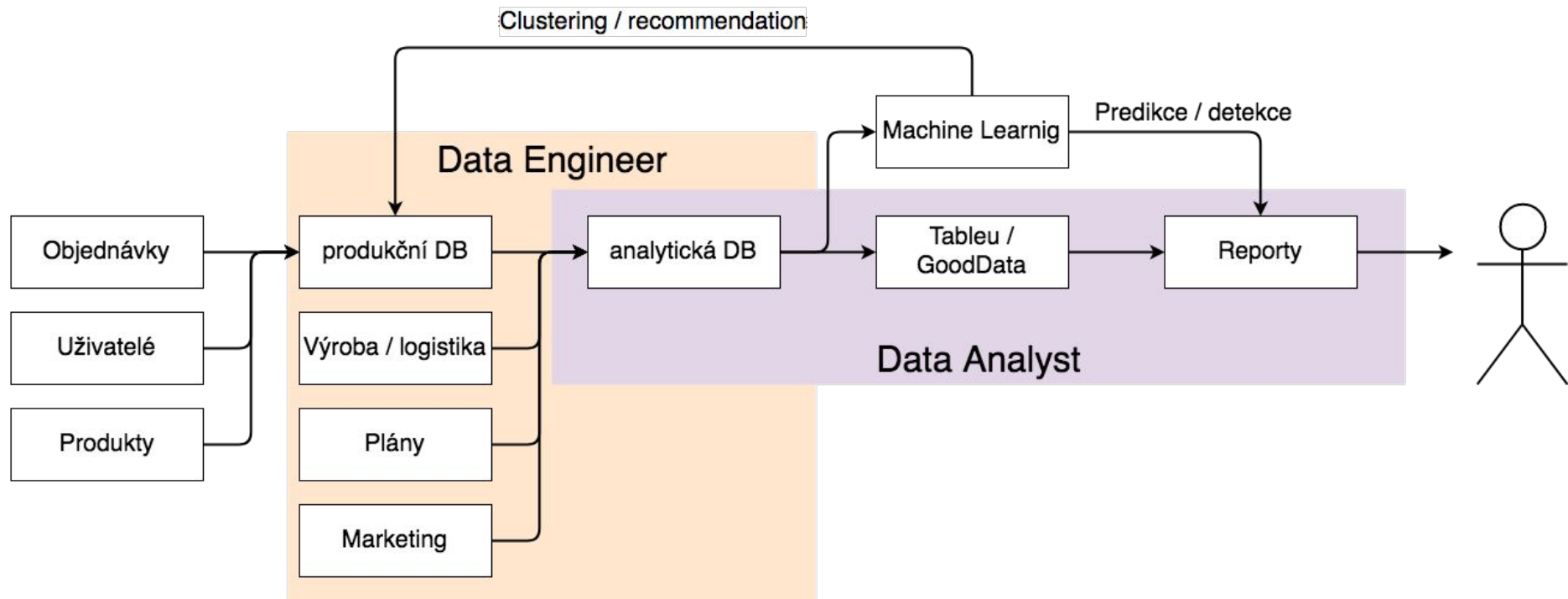
Co je BI?



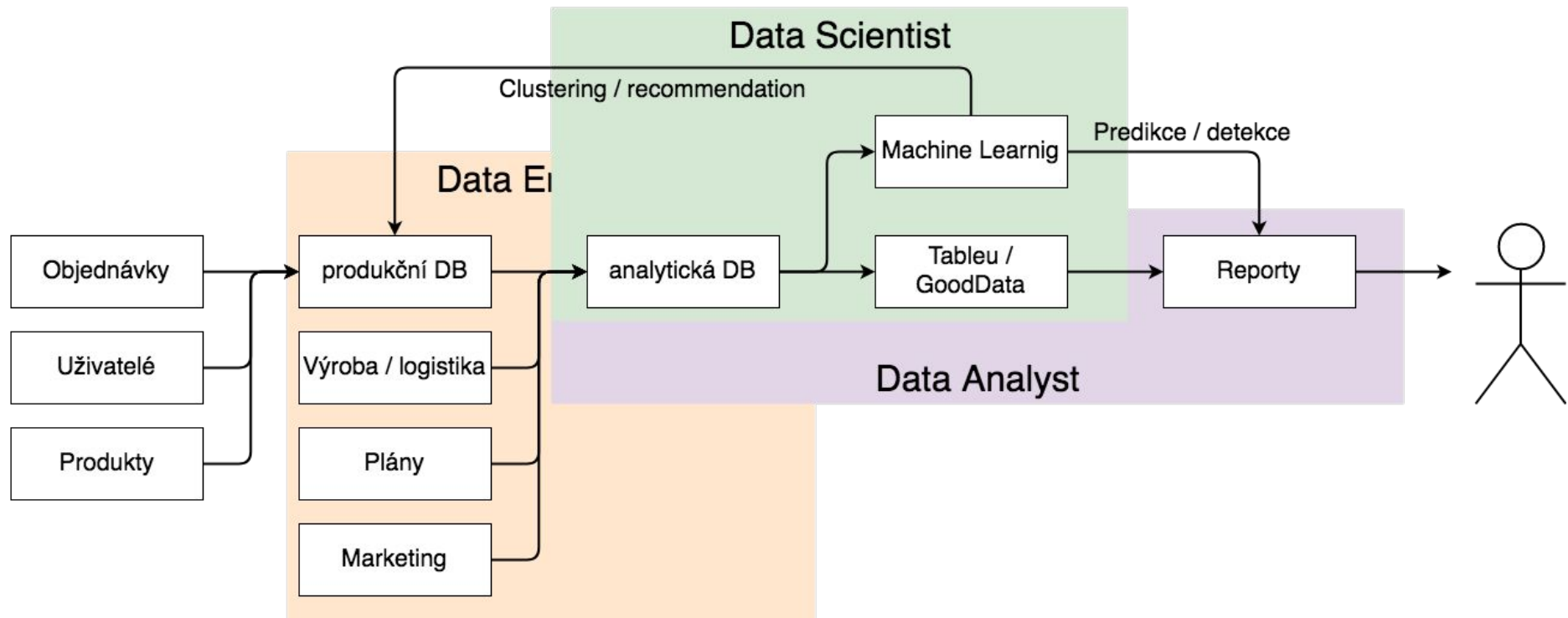
Co je BI?



