

DĚLAT
DOBRÝ SOFTWARE
NÁS BAVÍ

PROFINIT

B0M33BDT

Technologie pro velká data

Petr Paščenko, Jan Hučín

25. 9. 2019

The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are scattered across the white background, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, some appearing to float above others, which adds depth to the overall design.

Co?

Cíle předmětu

- › Představit technologie pro zpracování velkých dat
- › Architektonický pohled
 - systém jako celek a jeho součásti
- › Uživatelský pohled
 - jednotlivé technologie a jak je správně používat
- › Programátorský pohled
 - jak vyvíjet aplikace v prostředí velkých dat
- › Byznysový pohled
 - Data Science jako motivace pro rozvoj Big Dat
- › Ukázat skutečný život
na skutečných projektech



The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are scattered across the white background, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, some appearing to float above others, which adds depth to the overall design.

Jak?

Důležité odkazy

› Nástěnka na courseware

- <https://cw.fel.cvut.cz/b191/courses/b0m33bdt/start>
- vše na jednom místě
- přednášky

› Popis předmětu

- <https://www.fel.cvut.cz/cz/education/bk/predmety/47/73/p4773206.html>

› Úložiště ke cvičení

- GitHub, adresa bude zveřejněna
- zadání cvičení, vzorová řešení, odkazy na manuály

› Výpočetní cluster Metacentrum

- <https://www.metacentrum.cz/cs/Sluzby/Hadoop/>
- Požádat o přístup (včas!)

Přednášky

› Úvod, přehled, organizace, motivace

- Co jsou Big Data a kde se tady vzala
- vztah Big Data a Data Science
- aplikace Big Data technologií v průmyslu
- organizace, za co zápočet, představení zápočtových úloh

› Architektura clusteru

- Hadoop, distribuce Cloudera a Hortonworks, správa zdrojů YARN, atd.

› Storage

- HDFS, formáty ukládání a komprese dat, HIVE, Impala

› Map-reduce

- softwarové paradigma a implementace algoritmů, vztah k sql

› Apache Spark

- Distribuované výpočty v RAM a zpracování streamovaných dat (Kafka)

› Big Data Science a Datové Architektury

- page rank, kolaborativní filtrování, SNA
- typické architektury Big Data řešení



*"So what am I doing here—knighting,
beheading, or what?"*

CN
COLLECTION

Cvičení

› První kroky na clusteru

- Připojení, osahání, práce s více file systémy, základy HDFS

› Opakování potřebných dovedností

- Linux, shell utility, SQL, Python

› Hive

- Tabulky, int./ext., vytváření, správa, dotazy, partitions

› Spark

- Seznámení se Sparkem, transformace a akce v praxi
- Dva přístupy: RDD (map-reduce) a SQL (data frame)

› Praktický test

- zpracování praktických úloh na velkých datech pomocí jednotlivých technologií

Komplikovaný rozvrh

- › **Cíl: méně teorie, více cvičení**
 - Dotace předmětu 2/1 – posuneme směrem k 1,5/1,5
- › **Lichý týden**
 - přednáška středa: 9:15-10:45 [KN:E-127](#)
- › **Sudý týden varianta A**
 - přednáška středa: 9:15-10:45 [KN:E-127](#)
 - cvičení: 2 paralelky 11:00-12:30 a 12:45-14:15 [KN:E-307](#)
- › **Sudý týden varianta B**
 - supercvičení
 - 2 paralelky: 9:15-11:30 a 12:00-14:15 [KN:E-307](#)

Za co bude zápočet

› zisk aspoň 25 bodů z 50, dvě možné cesty

1. Vypracování zápočtové úlohy

- samostatná analýza velkého datového souboru
- zodpovězení předepsaných a vlastních analytických otázek
- výkonnostní měření v závislosti na velikosti vstupu
- dokumentace postupu a sepsání závěrečné zprávy (cca 10 stran)
- odevzdání dokumentovaných zdrojových kódů – github
- klasifikační kritéria:
 - Analytické výstupy – 40%
 - Programátorské postupy a kvalita kódu – 30%
 - Analytická zpráva a kvalita textu – 30%
- dohodnout: sergii.stamenov@profinit.eu

2. Testy a úkoly

- průběžný test z teorie, max. 10 bodů
- dva domácí úkoly, každý max. 5 bodů
- závěrečný praktický test, max. 30 bodů

Za co bude zkouška

- › Teorie, souvislosti, elementární praktické aplikace
 - za součet bodů ze zápočtu a samotné ústní zkoušky

The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are scattered across the entire frame, creating a complex, layered, and three-dimensional effect. The shapes vary in size and orientation, some appearing to float above others, which adds depth to the visual presentation.

Kdo?

Jan Hučín

› Zaměření

- Data Science a Machine Learning
- Big Data Byznys Aplikace

› Praxe

- V Profinitu 4 roky, konzultant
 - Data Science
 - klasifikační modely, survival analysis, časové řady
 - Big Data
 - popularizace a školení
- Scio
 - Psychometrika, datová analytika
- Institut informatiky
 - Výuka matematiky a statistiky



Marek Sušický



› Zaměření

- Databáze – Oracle, PostgreSQL, grafové db apod.
- Security, fraud, síť
- Big data

› Praxe

- V Profinitu 9 let
- Práce pro velkou banku, letiště, telekomunikačního operátora
- Praxe v mezinárodním prostředí

Stamenov Sergii

- › Zaměření
 - Data Science, Machine learning, Big Data
 - Vyvoj .NET
- › Praxe
 - V Profinitu 4 roky
 - Data science pro T-Mobile
 - Spark školení



Petr Paščenko

› Zaměření

- Data Science a Machine Learning
- Big Data Byznys Aplikace

› Praxe

- V Profinitu 7 let, konzultant, R&D
 - Data Science
 - vývoj vztahových a podobnostních modelů nad velkými daty
 - Grantový výzkum zaměřený na IT bezpečnost
 - Webmining, SNA, text-mining, grafová analytika, zabezpečení online kanálů
 - Big Data
 - aplikace a školení
 - Vývoj Java
- ČEZ, 2 roky
 - Analytik trhů
- Eccam, 2 roky
 - Vývoj c++



<http://cz.linkedin.com/in/pascenko/>

The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray. They are arranged in a way that creates a sense of depth and movement, as if they are floating or falling from the top right towards the bottom left. The overall effect is a modern, digital aesthetic.

