# Big Data Problems

Per Persson

# What is Big Data?

and what's the problem?

Too many bytes — Volume

Too high a rate — Velocity
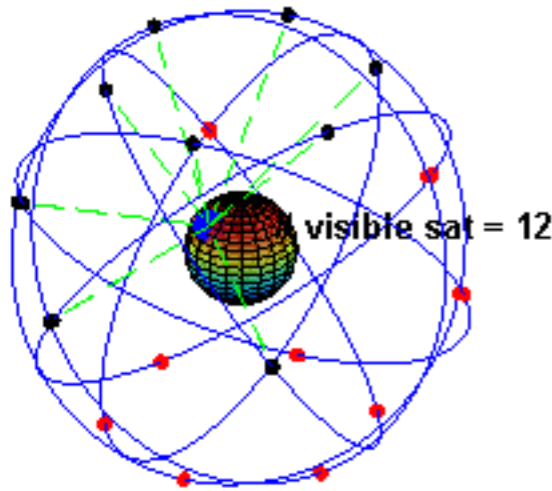
Too many sources — Variety

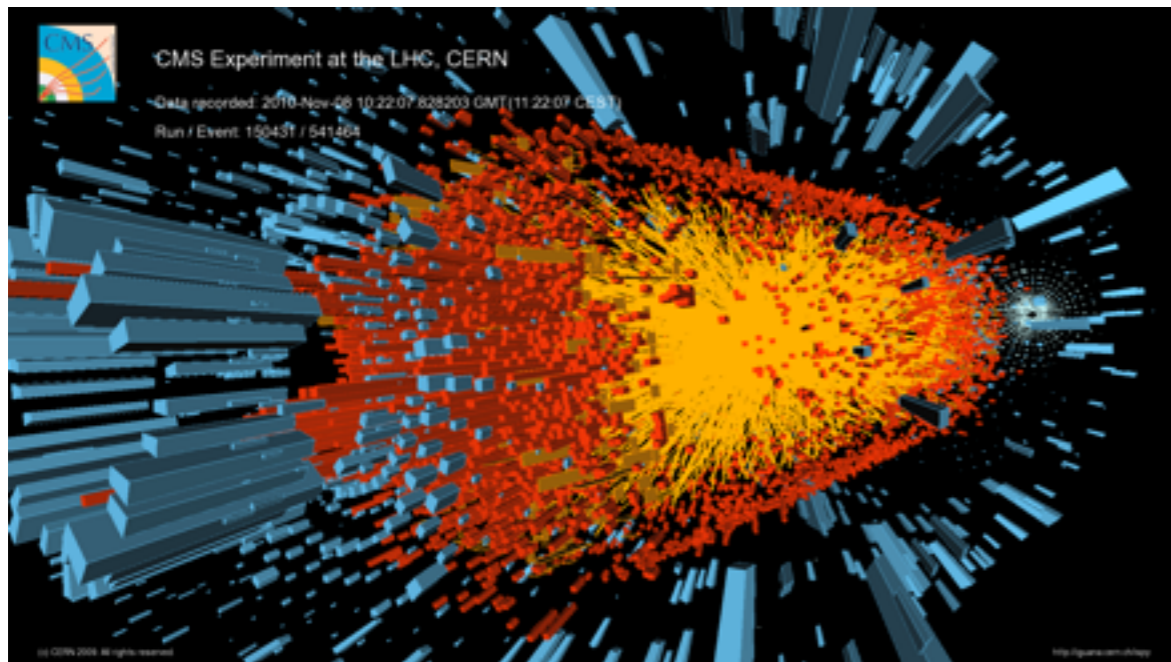Non-scalable analysis, a.k.a simply hard problems

# Volume



visible sat = 12

Too many bytes

GPS data is ≈100 bytes @ 0.1Hz
⇒ 40kB/h ⇒ ≈**1MB/day**

Phones w. GPS: ≈$3 \times 10^9$ to date
⇒ $3 \times 10^9 \times 1MB$ ⇒ ≈**3PB/day**

Moving 3PB @ 1Gbps ⇒ ≈**1 year**

References: Feldman13

# Velocity



The **Large Hadron Collider** (LHC) generates data at a rate of 1PB/s.

Fast electronics selects one in 10000 events in a first pass.

15000 core cluster select 1% of the remaining events for analysis.

Single Tier-0 DC with 73000 cores does reconstruction and storage.

Tier-0 DC distributes data to 11 Tier-1 and 140 Tier-2 DCs.