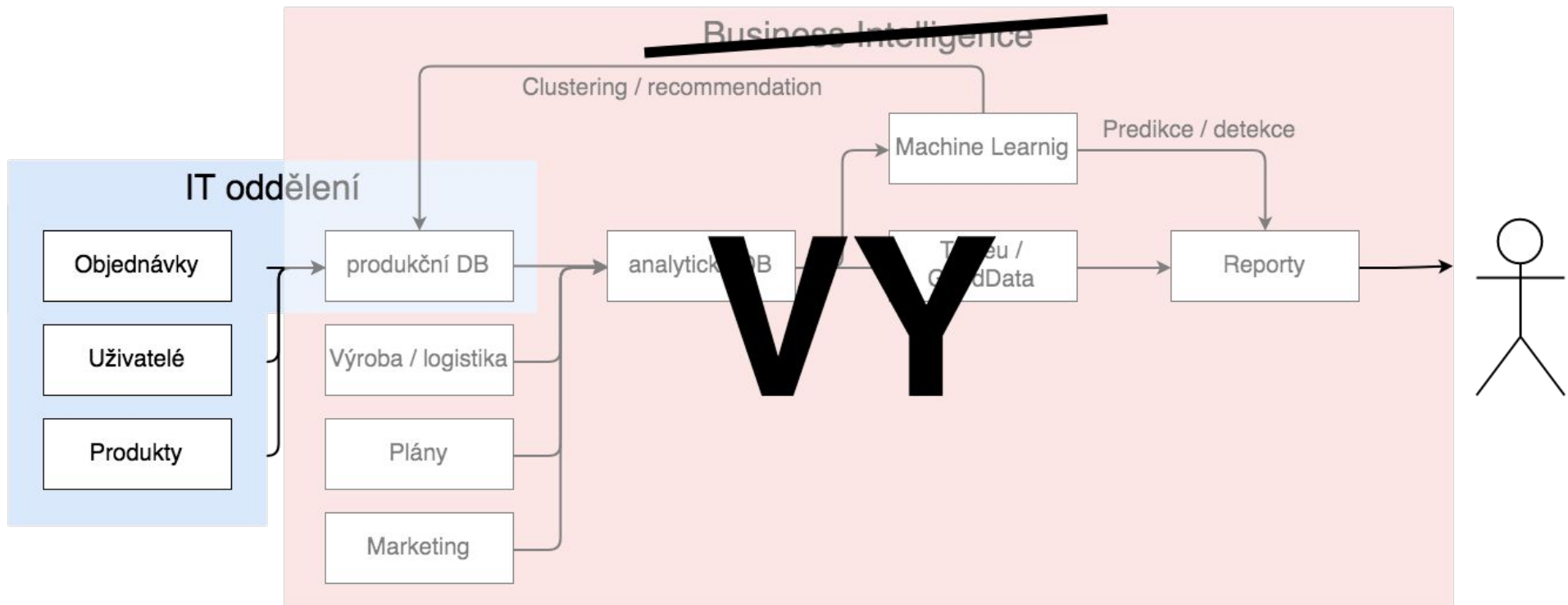
Co je BI?

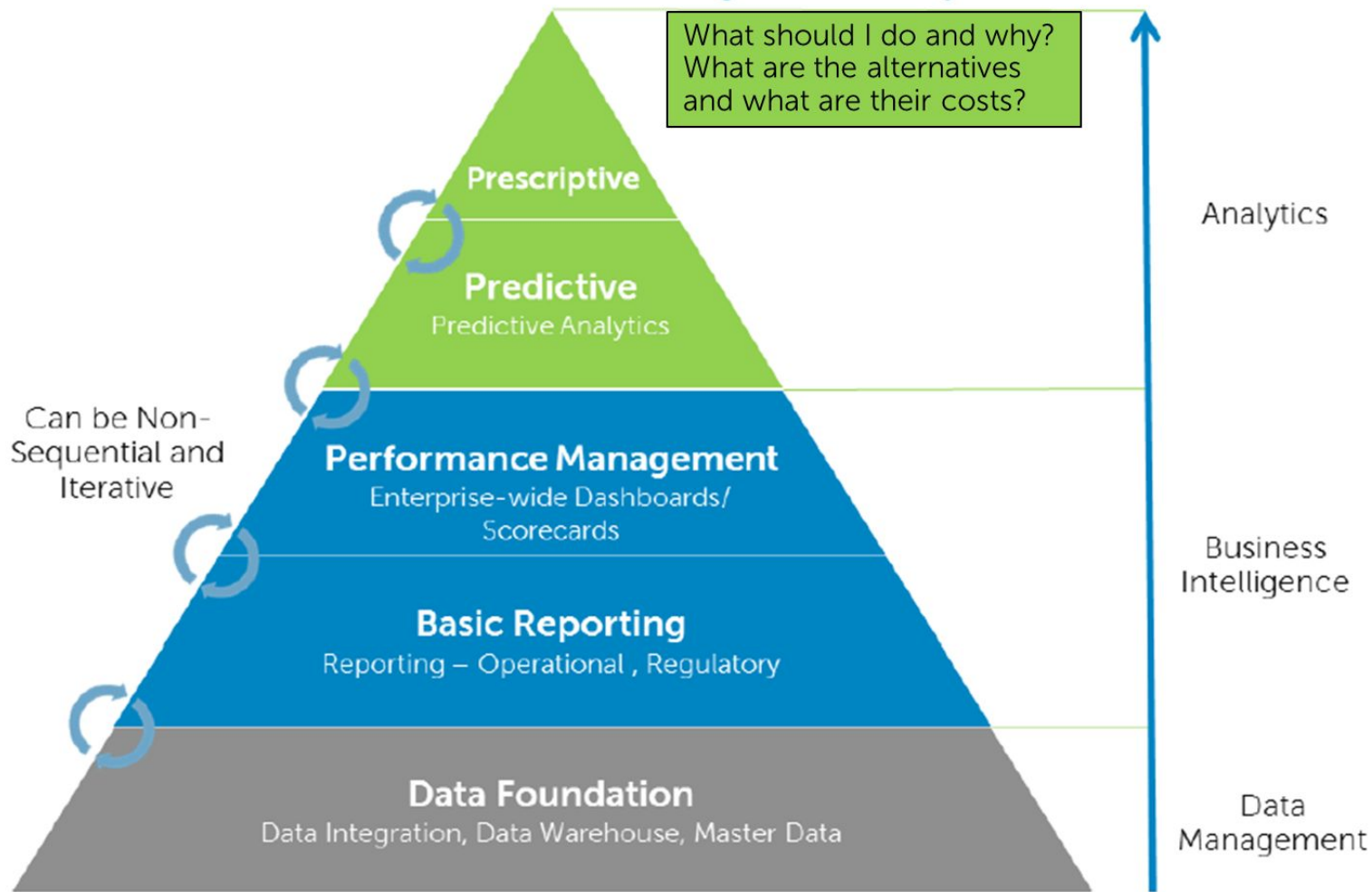


Co je BI?