Big Data?

Co jsou Big Data?



› Je to složité.

Co jsou Big Data?

- › Prodloužení několika trendů:
- › Data
 - od malých dat k velkým
 - od jednoduchých ke složitým
- › Databáze
 - od souborů přes relační paradigma k distribuovaným clusterům
- › Programování
 - od procedurálního programování k funkcionálním frameworkům
- › Data Science
 - od výběrových statistik k detailní kontextové analýze
- › Postupně si je projdeme



The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include rectangles, parallelograms, and trapezoids, are rendered in various shades of light gray and white. They are arranged in a way that creates a sense of depth and movement, as if they are floating or shifting in a three-dimensional space. The overall effect is a modern, minimalist aesthetic that suggests a digital or data-driven environment.

Data?

A Very Short History Of (Big) Data

Gil Press pro [forbes](#)



- › **1944:** Fremont Rider, Wesleyan University Librarian
 - „American university libraries were doubling in size every sixteen years.“
 - „Yale Library in 2040 will have approximately 200,000,000 volumes, which will occupy over 6,000 miles of shelves“
- › **1961:** Derek Price, Science Since Babylon
 - „the number of new journals has grown exponentially rather than linearly, doubling every fifteen years and increasing by a factor of ten during every half-century.“
 - „Each scientific advance generates a new series of advances at a reasonably constant birth rate, so that the number of births is strictly proportional to the size of the population of discoveries at any given time. “
- › **1975:** The Ministry of Posts and Telecommunications in Japan
 - „information supply is increasing much faster than information consumption“
 - „the demand for information provided by mass media, which are one-way communication, has become stagnant, and the demand for information provided by personal telecommunications media, which are characterized by two-way communications, has drastically increased.“

A Very Short History Of Big Data

Gil Press pro [forbes](#)

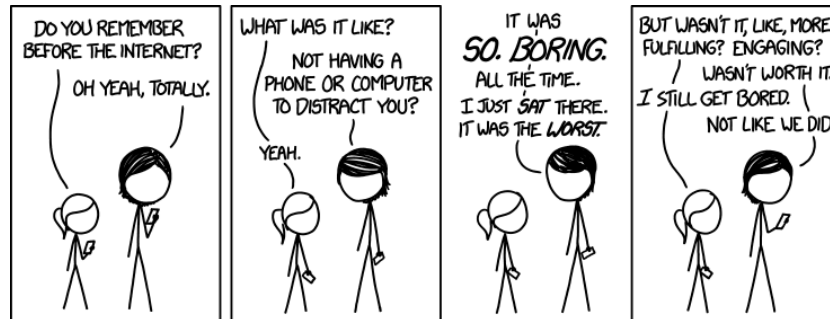


- › **1980:** I.A. Tjomsland, Fourth IEEE Symposium on Mass Storage Systems
 - „Parkinson’s 1st Law paraphrased: Data expands to fill the space available.“
 - „The penalties for storing obsolete data are less apparent than are the penalties for discarding potentially useful data.“
- › **1986:** Hal B. Becker, Can users really absorb data at today’s rates?
 - „The recoding density achieved by Gutenberg was approximately 500 symbols per cubic inch – 500 times the density of [4,000 B.C. Sumerian] clay tablets. By the year 2000, semiconductor random access memory should be storing 1.25×10^{11} bytes per cubic inch.“ – po pravdě o řád přestřelené v roce 2017
- › **1996:** B.J. Truskowski, The Evolution of Storage Systems
 - Digital storage becomes more cost-effective for storing data than paper.
- › **1997** Michael Cox and David Ellsworth
 - „Data sets are generally quite large, taxing the capacities of main memory, local disk, and even remote disk. We call this the problem of *big data*. When data sets do not fit in main memory (*in core*), or when they do not fit even on local disk, the most common solution is to acquire more resources. “

A Very Short History Of Big Data

Gil Press pro [forbes](#)

- › **1997** Michael Lesk, How much information is there in the world?
 - „In only a few years, (a) we will be able [to] save everything—no information will have to be thrown out, and (b) the typical piece of information will never be looked at by a human being.“
- › **1998:** K.G. Coffman and Andrew Odlyzko
 - „ the growth rate of traffic on the public Internet, while lower than is often cited, is still about 100% per year, much higher than for traffic on other networks.“