Continuously 1.5 million jobs and 10GB/s transfer rate globally.
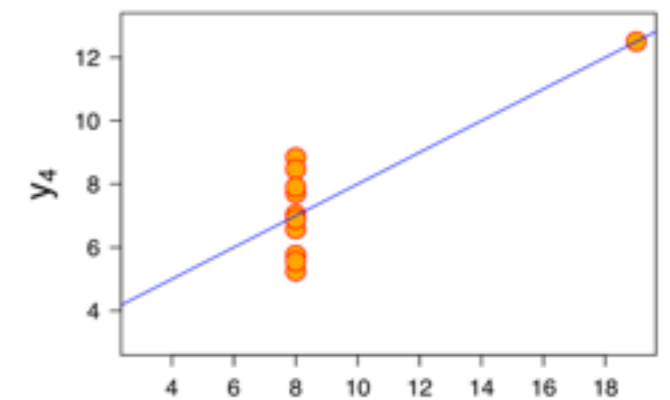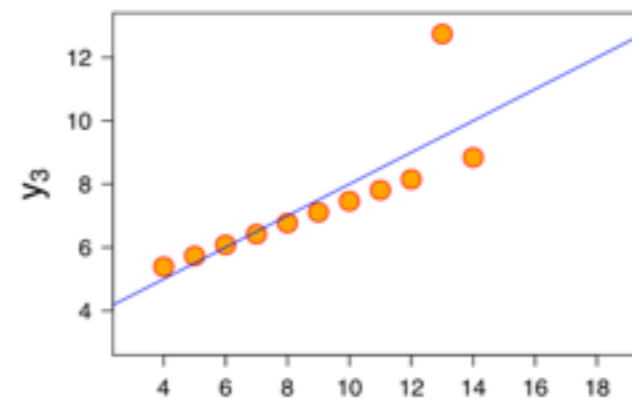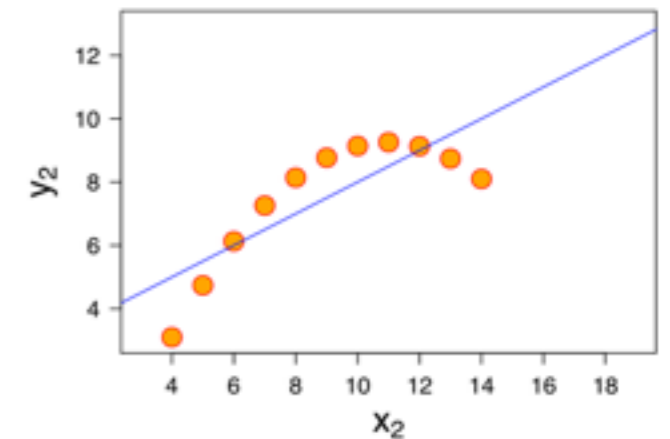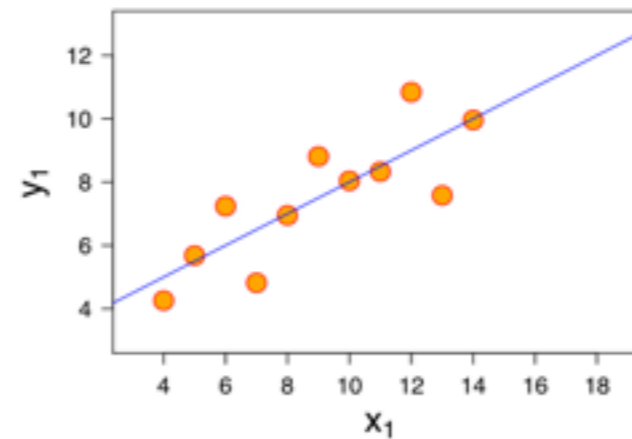
Too high a rate

References: cern13

# Variety

Too many sources

- ETL: Extract, Transform, and Load (labor intensive)
  - OK for up to 10-20 sources, doable up to 30 sources
  - Prohibitively expensive at 50 sources
- Data Curation (with help from tools)
  - Ingest – from an alien source
  - Validate – if bad data gets into your store…, it stays there
  - Transform – align with your schema/ontology
  - Clean – real data is invariably dirty
  - Consolidate – merge it with your previous data
  - Visualize – this is important!
- Example: Novartis
  - Consolidate 8000 electronic lab journals
  - No common schema, no common language, no rules whatsoever…