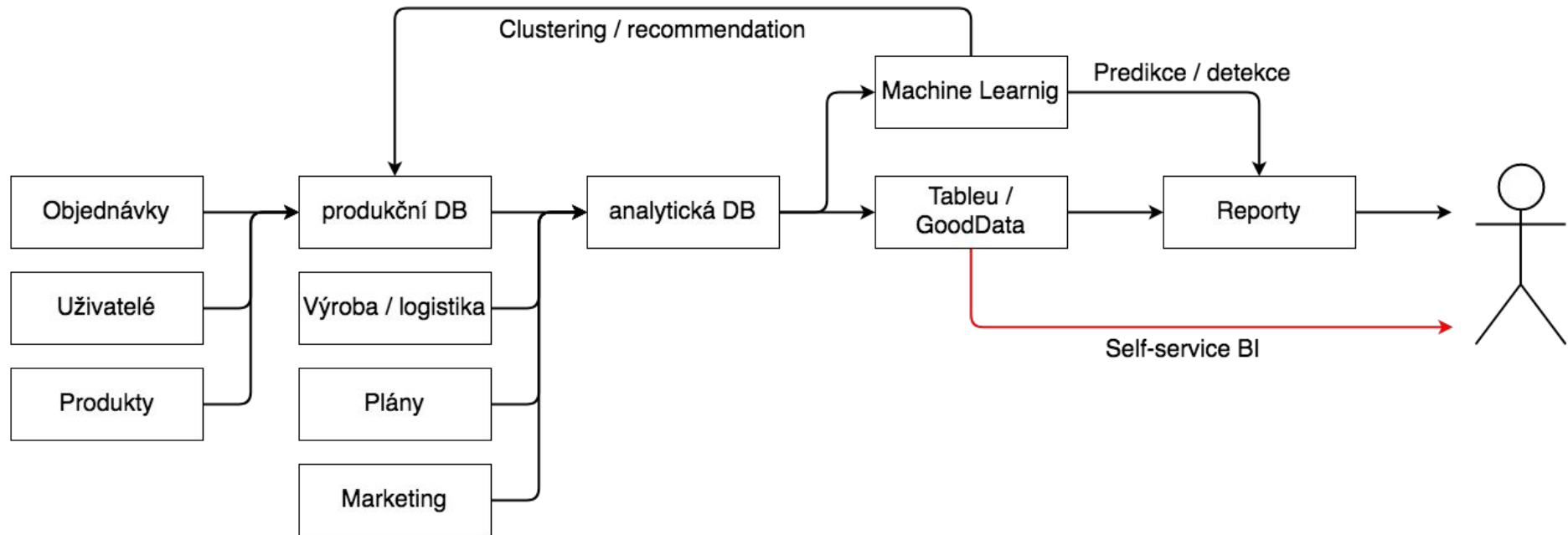


Co je BI?



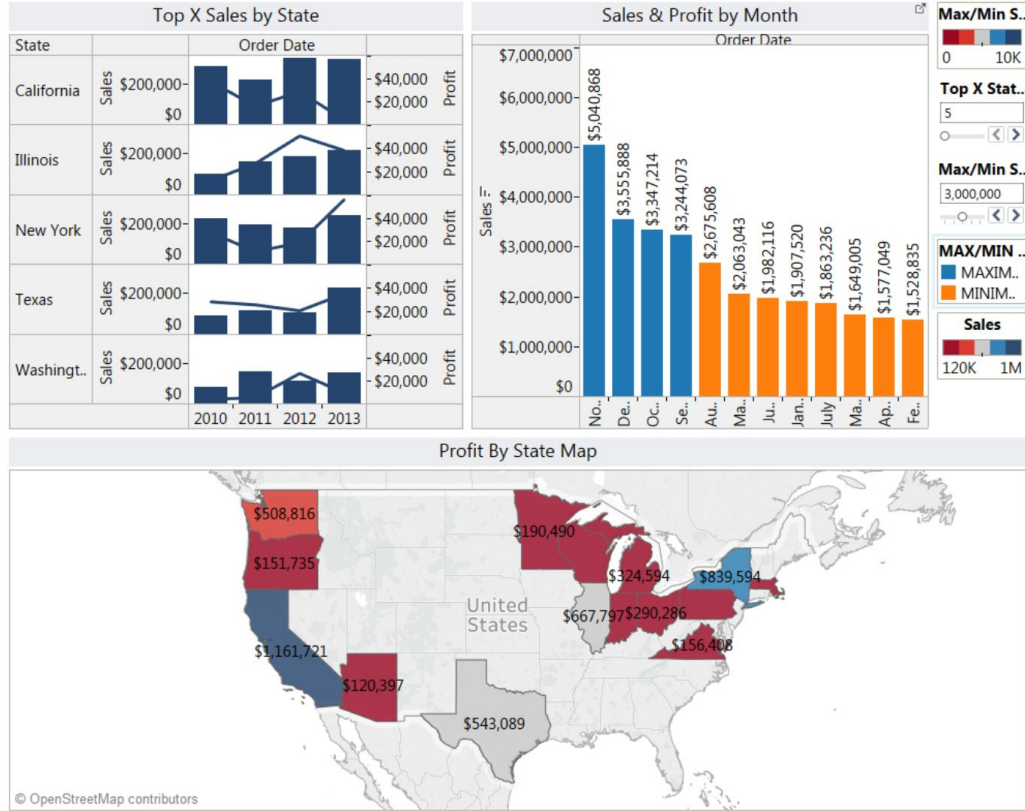


Co je Self-service BI?