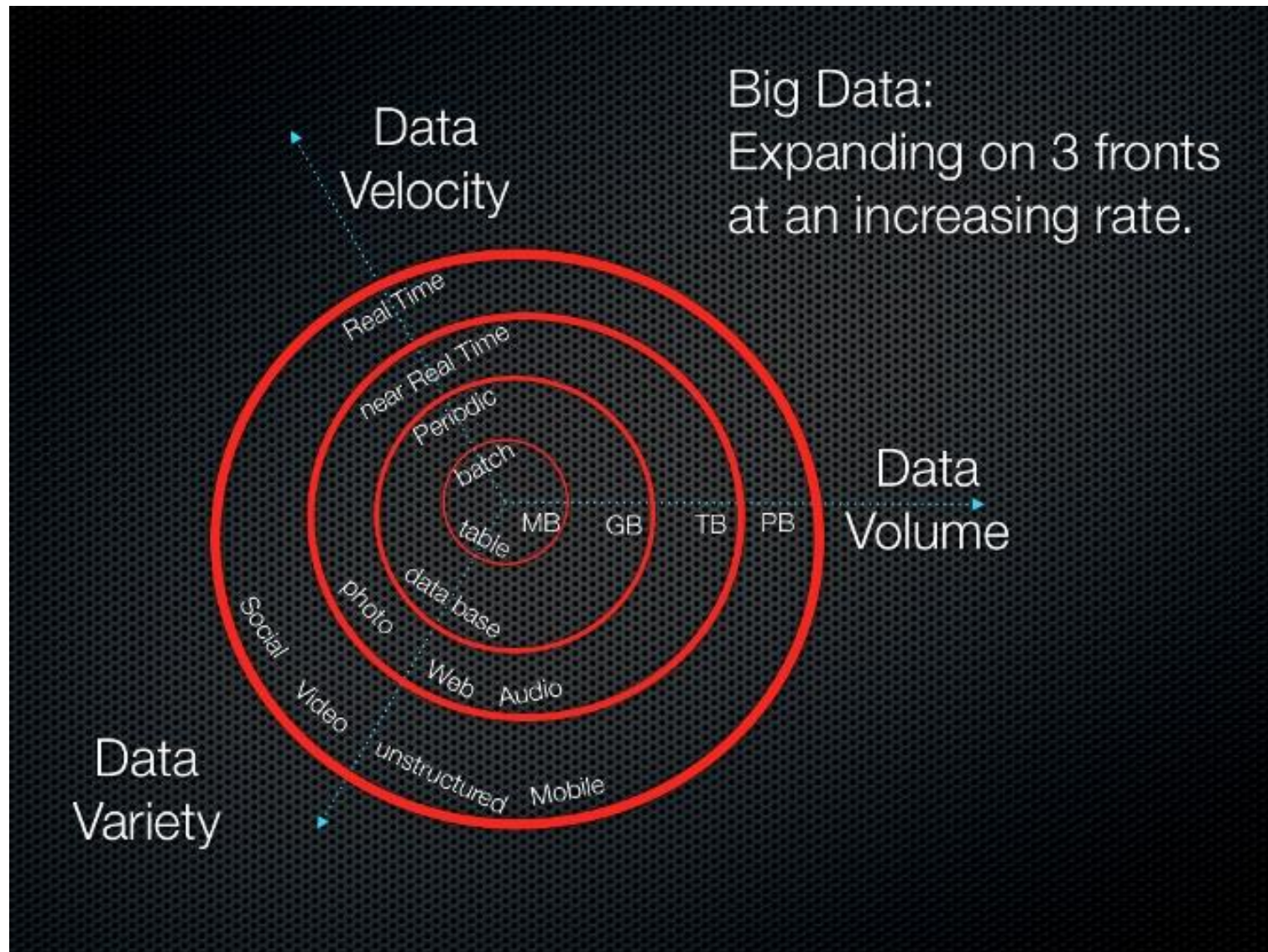


- › **2000: Big Data Era**
 - „the world produced about 1.5 exabytes of unique information, or about 250 megabytes for every man, woman, and child on earth. It also finds that “a vast amount of unique information is created and stored by individuals” (what it calls the “democratization of data”) and that “not only is digital information production the largest in total, it is also the most rapidly growing.“

A Very Short History Of Big Data – Shrnutí

- › Množství dat exponenciálně roste, stejně jako přenosová kapacita
- › Nabídka dat roste rychleji než poptávka
- › V komunikaci převládá obousměrnost a aktivní role uživatelů
- › Je levnější data skladovat než je třídit a vyhazovat
- › Opačná perspektiva: roste úložná kapacita a data ji jen celou zaplňují – poprvé můžeme uložit (skoro) všechny informace.
- › Datové soubory překračují hranici jednoho zařízení / média
- › Průměrná informace není nikdy přečtena člověkem