References: Stonebraker13

# Visualization

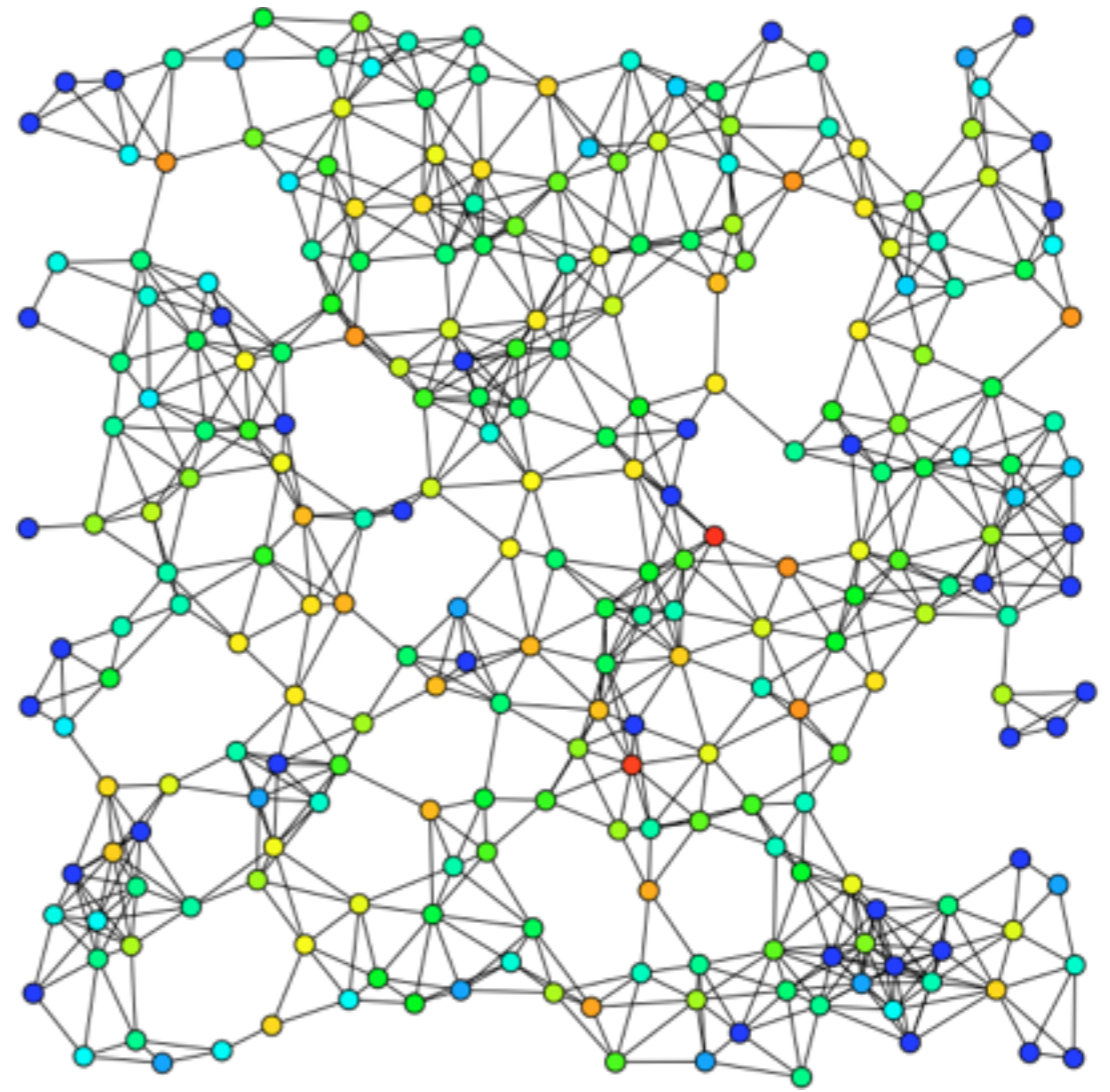| | 1 | | 2 | | 3 | | 4 | |
|---|---|---|---|---|---|---|---|---|
| | x | y | x | y | x | y | x | y |
| | 10,00 | 8,04 | 10,00 | 9,14 | 10,00 | 7,46 | 8,00 | 6,58 |
| | 8,00 | 6,95 | 8,00 | 8,14 | 8,00 | 6,77 | 8,00 | 5,76 |
| | 13,00 | 7,58 | 13,00 | 8,74 | 13,00 | 12,74 | 8,00 | 7,71 |
| | 9,00 | 8,81 | 9,00 | 8,77 | 9,00 | 7,11 | 8,00 | 8,84 |
| | 11,00 | 8,33 | 11,00 | 9,26 | 11,00 | 7,81 | 8,00 | 8,47 |
| | 14,00 | 9,96 | 14,00 | 8,10 | 14,00 | 8,84 | 8,00 | 7,04 |
| | 6,00 | 7,24 | 6,00 | 6,13 | 6,00 | 6,08 | 8,00 | 5,25 |
| | 4,00 | 4,26 | 4,00 | 3,10 | 4,00 | 5,39 | 19,00 | 12,50 |
| | 12,00 | 10,84 | 12,00 | 9,13 | 12,00 | 8,15 | 8,00 | 5,56 |
| | 7,00 | 4,82 | 7,00 | 7,26 | 7,00 | 6,42 | 8,00 | 7,91 |
| | 5,00 | 5,68 | 5,00 | 4,74 | 5,00 | 5,73 | 8,00 | 6,89 |
| Mean | 9,00 | 7,50 | 9,00 | 7,50 | 9,00 | 7,50 | 9,00 | 7,50 |
| Variance | 11,00 | 4,13 | 11,00 | 4,13 | 11,00 | 4,12 | 11,00 | 4,12 |
| Corr. | 0,82 | | 0,82 | | 0,82 | | 0,82 | |
| Lin.reg. | y=3+0.5x | | y=3+0.5x | | y=3+0.5x | | y=3+0.5x | |