From:

Sales Dashboard



Source: Tableau

Too many places to look in the data

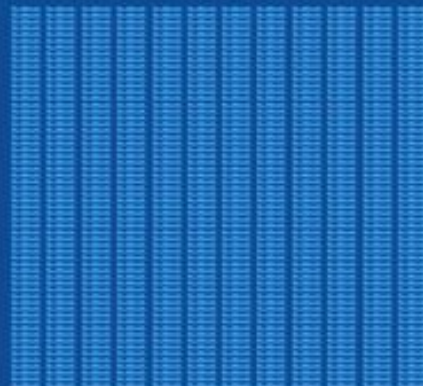
If there are just 10 Variables:
Graphs with
1 Variable at a Time

10



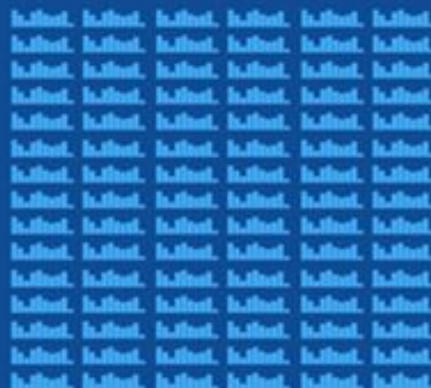
If there are just 10 Variables:
Graphs with
3 Variables at a Time

720



If there are just 10 Variables:
Graphs with
2 Variables at a Time

90



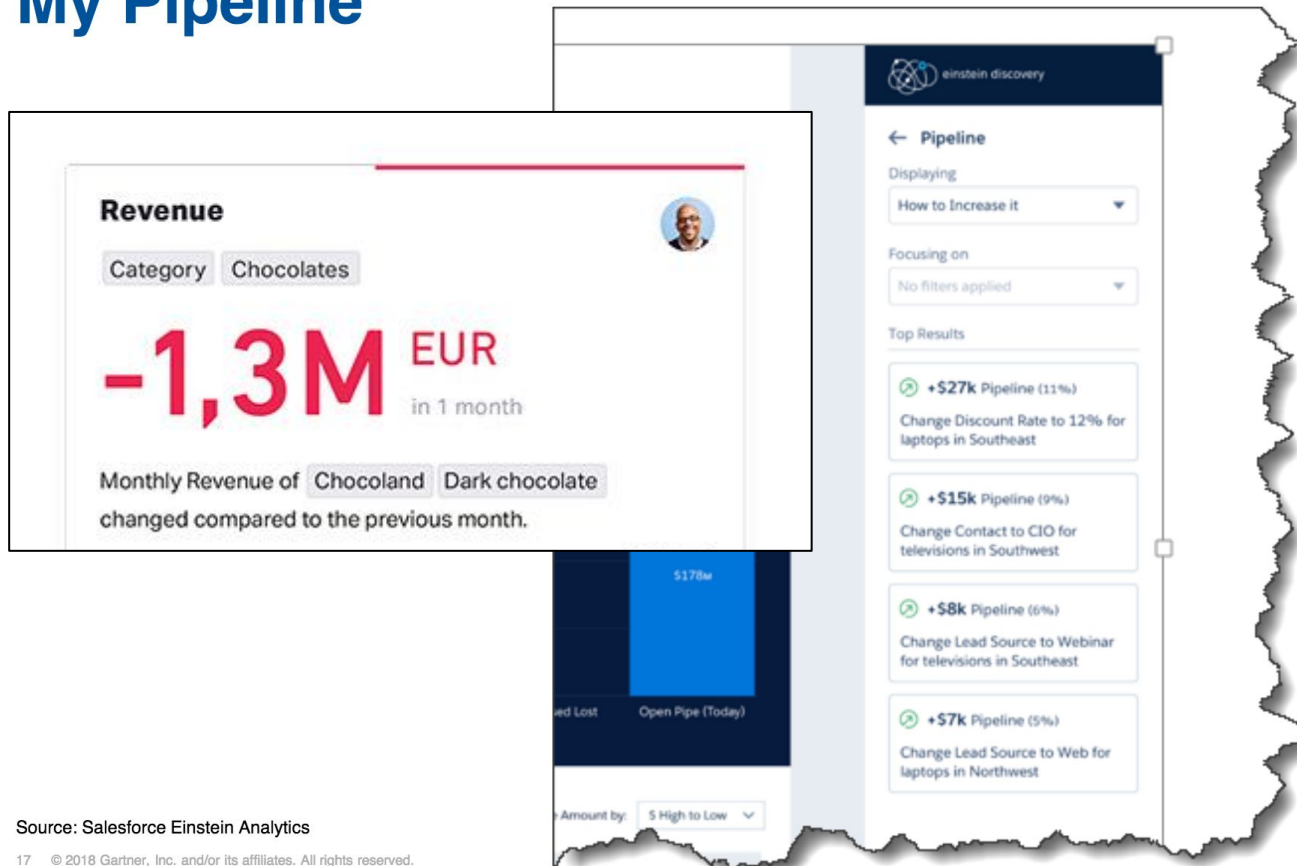
Who wants to review

820

Graphs?



To: Here Are the Top Four Things I Can Do to Increase My Pipeline



Source: Salesforce Einstein Analytics

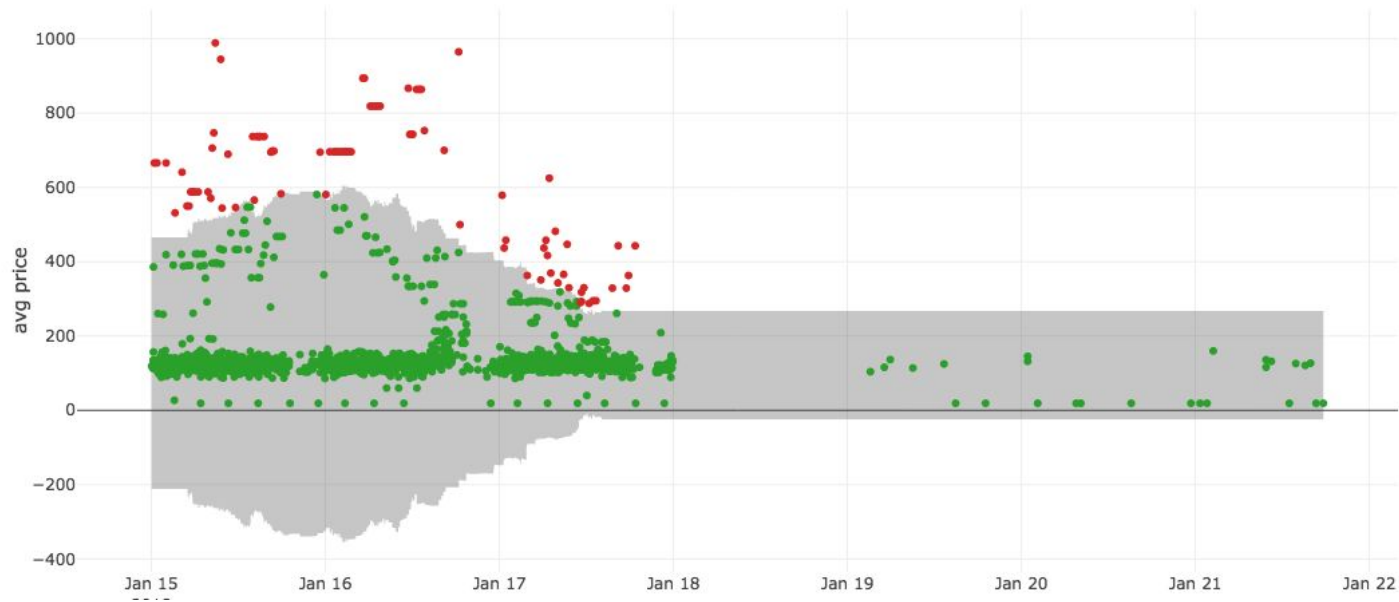
SPAs: What Will the ABI Market Impact Be?