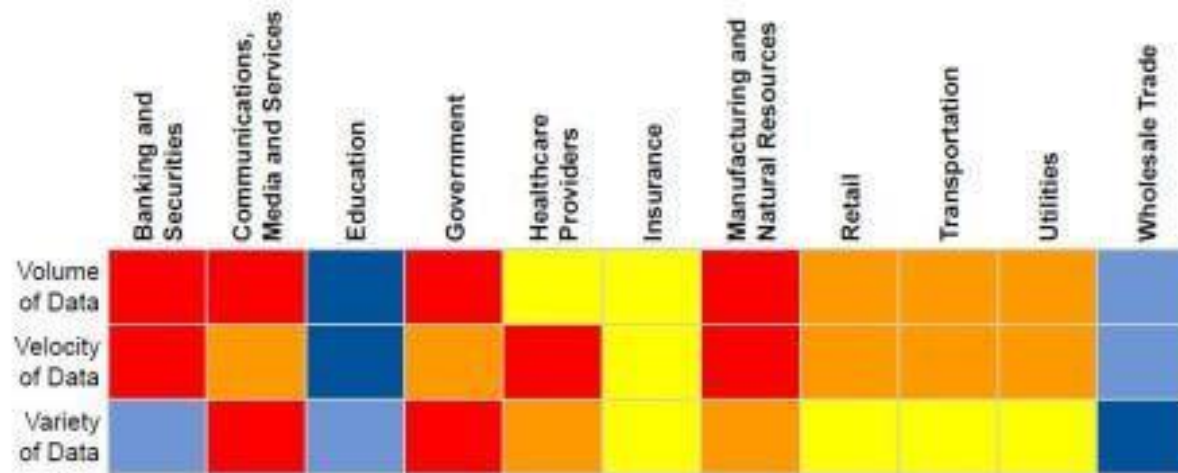
Big Data Definice – 3V



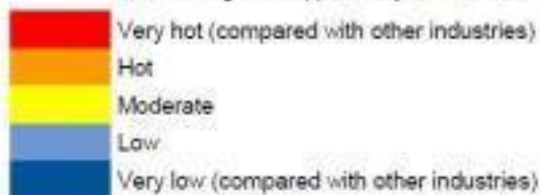
Big Data Definice – 3V

- › Volume, Velocity, Variety

Comparison of Data Characteristics by Industry



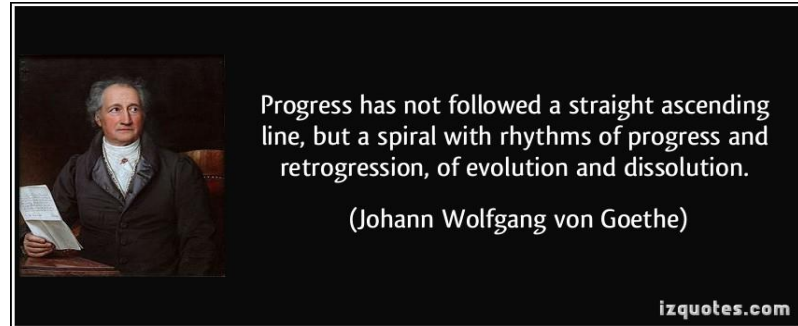
Potential big data opportunity on each dimension is:



The background of the slide is composed of numerous overlapping, translucent, light-gray geometric shapes. These shapes, which include rectangles, trapezoids, and irregular polygons, are arranged in a way that creates a sense of depth and movement, as if they are floating or falling from the top right towards the bottom left. The overall effect is a complex, layered pattern that provides a modern and abstract visual context for the text.

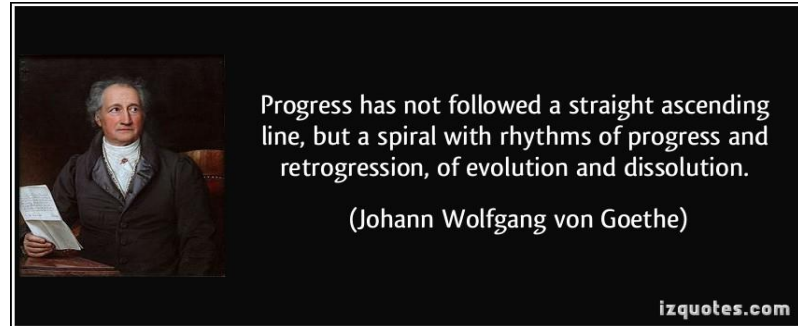
Databáze?

Prehistorie pojmu Databáze



- › Původně – datová základna, slovník faktů
- › **40. léta – Ad hoc proudy dat**
 - Děrné štítky a pásky
- › **50. léta – File System**
 - Soubory a složky s hierarchickou strukturou a unikátní cestou
 - Nezávislost na použitém médiu (páska, disk, ...)
- › **60. léta – DBMS**
 - Tabulky a Indexování: hashování, B-stromy. Klient-server architektura
- › **70. léta – R-DBMS**
 - Relační paradigma: normální formy, relační algebra, selekce, projekce atd.

Historie pojmu Databáze



- › **80. léta – SQL**
 - Jednotný dotazovací jazyk, rozvoj proprietálních databází
- › **90. léta – Datové sklady**
 - Datová Integrace, jedna pravda, analytický reporting.
- › **0. léta – NoSQL**
 - Webové programování, rozvolňování normálních forem, grafové databáze, cloud – částečný návrat ke vzdáleným architekturám
- › **10. léta – Big Data**

RDMS vs Big Data

	Relační DB	Hadoop
Velikost	GB, TB	TB, EB, PB
Přístup	interaktivní i batch	jen batch
Dotazy	SQL a program. návody	Map-Reduce a SQL emulace
Úpravy	opakované čtení i zápis	zápis jednou, čtení opakovaně
Struktura	statické databázové schéma	dynamické schéma – volba při načtení
Integrita	ACID	není, konzistence pomocí redundance
Výkon	limitovaný shora optimalizace na straně DB	lineární škálování optimalizace přidáváním výkonu
Latence	minimální (ms)	značná (desítky sekund i více)
HW	špičkově vyladěné stroje	běžný hardware
Licence	komerční, velmi drahé	open source + podpora
Paralelizace	limitovaná a extrémě drahá cena per jádro	základní kámen architektury

The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray. They are arranged in a way that creates a sense of depth and three-dimensional space, with some shapes appearing to float above others. The overall effect is a modern, architectural, and somewhat crystalline aesthetic.

Programování?

Programování – stručná historie aneb stále větší abstrakce



- › Procedurální – Imperativní programování
 - Programátor tvoří program tím, že přímo specifikuje posloupnost příkazů
- › **Nestrukturované paradigma** (cca do 60.let)
 - Asemblerové instrukce a podmíněné skoky – goto era
- › **Strukturované paradigma** (cca 60.-80. léta)
 - Ústřední prvek programu jsou funkce sdružené do knihoven
- › **Objektové paradigma** (cca 90. léta)
 - Propojení funkcí a dat do objektů, dědičnost a polymorfismus
- › **Virtualizace** (cca 0. léta)
 - Oddělení programátora od specifického OS a HW
- › Nosná myšlenka
 - Zvyšování úrovně abstrakce (kompilátor, linker, vm)
 - Odstraňování complexity ležící dole – knihovny obstarají nízkoúrovňové úlohy
- › Problém: komplexita neleží jen dole, ale i nahoře.

