# cloudera®

# Data Science in Action

Kamal Maheshwari, Boulder, CO

[On loan from - Sean Owen| Director, Data Science]

# Data Science Deluge..



Will data science prevent deadly police shootings?

A Data Scientist's Real Job: Storytelling

DATA SCIENTIST: THE SEXIEST JOB OF THE 21ST CENTURY
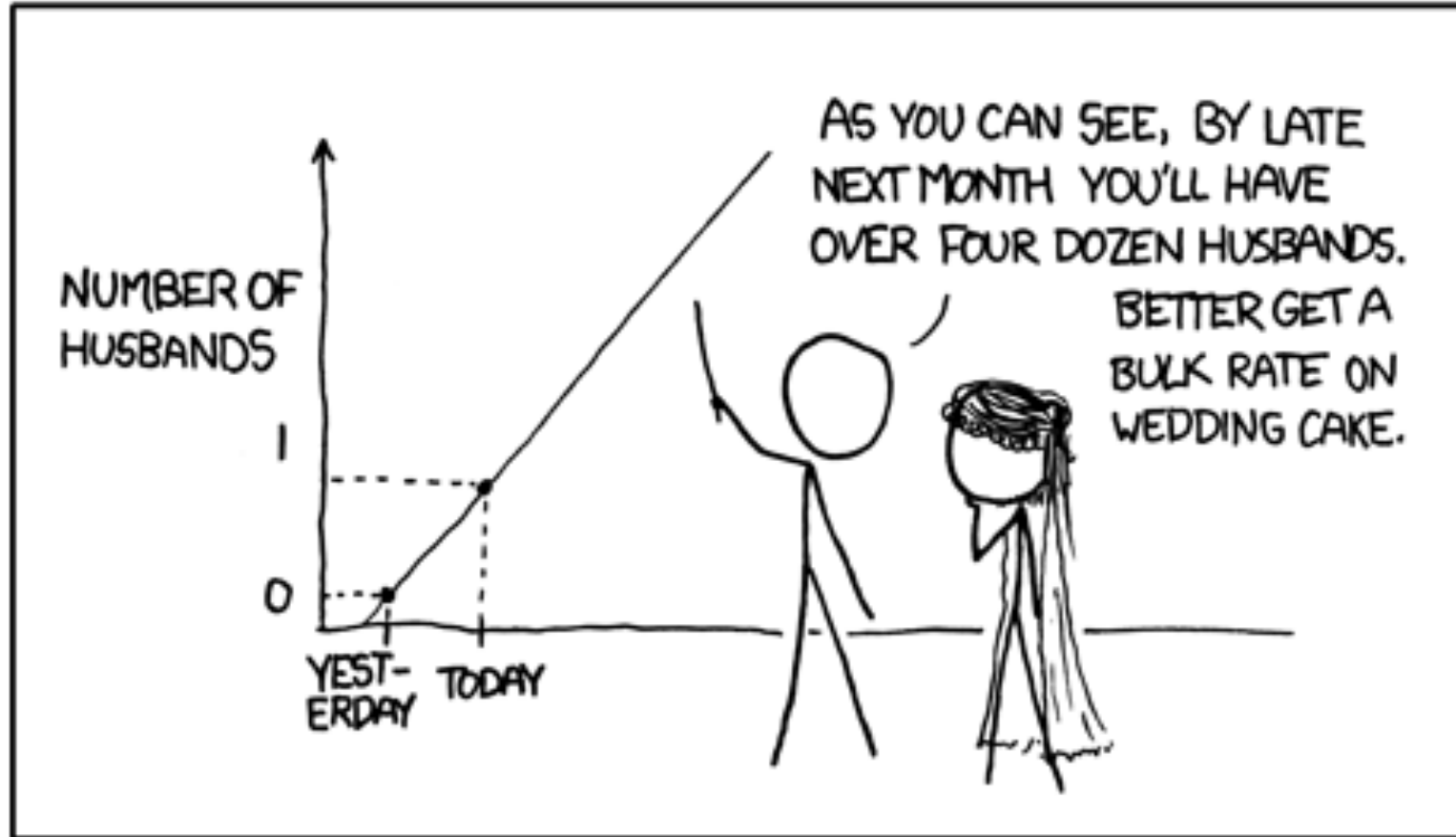
Data Science Is Dead

# Talk Overview

- What is a Data Scientist?

- Introduction to Data Products

- Understanding Data Processing Engines: SQL, MapReduce, and Spark

- Data Modeling for Data Science

- Algorithms Every Data Scientist Should Know

- Executing a Data Product Strategy

# What is a Data Scientist?

cloudera

# One Definition…

# ...and Another...



**Zvi**
@nivertech

"Data Scientist" is a Data Analyst who lives in California.

RETWEETS
**143**

FAVORITES
**44**

6:55 PM - 14 Mar 2012

# …and Another…

# …and Another…

# ...and Another...

Phase 1. Collect Data

Phase 2. Data Science?

Phase 3. Profit!

# …versus Another

**Josh Wills**
@josh_wills

Data Scientist (n.): Person who is better at statistics than any software engineer and better at software engineering than any statistician.

RETWEETS
976

FAVORITES
464

cloudera

# Data Scientist vs. Software Engineer

- Most software engineers do not understand the assumptions behind statistical models
  - Independence
  - Normally distributed errors
- Common errors/pitfalls when conducting data analysis
  - Correlation != causation
  - Simpson's Paradox

# Introduction to Data Products

# What is a Data Product?



- Created from data

- Using the product generates more data

- This data can be used to improve the functionality of the product
  - Ideally in an automated fashion

# Classic Data Products

- Data products aren't a new idea, just a new phrase

- Credit scorecards and rules-based fraud detection engines have been around for a long time

- What has changed?