Anscombe's quartet



References: Anscombe73

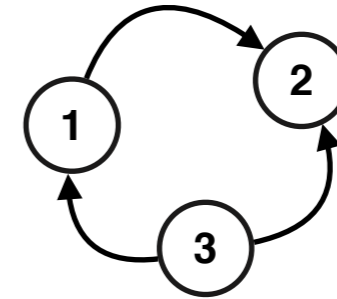# Simply hard problems

- Graph queries

- Collaborative filtering

- k-means clustering

- Logistic regression



References:

# PageRank(ing)

$$\mathbf{A} = \begin{bmatrix} 0 & 0 & 1 \\ 1 & 0 & 1 \\ 0 & 0 & 0 \end{bmatrix}$$

Count in-degree
$$x_i = \sum_j A_{ij} \Rightarrow \mathbf{x} = \begin{bmatrix} 1 \\ 2 \\ 0 \end{bmatrix}$$
In-degree centrality

Weight by node rank
$$x_i' = \sum_j A_{ij} x_j \Leftrightarrow \mathbf{x}' = \mathbf{A}\mathbf{x}$$
Weighted in-degree centrality

Recursively refine
$$\mathbf{x}^{t+1} = \mathbf{A}^t \mathbf{x} \Rightarrow \mathbf{A}\mathbf{x} = \lambda \mathbf{x}$$
Eigencentrality

Fix the contagious 0
$$\mathbf{x} = \alpha \mathbf{A}\mathbf{x} + \beta \mathbf{1} \Rightarrow \mathbf{x} = [\beta = 1] = (\mathbf{1} - \alpha \mathbf{A})^{-1} \mathbf{1} = \begin{bmatrix} 1.5 \\ 2.25 \\ 1 \end{bmatrix}$$
Katz centrality

Compensate for Yahoo!
$$x_i' = \alpha \sum_j A_{ij} \frac{x_j}{k_j^{\text{out}}} + \beta, \quad \text{where} \quad k_j^{\text{out}} = \sum_i A_{ij}$$

Determine $\alpha$
$$\mathbf{x} = \mathbf{D}(\mathbf{D} - \alpha \mathbf{A})^{-1} \mathbf{1}, \quad \text{where} \quad D_{ii} = \max\left(k_i^{\text{out}}, 1\right) \Rightarrow 0 < \alpha < 1$$
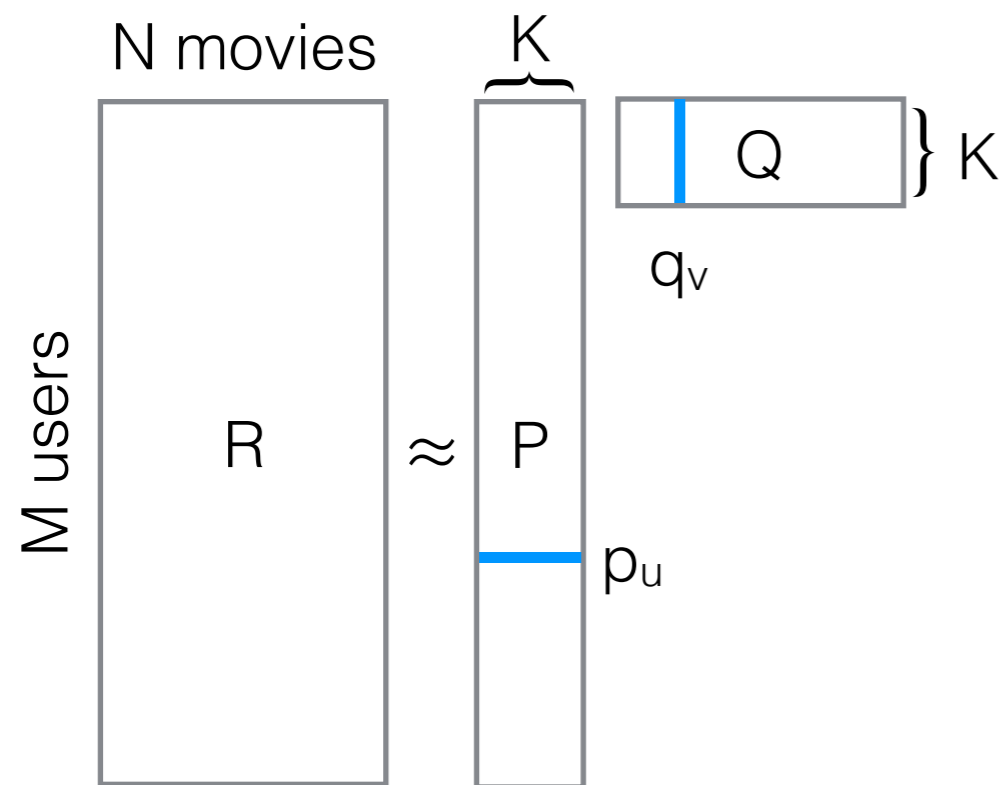
With $\alpha$=0.85
$$\mathbf{x} = \begin{bmatrix} 1.4250 \\ 2.6362 \\ 1.0000 \end{bmatrix}$$
PageRank