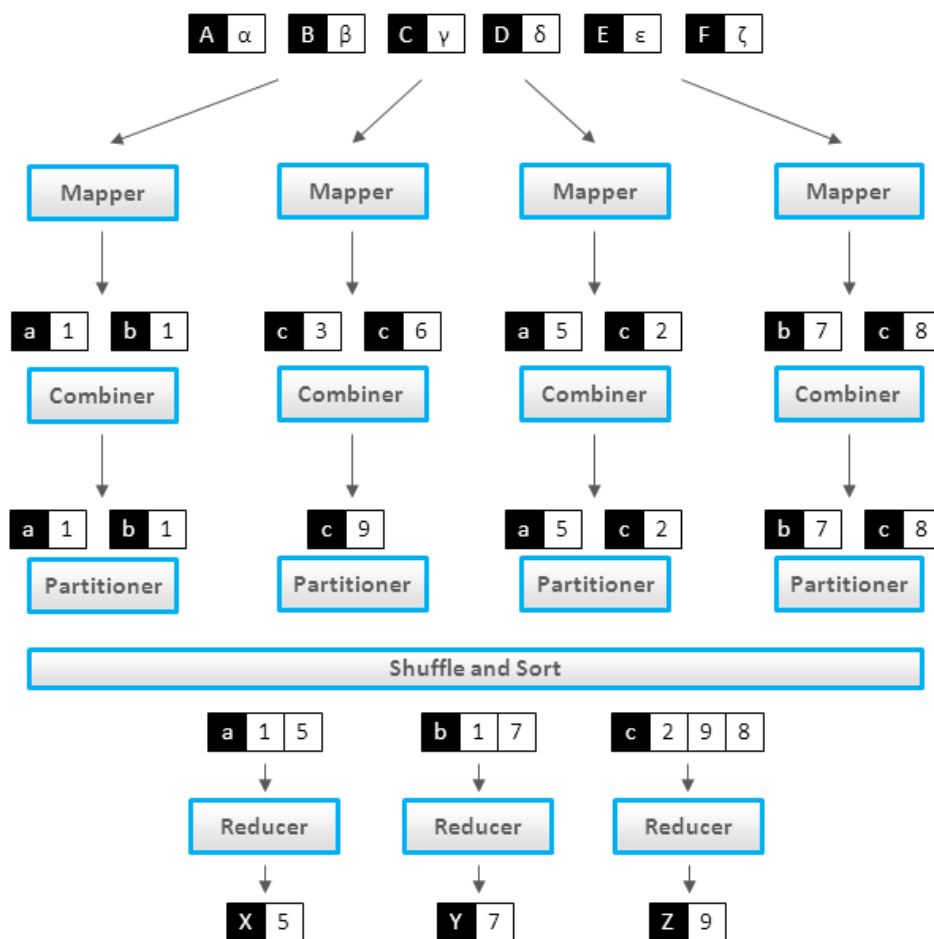
Paralelní framework



- › Problém paralelního programu – komplexita přichází svrchu
 - Je potřeba formulovat algoritmus rovnou jako paralelní a to je velmi těžké
 - Ještě těžší je program udržovat, debugovat, optimalizovat a ladit
- › **Rozvoj frameworků** (současnost)
 - Framework je programová konstrukce popisující hrubou strukturu algoritmu
 - Například quick-sort s uživatelským komparátorem, struktura rozděl a panuj,...
 - Programátor tvoří program tím, že vkládá vlastní kód na předem připravená místa – obvykle se využívá dědičnosti, šablonování, lambda funkcí atd.
 - Odstínění programátora od komplexity problému, zejména paralelizace
 - Událostmi řízené programování, webové frameworky (apache tomcat), Tensorflow pro neuronové sítě, Map-Reduce paradigma atd.
- › Paralelní knihovny (MPI, BSP, OpenMP ...)
 - Univerzální frameworky mají plnou sílu, ale neřeší komplexitu algoritmu
- › **Map-Reduce** paradigma
 - Limitovaný rozsah funkcionality na zpracování dat
 - výměnou za jednoduchou funkcionální strukturu

Map-Reduce

- › Schéma algoritmu dáno
- › Doplnujeme implementaci metod
- › Map
 - sestaví páry <key,val>
- › Combine*
 - lze sloučit páry se stejným klíčem pro omezení síť. přenosu
- › Shuffle and Sort
 - Předun dat mezi uzly
- › Reduce
 - Agregace výsledků
- › Funkce mohou být obecně komplexní
- › Algoritmus = zřetězení Map a Reduce fází



The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray and white. They are arranged in a way that creates a sense of depth and three-dimensional space, with some shapes appearing to float above others. The overall effect is a modern, architectural, and somewhat crystalline aesthetic.