By 2020, 50% of the analytic queries will be generated using search, NLP, voice or autogenerated.

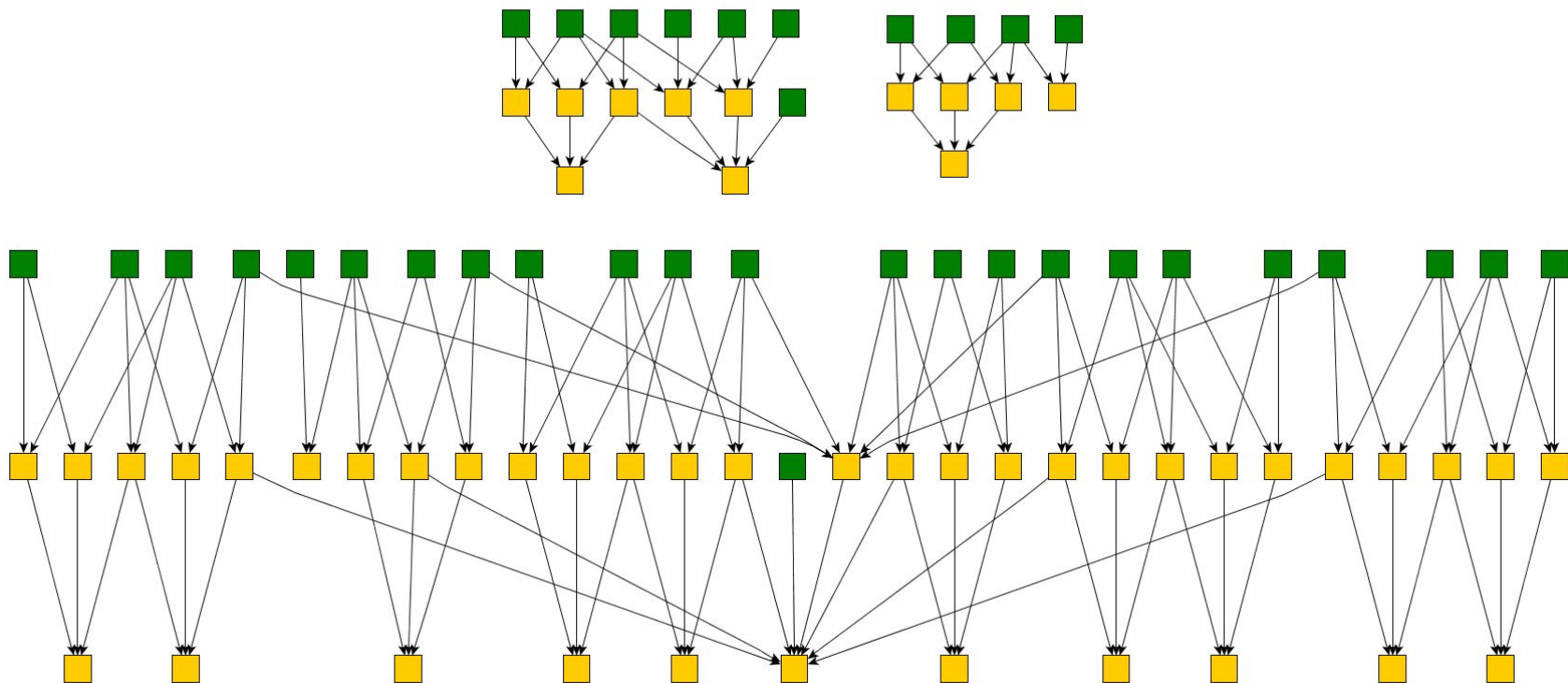
By 2020, 30% of today's data scientist tasks can be automated.

Co děláme?

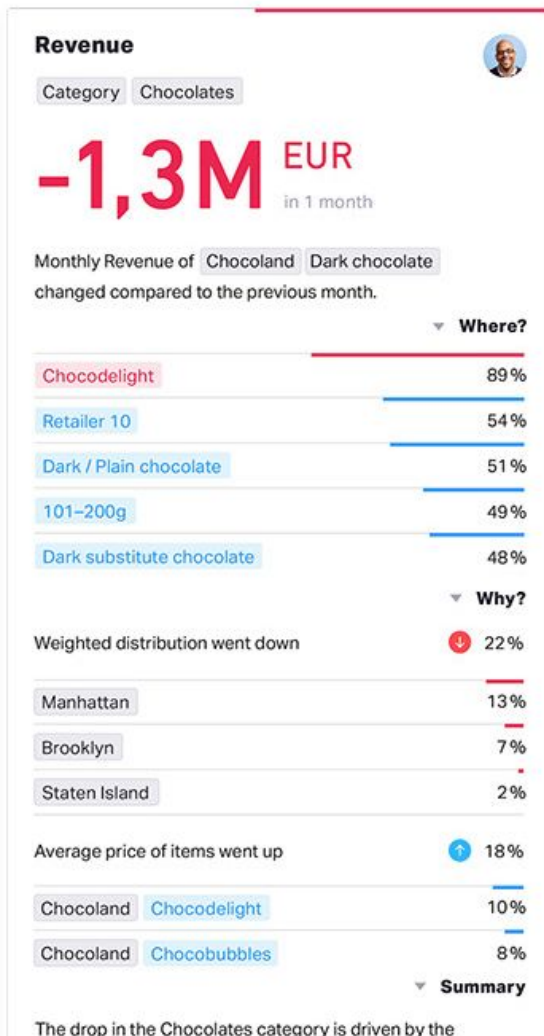
Uděláme z toho tohle:



A tohle:



Abychom zobrazili tohle:

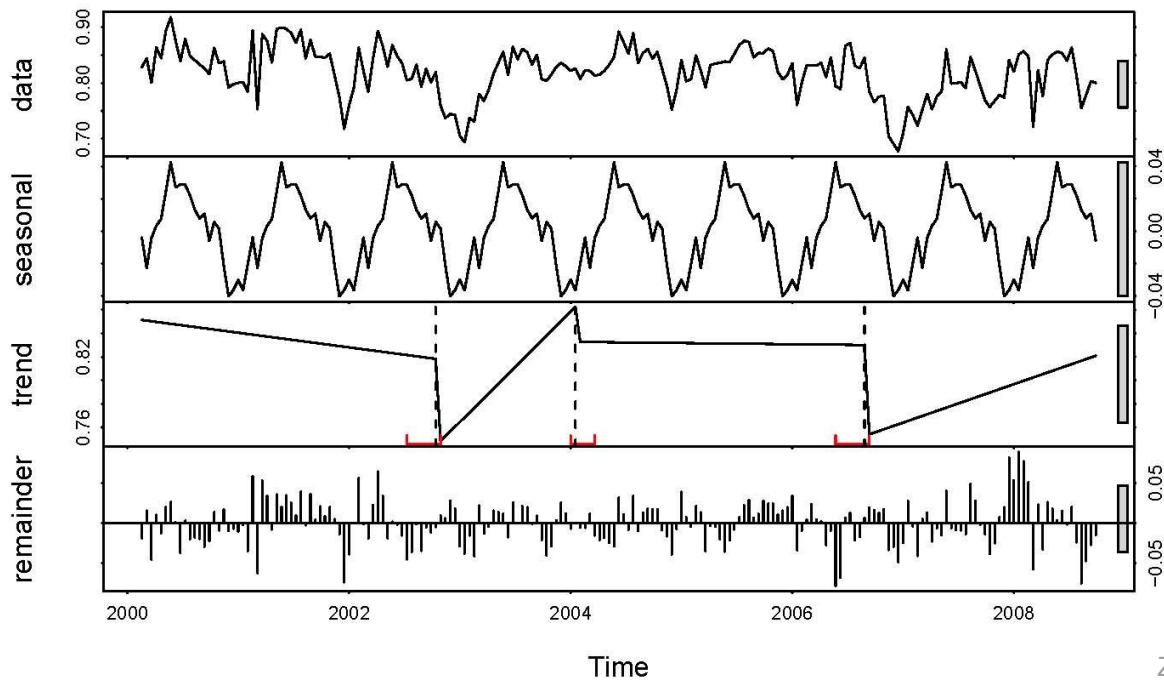


1. úloha: Detekce změny v časových řadách

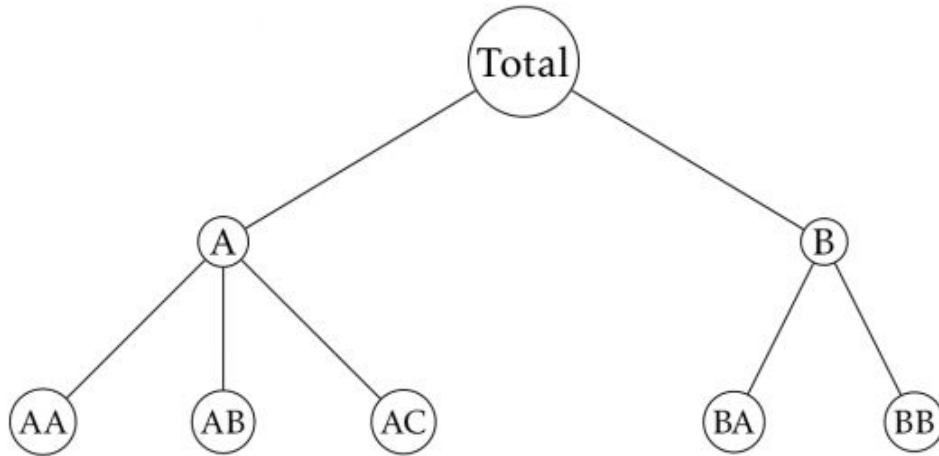
Businessové otázky

- Vypadly prodeje (hluboký propad),
- Klesají marže (klesající dlouhodobý trend + kde začal),
- Anomální událost,
- Neplním plán (odchylka od predikce).

Sezónnost, šum

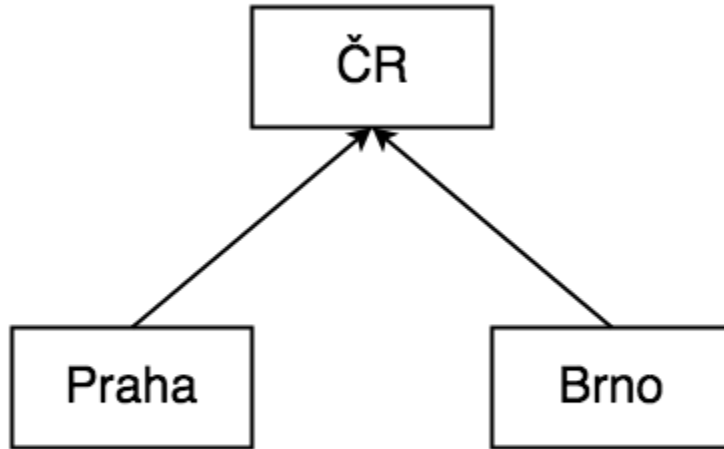


Jak využít strukturu v datech?