References: Page99, Newman10 pp.168-178

# Recommendations

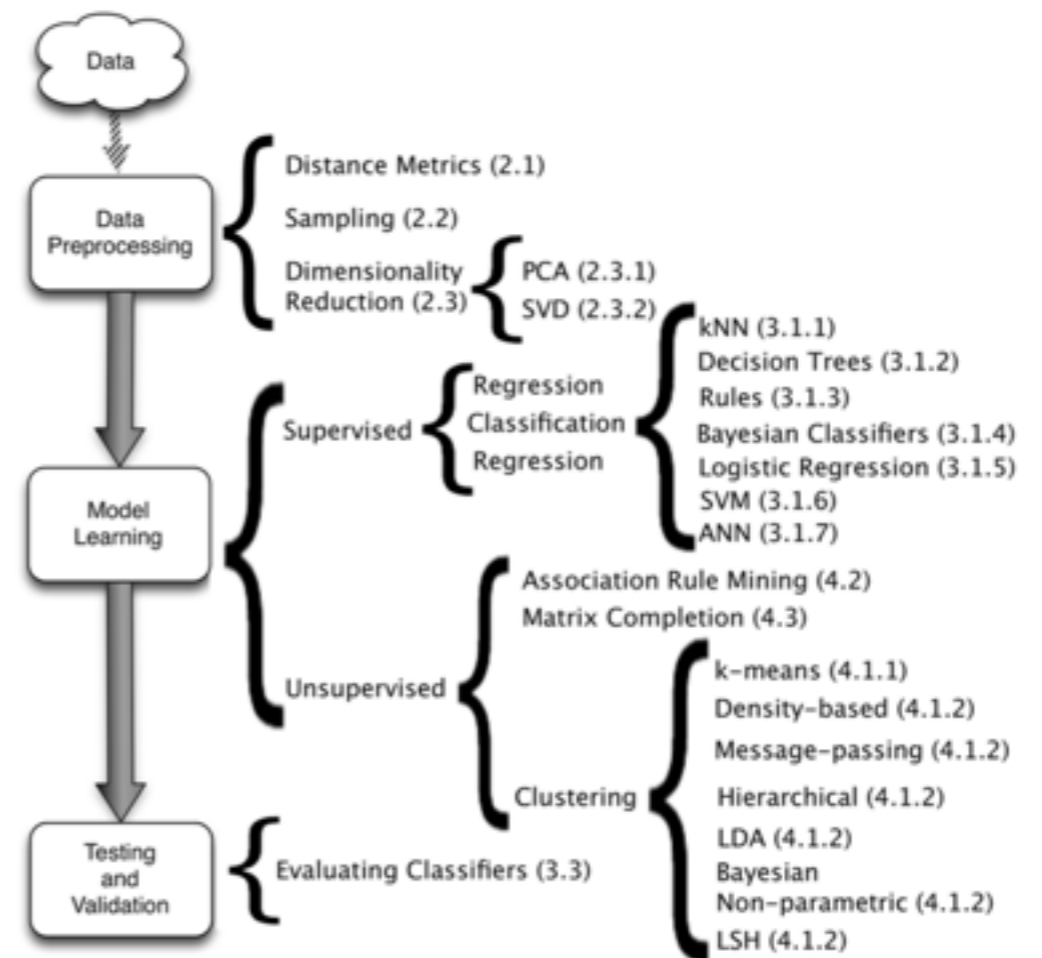## Naïve approach

N movies

$$R \approx P$$



M users

K

Q

} K

$q_v$

$p_u$

R is **sparse** (99.9% empty) matrix of ratings,
M=500000, N=17000
Find **dense** matrices P and Q, such that
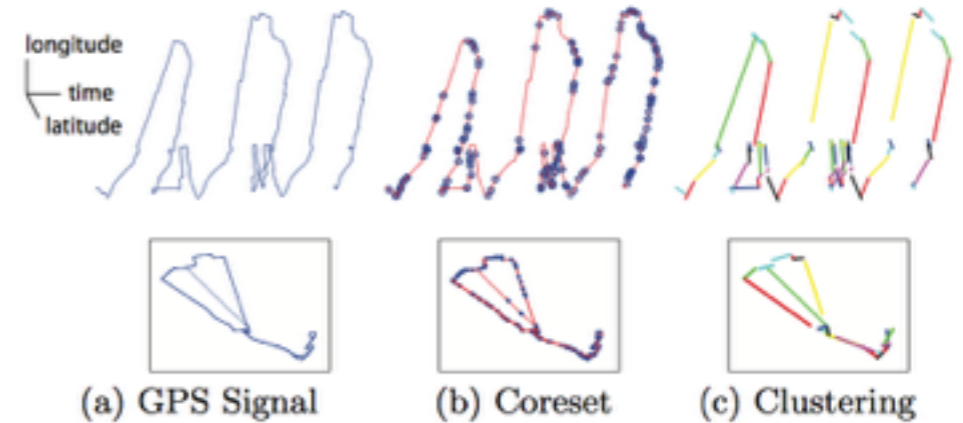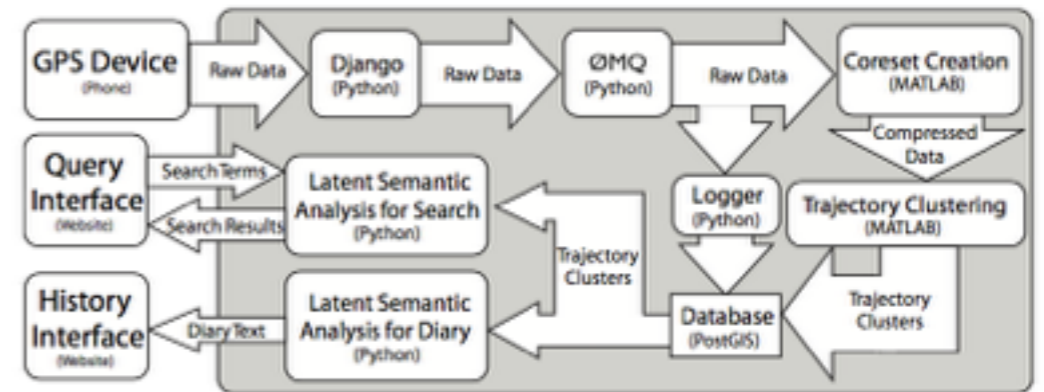$p_u$ & $p_v$ accurately estimates u:s rating of v.

## Realistic approach



Data

Distance Metrics (2.1)
Data Preprocessing
Sampling (2.2)
Dimensionality Reduction (2.3)
PCA (2.3.1)
SVD (2.3.2)

kNN (3.1.1)
Decision Trees (3.1.2)
Regression
Classification
Rules (3.1.3)
Bayesian Classifiers (3.1.4)
Supervised
Logistic Regression (3.1.5)
Regression
SVM (3.1.6)
Model Learning
ANN (3.1.7)

Association Rule Mining (4.2)
Matrix Completion (4.3)

k-means (4.1.1)
Unsupervised
Density-based (4.1.2)
Message-passing (4.1.2)
Clustering
Hierarchical (4.1.2)
LDA (4.1.2)
Testing and Validation
Evaluating Classifiers (3.3)
Bayesian Non-parametric (4.1.2)
LSH (4.1.2)

From: Amatriain et al. "Data Mining Methods for Recommender Systems" in "Recommender Systems Handbook"

References:Satish14, Amatriain14

# Life-logging



- Combines GPS, Yelp, maps, search, and semantic analysis



(a) GPS Signal    (b) Coreset    (c) Clustering





References: Feldman13

# The law of the instrument

I suppose it is tempting, if the only tool you have is a hammer, to treat everything as if it were a nail