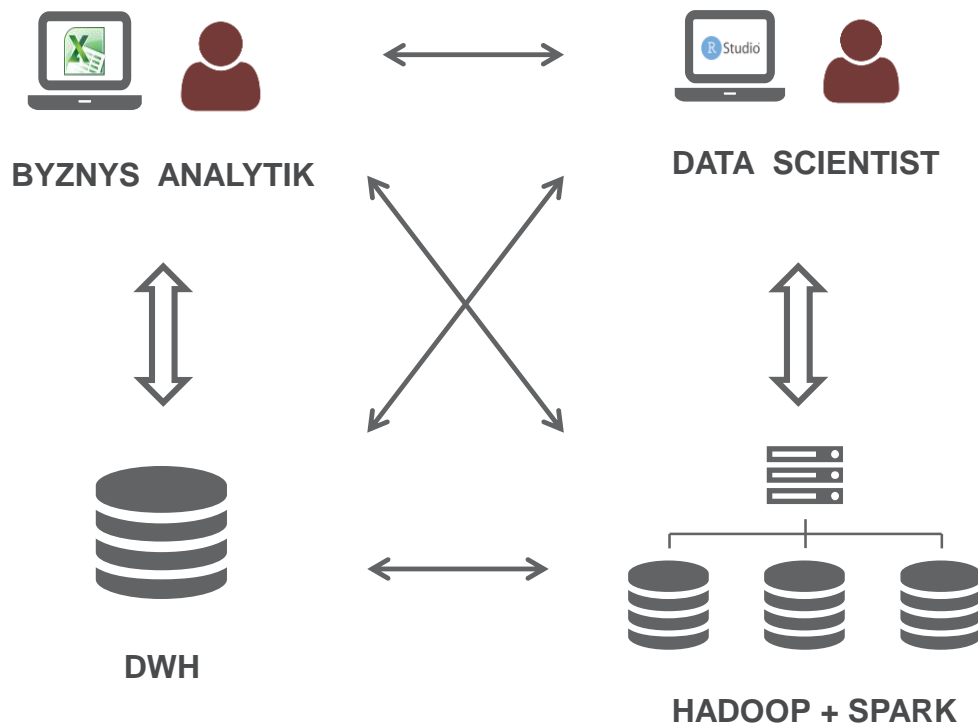
Data Science?

Big Data a Data Science

DWH se má k Business Intelligence

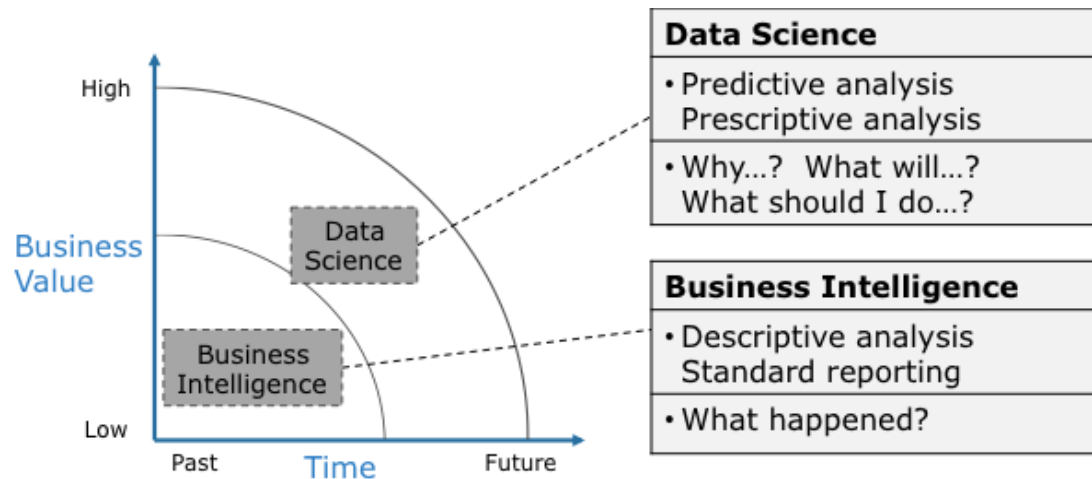
jako

Big Data k Data Science

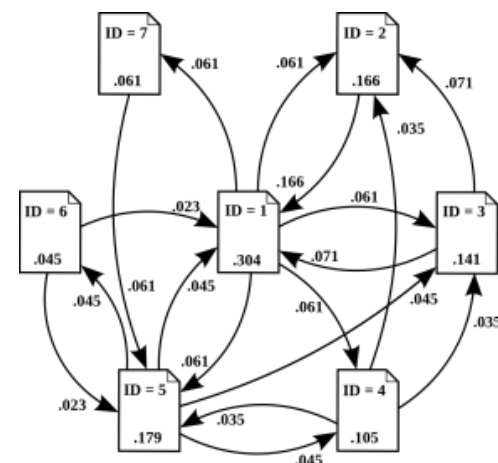


Data Science

- › Spojení oborů
 - Statistika, Informatika, Data mining, Strojové učení, Umělá inteligence
- › Rozdíl oproti Business Intelligence
 - BI: kolik tužek prodaly jednotlivé pobočky v září?
 - DS: kolik jich prodají v říjnu?
- › Klíčové kritérium je práce s nejistotou, pravděpodobnostní výsledek
 - Prediktivní modelování
 - Segmentace, shlukování
 - Podobnostní modelování, kolaborativní filtrování, doporučovací systémy
 - Detekce anomalií
 - Text-mining,
 - Web-mining,
 - Image processing,
 - SNA,
 - atd.



