# Lec 1

5 сентября 2024 г.     18:24
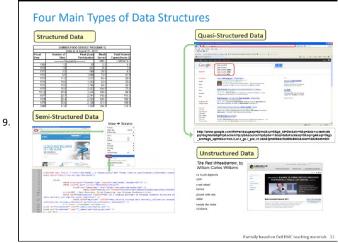
Notes:
1. Big data not new. It's just raw data of different scale and variety
2. 4V: volume(storage problems), velocity(real time analys challenge), variety(very different data), varacity( quality data challenge), variability, value
3. Analys models can be hacked by being given wrong data in learning process in order to get specific output on specific dataset.(For example: hacking banking credit analys model to get credit on big money aprooved)
4. Data can be structured and not structured.
5. Structured - database-like. Transactions and etc
6. Unstructured - videos, images, text. Data that should be analysed as a whole.
7. Semi structured - xml, json.
8. Quasi structured data - very exotic data formats. Not full or none documentation at all. Can be encrypted. Exp: clickstream data from web browser, that can be used to optimize site structure.
9. 



10. Data repositories problem. How to gather data? If every analyst have its own data extract, than developed model can be inaccurate on whole data. It's called data islands
11. Data islands - decentralized data marts.
12. Data warehouses - centralized data centers that store gathered and structured data. Can store BIG amount of data, but hard to maintain and expensive in resources. Not everything should be stored despite scale of storage. Resources are limited.
13. Data analysis and storage technologies are developing together
14. Many areas like banking prefer to use easy-interpreted models. Like regressive. Neural networks are hard to predict and rely on in some ways
15. Data lake - exact copy or near-exact copy of data. Unstructured. Organized and classified only when accessed. Convenient for analysis. Allows to use more advanced analysis technics.
16. Server oriented solutions are not good anymore due to petabytes of data that need to be send to data centre.
17. Decentralized storage is more convenient now. Files are stored in parts distributed between several machines.
18. Every machine analyze only its own part. Then all analysis results combined. Batch processing.
19. Query - user-friendly interface and language for non-programmers to access analytics system and get results. Enhance productivity among personnel.