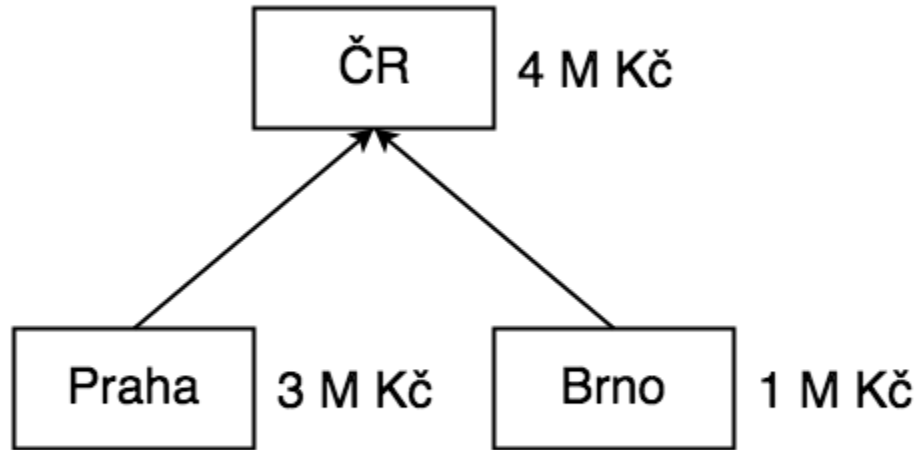


$$\begin{bmatrix} y_t \\ y_{A,t} \\ y_{B,t} \\ y_{AA,t} \\ y_{AB,t} \\ y_{AC,t} \\ y_{BA,t} \\ y_{BB,t} \end{bmatrix} = \begin{bmatrix} 1 & 1 & 1 & 1 & 1 \\ 1 & 1 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 1 \\ 1 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix} \begin{bmatrix} y_{AA,t} \\ y_{AB,t} \\ y_{AC,t} \\ y_{BA,t} \\ y_{BB,t} \end{bmatrix}$$

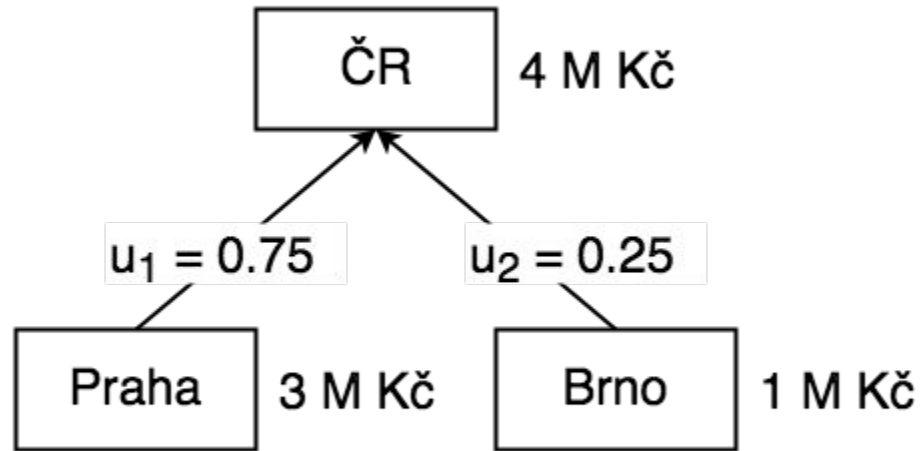
2. úloha: Kauzalita mezi detekovanými signály