- › Google není vyhledávač
- › Google je řadič
- › 1G stránek (webů) s řádově 10x více odkazy
- › Fulltextové vyhledávání
 - Roboti procházejí web a indexují obsah – inženýrská úloha
 - Uživatel zadá dotaz (slovo, frázi) a výsledkem je seznam stránek
- › V jakém pořadí zobrazit výsledky vyhledávání
 - Většinou je hledaný výraz hned ten první nebo na první stránce
 - Jak je to možné?
- › Google Pagerank
 - Iterativní metoda pro ohodnocení stránek
 - Stránky jsou uzly, odkazy jsou hrany
 - Markovovské řetězce s okrajovými podmínkami
 - Násobení řídkých matic o hraně 1G
 - Implementace ve sparku ;-)



Amazon, Netflix, YouTube

- › Doporučování obsahu
 - Viděl jste těchto deset filmů, podívejte se na jedenáctý
 - Kdo si koupil Babičku, ten si koupí Broučky

› Dva přístupy

- Podobnost zboží
- Podobnost zákazníků

› Kolaborativní filtrování

- Singular Value Decomposition
- Alternating Least Squares (Spark)

› A – matice klient x produkt

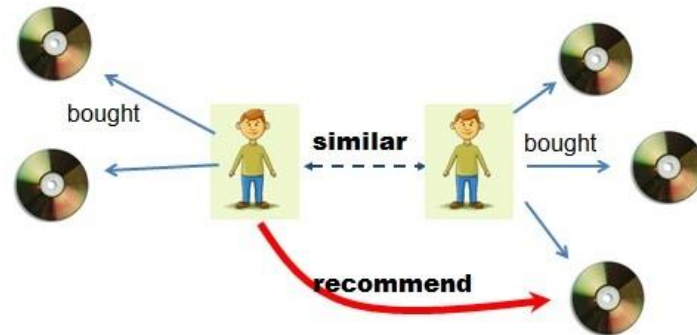
› U – matice klient x faktor

› L – matice faktorů

› V – matice faktor x produkt

› Opětovné vynásobení

- Doporučení chybějících

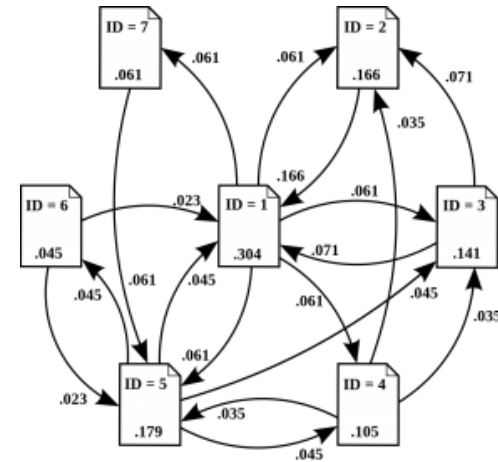


$$\begin{array}{c} \mathbf{A} \end{array} = \begin{array}{c} \mathbf{U} \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array} \begin{array}{c} \mathbf{L} \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array} \begin{array}{c} \mathbf{V}^T \\ \begin{array}{|c|} \hline \text{gray bar} \\ \hline \end{array} \end{array}$$

Big Data Science

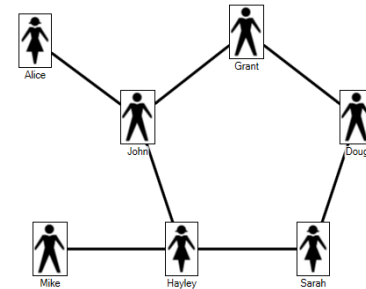
› Google

- Google není vyhledávač, Google je řadič
- 1G stránek s řádově více odkazy
- V jakém pořadí zobrazit výsledky vyhledávání
- Big Data Algoritmus PageRank
 - Hledání vlastních vektorů velké matice



› Amazon, Netflix, YouTube

- Kdo si koupil Babičku, ten si koupí Broučky
- Big Data Algoritmus Kolaborativní filtrování
 - Singular Value Decomposition



› Facebook

- Komu zobrazit jaký obsah
- Kombinace
 - SNA – přátelé
 - Kolaborativní filtrování – like

› A co banky?

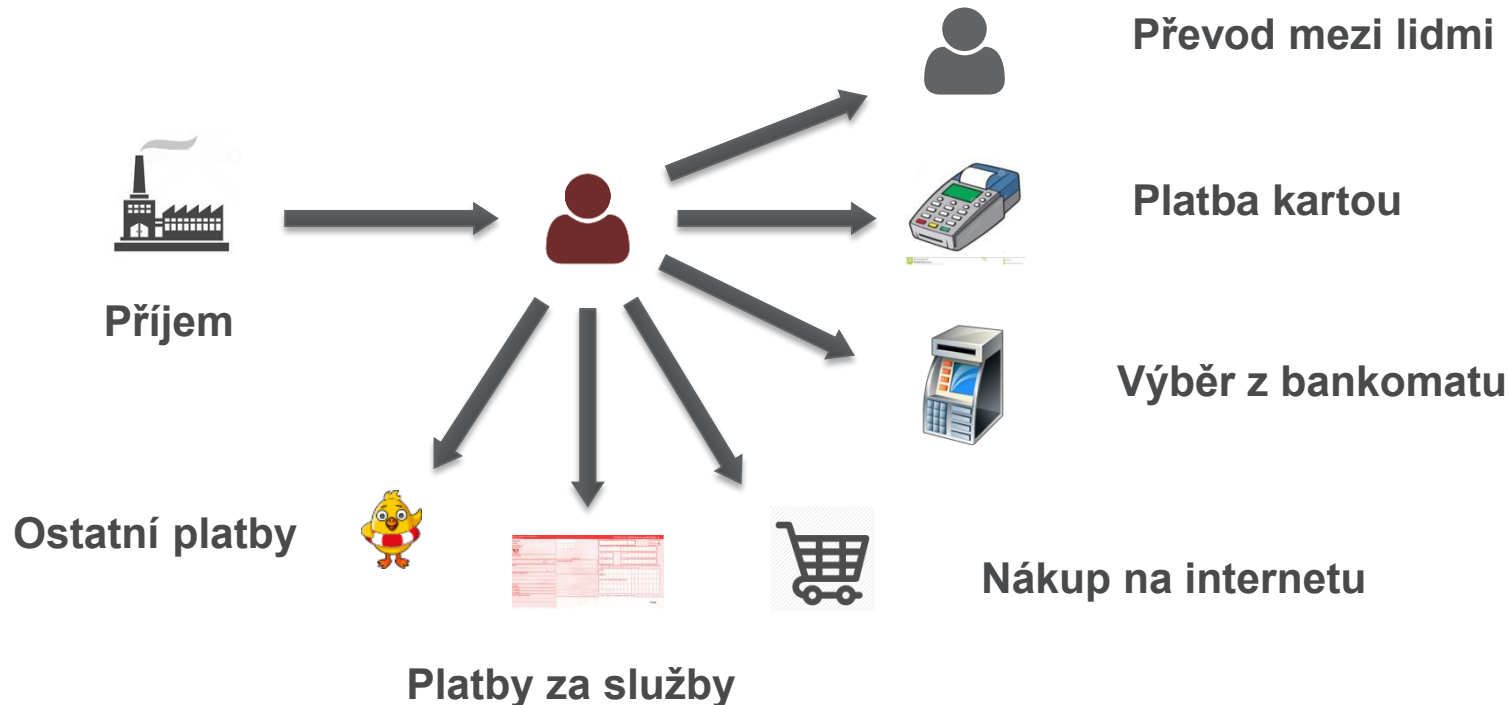
$$\mathbf{A} = \mathbf{U} \mathbf{L} \mathbf{V}^T$$