*Abraham Maslow*

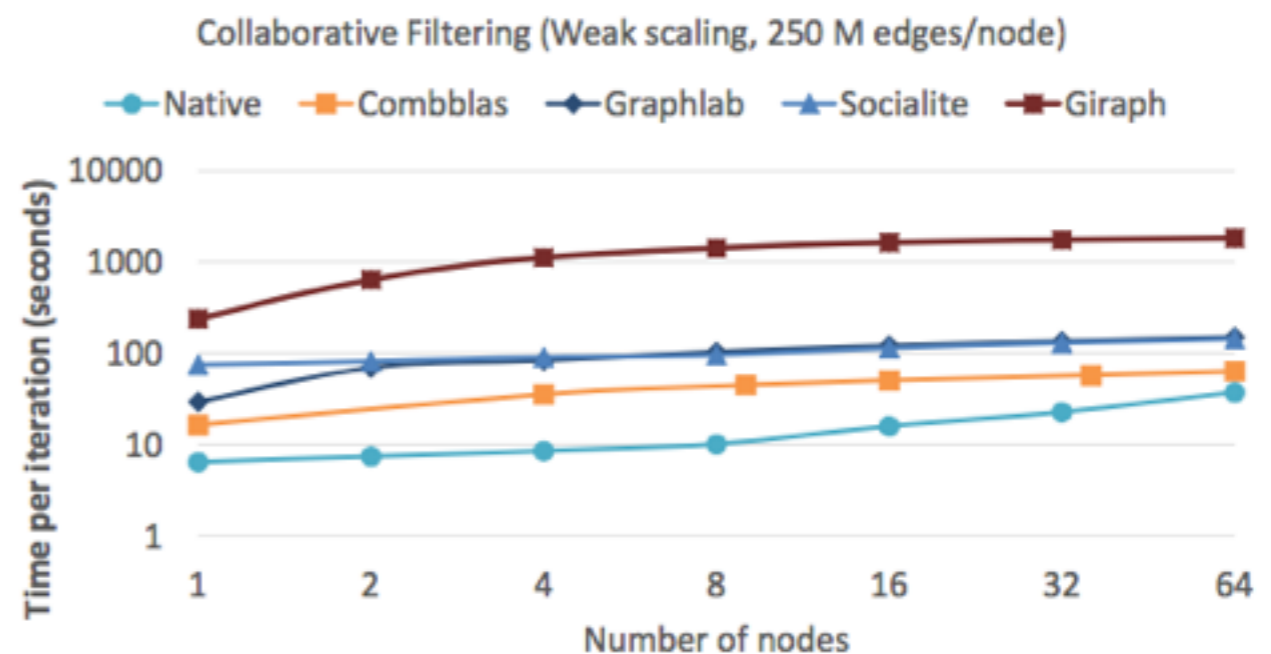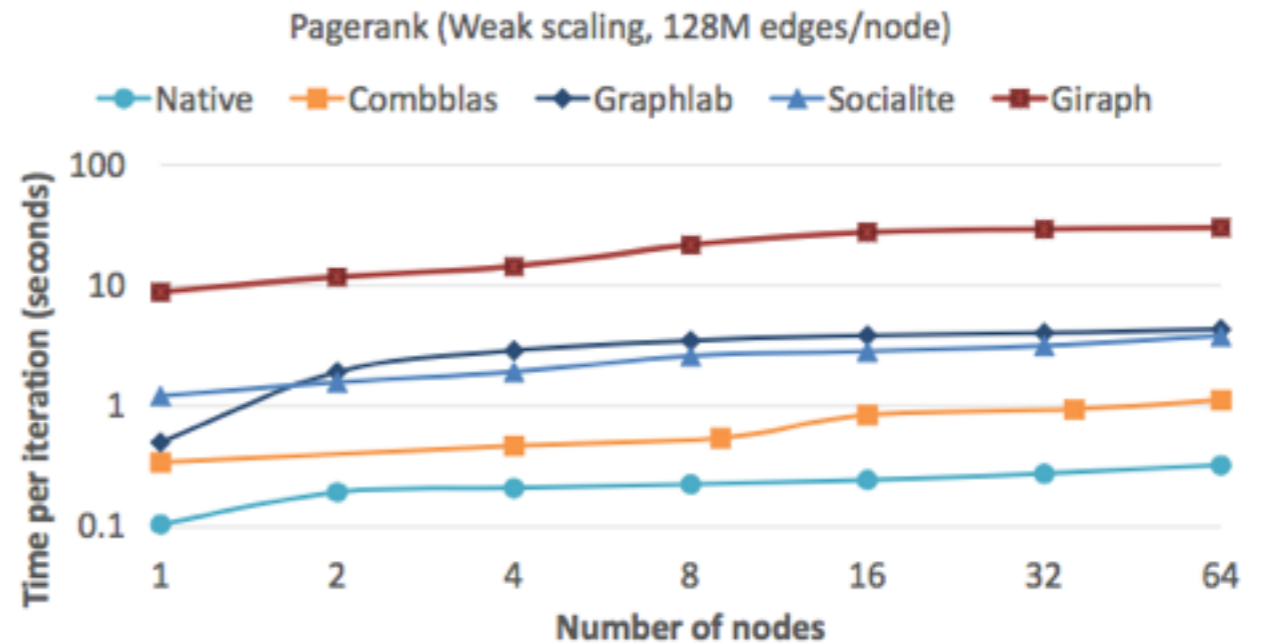- MapReduce for everything?
  - Hive/Pig/…
  - Dryad/Spark/…
  - Pregel/Giraph/…
- Some people actually think so…
  - Benefits of a familiar tool outweighs drawbacks