2. úloha: Kauzalita mezi detekovanými signály

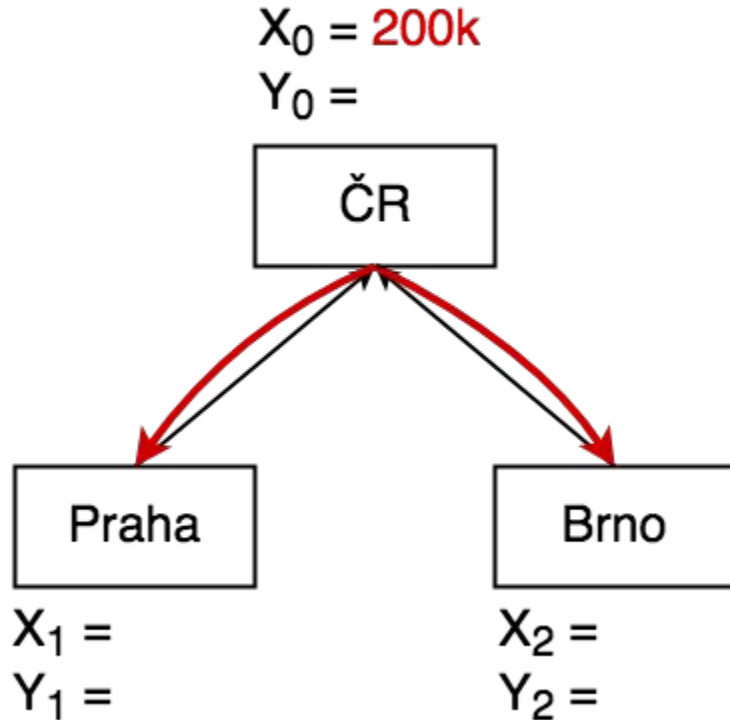


2. úloha: Kauzalita mezi detekovanými signály

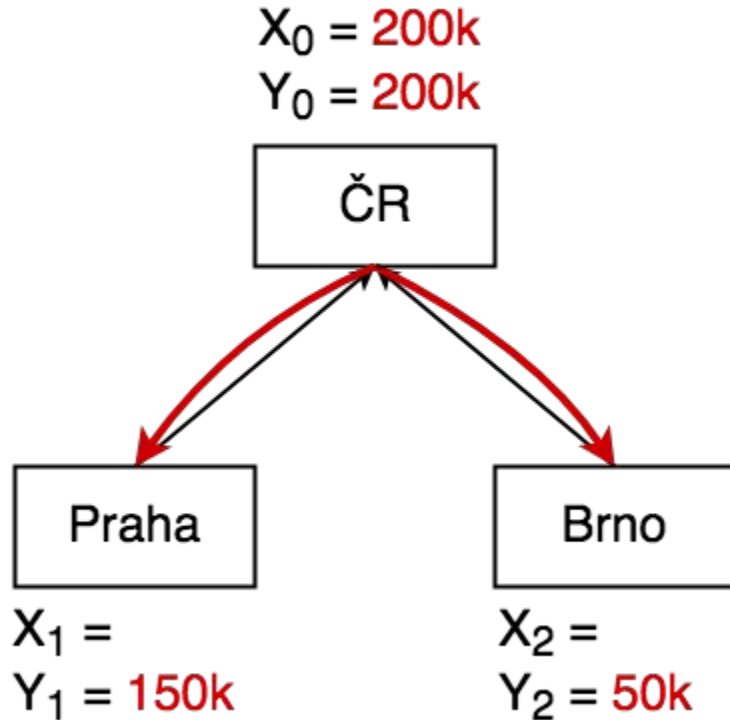


$$u_1 + u_2 = 1$$

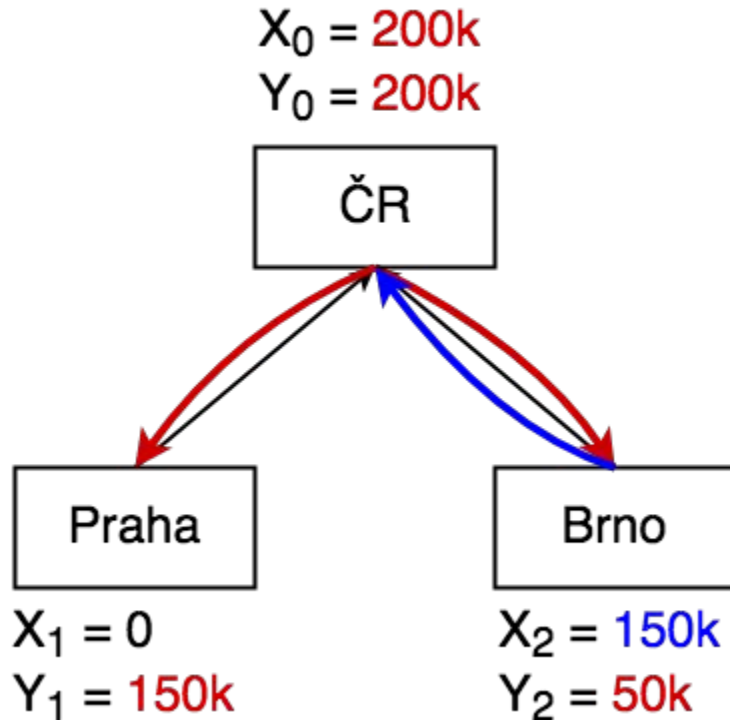
2. úloha: Kauzalita mezi detekovanými signály



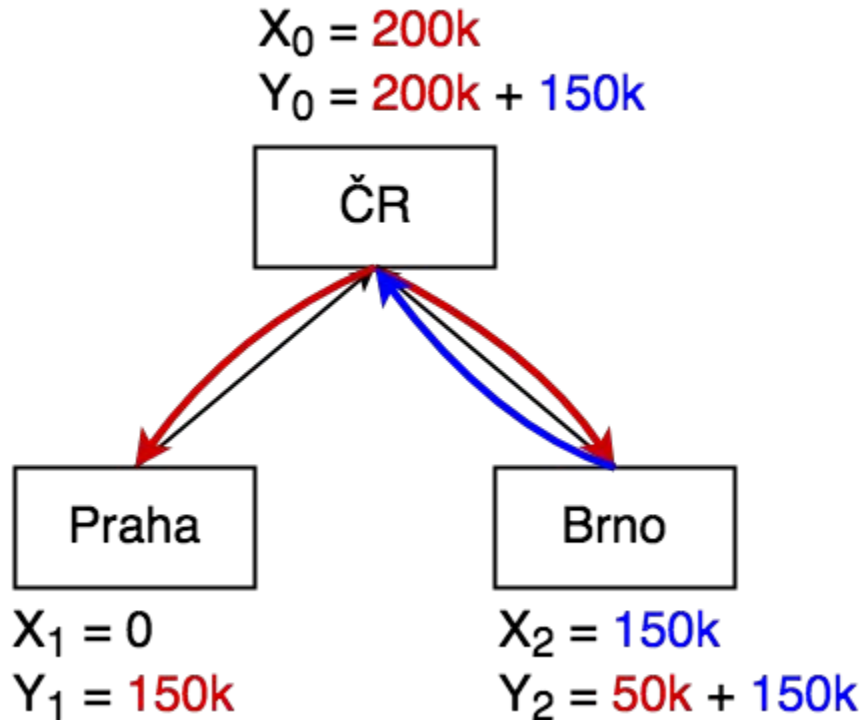
2. úloha: Kauzalita mezi detekovanými signály



2. úloha: Kauzalita mezi detekovanými signály



2. úloha: Kauzalita mezi detekovanými signály



$$u_1 + u_2 = 1$$

$$Y_0 = X_0 + X_1 + X_2$$

$$Y_1 = X_1 + u_1 X_0$$

$$Y_2 = X_2 + u_2 X_0$$

$$\begin{pmatrix} 1 & 1 & 1 \\ u_1 & 1 & 0 \\ u_2 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} \quad (1)$$

$$\begin{pmatrix} 1 & 1 & 1 \\ u_1 & 1 & 0 \\ u_2 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} \quad (1)$$

$$Ax = y \quad (2)$$

$$\begin{pmatrix} 1 & 1 & 1 \\ u_1 & 1 & 0 \\ u_2 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} \quad (1)$$

$$\mathbb{A}x = y \quad (2)$$

$$D(x, y) = \|\mathbb{A}x - y\|_2^2 \quad (3)$$

$$\begin{pmatrix} 1 & 1 & 1 \\ u_1 & 1 & 0 \\ u_2 & 0 & 1 \end{pmatrix} \begin{pmatrix} x_0 \\ x_1 \\ x_2 \end{pmatrix} = \begin{pmatrix} y_0 \\ y_1 \\ y_2 \end{pmatrix} \quad (1)$$

$$\mathbb{A}x = y \quad (2)$$

$$D(x, y) = \|\mathbb{A}x - y\|_2^2 \quad (3)$$

$$\min_x \{D(x, y) + \lambda R(x)\} \quad (4)$$

3. Úloha: Jak signály seřadit podle důležitosti?

12M Kč

37 %

14k kusů

125 hodin

...

Ostatní analytické úlohy

“Most firms that thinks they want advanced AI/ML really just need linear regression on cleaned-up data.” -- Robin Hanson (@Twitter)

- predikce (zpoždění kamionů, prodeje v maloobchodech), klasifikace (analýza košíku)
- rychlé iterování metod (vždy začnět něčím jednoduše interpretovatelným)
- zohledňování business vazeb a jejich dopad na model
- pro high-end metody musí firma vyvinout velký závazek (fine-tuning parametrů, A/B testování, nasazení v produkci a garance správnosti celé pipeline)

Co čekat od práce s daty ve firmách?

- nebojte se vracet zadání
- připravte se na zahazování vašich projektů
- opakovatelná použitelnost
- data-quality check
- vždycky začněte nejjednodušším (správným) řešením, pokud nebude fungovat, ušetříte si spousty budoucího smutku
- lidi na to nekoukaj (a nerozumí)



Stories.

www.stories.bi

{viktor, hynek} @ stories.bi