The background of the slide is a complex, abstract composition of numerous overlapping, translucent geometric shapes. These shapes, which include various polygons and rectangular prisms, are rendered in different shades of light gray. They are arranged in a way that creates a sense of depth and movement, with some shapes appearing to float above others. The overall effect is a textured, crystalline surface that changes as the viewer's perspective shifts.

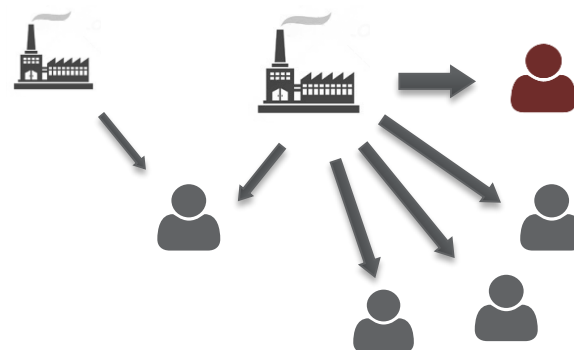
Co děláme my?

Analýza finančních transakcí pomocí BD

- › Vytváříme vyladěné modely pro retailové banky
- › Vstup – finanční transakce
- › Výstup – využitelné informace o klientovi, příznaky, události,
- › Cílem je obohatit stávající obchodní proces o novou znalost



Salary detector



- › Vstup
 - Finanční transakce typu firma - klient
- › Výstup: Identifikované vztahy zaměstnavatel – zaměstnanec
- › Business case
 - Rizikové skóre, detekce událostí, podobnosti (c2c/b2b),...
- › Principy
 - Detekce transakčních vzorců, text mining, pokročilá statistika
- › Vysoká přesnost i pro
 - Krátké úvazky – délka nepřesahující 3 měsíce
 - Nestandardní úvazky (částečné úvazky, práce na živnost, atd.)
 - Firmy s malým počtem zaměstnanců

Detekce domácnosti – Banka/Telco

› Vstup

- Klientské transakce – banka (c2c, karetní operace,...)
- Informace ze sítě – telco (cdr, lokace, billing)
- Základní demografie (věk, pohlaví, adresa, příjmení,...)

› Výstup

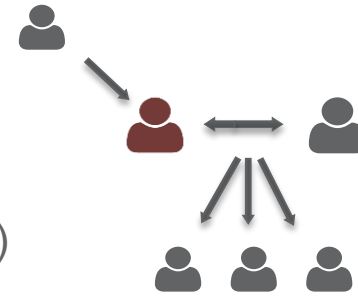
- Identifikace členů domácnosti a rodinných vztahů

› Obchodní využití

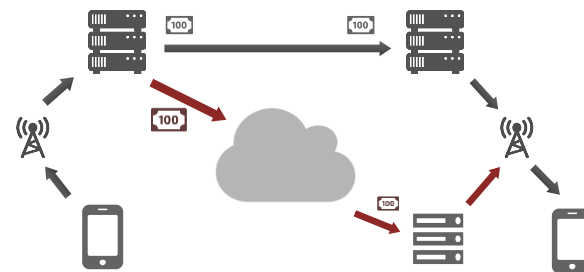
- Rodinný marketing, robustní rizikové skóre,...

› Principy

- Detekce transakčních vzorců, analýza interakcí, text mining



Telco Big Data SimBox Fraud



› Scénář

- Zahraniční operátor/subjekt obchází standardní mezistátní hovor přes internet s cílem ušetřit na mezinárodním propojovacím poplatku

› Vstup

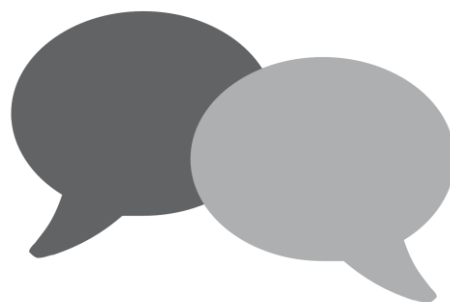
- Telco, síťová data (cdr, location, billing)

› Výstup

- Identifikované podezřelé sim karty

› Principy

- Detekce specifických typů neobvyklého chování
- Rozpoznání skupin s podobným chováním
- Automatická detekce pomocí roamingových dat



Dotazy