References: Kaplan64, Maslow66, Lin13

# Speed bumps and Ninjas

- How bad can using a golden hammer be?
- How much can you benefit from a Ninja programmer?
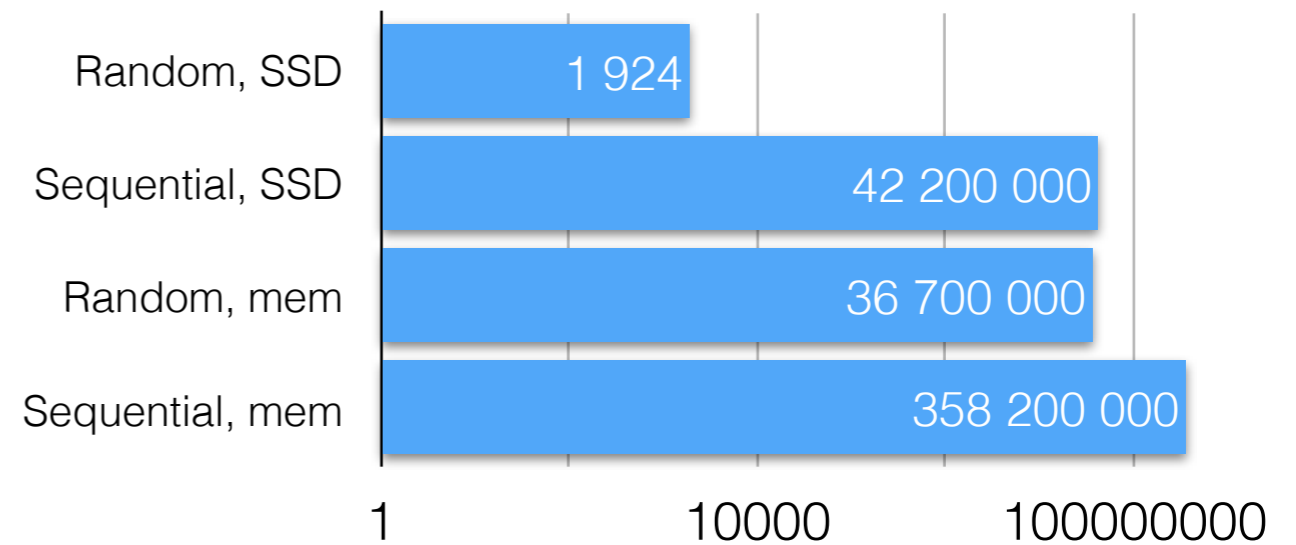- Do you pay your cloud provider as you go?

# Releasing the Ninjas

- How much performance is lost with a hammer?
- Custom code vs. standard graph tools
- Apparently 2x-30x (500x) depending on problem/tool



Pagerank (Weak scaling, 128M edges/node)



Collaborative Filtering (Weak scaling, 250 M edges/node)

References: Satish14

# Speedbumps ahead

- Beware of RDBMS
  - Custom: 15min
  - PostgreSQL: fail after 6h
- Row vs column store
- ≈90% wasted on locks and queue management
- 10-100x faster with column, in-memory, single-threaded, lock-free implementation



| | |
|---|---|
| Random, SSD | 1 924 |
| Sequential, SSD | 42 200 000 |
| Random, mem | 36 700 000 |
| Sequential, mem | 358 200 000 |

Source: Jacobs09



Locking 23 %
Latching 23 %
Useful work 8 %
Recovery 23 %
Buffer pool 23 %

Source: Stonebraker14

References: Jacobs09, Stonebraker14

# Conclusions

- Understand your data (curation, compression)
- Understand your questions (relax, property testing)
- Understand your algorithms
- Understand your tool(s) (cost/benefit analysis)
- Cloud can "hide" inefficiencies in algorithm
  - pay as you go could mean wasting money

# References

[cern13] http://home.web.cern.ch/about/updates/2013/04/animation-shows-lhc-data-processing, 04 2013.

[Amatriain14] Xavier Amatriain, http://www.slideshare.net/xamat/kdd-2014-tutorial-the-recommender-problem-revisited, August 2014.

[Anscombe73] Francis J Anscombe, "Graphs in statistical analysis", The American Statistician 27 (1973), no. 1.

[Brin98] Sergey Brin and Lawrence Page, "The anatomy of a large-scale hypertextual web search engine", Computer networks and ISDN systems 30 (1998), no. 1

[Feldman13] Dan Feldman, et al., "iDiary: From gps signals to a text-searchable diary", Proc. ACM Conference on Embedded Networked Sensor Systems (New York, NY, USA), SenSys '13, ACM, 2013

[Jacobs09] Adam Jacobs, "The pathologies of big data", Commun. ACM 52 (2009), no. 8

[Kaplan64] Abraham Kaplan, "The conduct of inquiry: Methodology for behavioral science", Chandler Publishing Company, 1964.

[Lin13] Jimmy Lin, "MapReduce is good enough? if all you have is a hammer, throw away everything that's not a nail!", Big Data 1 (2013), no. 1

[Maslow66] Abraham Maslow, "The psychology of science", the john dewey society lectureship series, 1966.

[Newman10] Mark E. J. Newman, "Networks: An introduction", Oxford University Press, March 2010.

[Page99] Lawrence Page, et al, "The PageRank citation ranking: Bringing order to the web"., Technical Report 1999-66, Stanford InfoLab, November 1999.

[Satish14] Nadathur Satish, et al, "Navigating the maze of graph analytics frameworks using massive graph datasets", Proc. ACM SIGMOD International Conference on Management of Data (New York, NY, USA), SIGMOD '14, ACM, 2014

[Stonebraker13] Michael Stonebraker, et al, "Data curation at scale: The data tamer system"., CIDR, 2013.

[Stonebraker14] Michael Stonebraker, et al, "Enterprise database applications and the cloud: A difficult road ahead", Proc. IEEE International Conference on Cloud Engineering (Washington, DC, USA), IC2E '14, IEEE Computer Society, 2014