**Sample Credit Decision Scorecard Report**

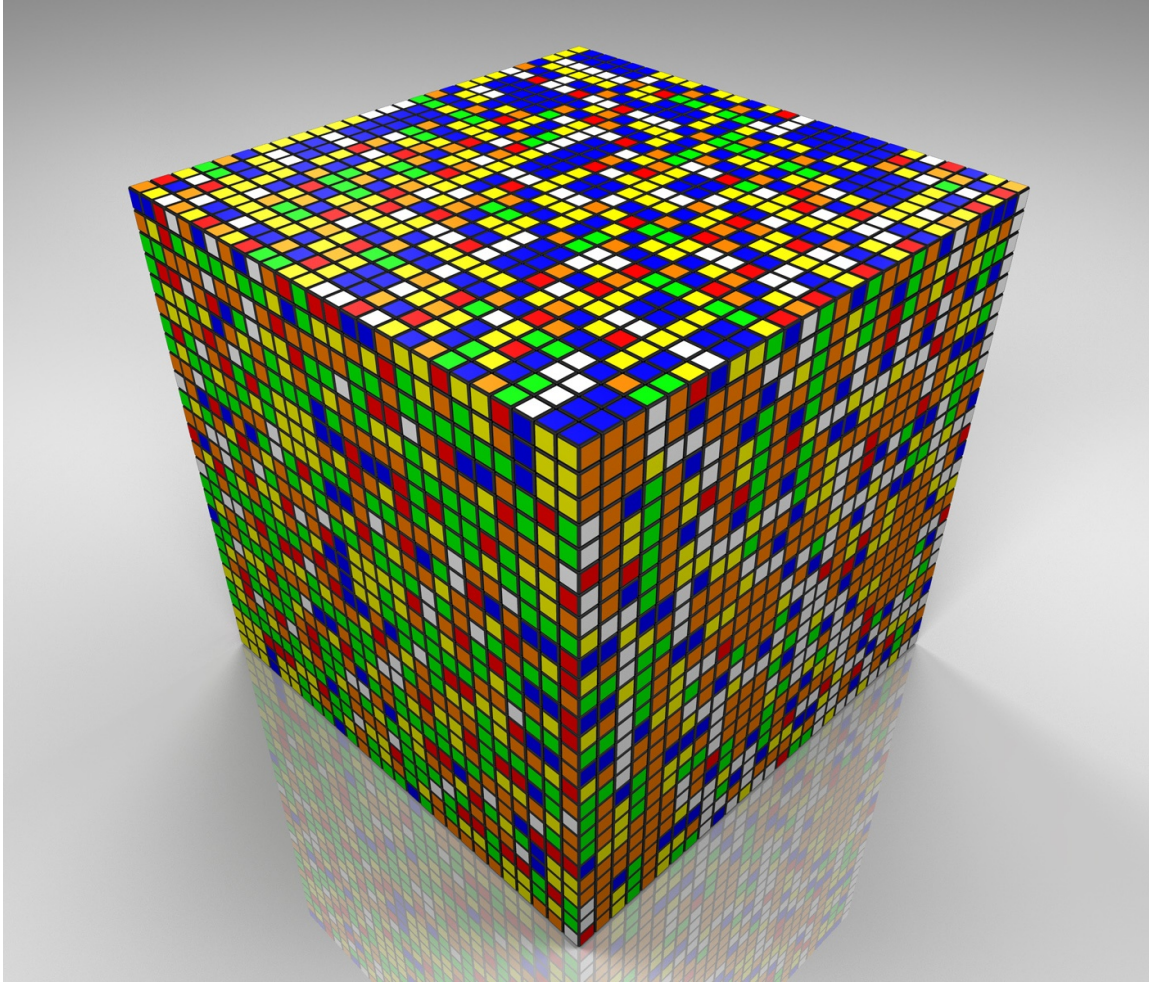| Scorecard | results | recommendation |
|---|---|---|
| PROPOSED RENT | 1000 | N/A |
| STATED MONTHLY INCOME | 4500 | N/A |
| INCOME TO RENT | 4.5:1 | PASS |
| INCOME TO DEBT | 2.75:1 | PASS |
| INCOME TO DEBT INCL RENT | 1.7:1 | PASS |
| CREDIT SCORE    PASS above 700 | FAIL below 600 | CONDITIONAL |
| DELINQUENT ACCOUNTS    FAIL if more than 2 | | PASS |
| COLLECTION/CHARGE OFF    FAIL if more than 2 | | CONDITIONAL |
| BANKRUPTCY RECORDS    FAIL if within 3 years | | PASS |
| OVERALL COMPOSITE: | | CONDITIONAL |

# Big Data and Data Products



- What makes "big data" big?
  - Volume?
  - Variety?
  - Velocity?

- Data becomes big when we take one or more large data sets and start to analyze *relationships* between observations

**cloudera**

# Recommendation Engines

**Customers who viewed this item also viewed these products**



**Dualit Food XL1500 Processor**

$560

🛒 Add to cart

**Kenwood kMix Manual Espresso Machine**

★★★★☆

$250

🔧 Select options

**Weber One Touch Gold Premium Charcoal Grill-57cm**

$225

🛒 Add to cart

**NoMU Salt Pepper and Spice Grinders**

$3

🛒 View options

# Search Engines

# Data Products Are The New Competitive Advantage

# Designing for Iteration

cloudera
Thanks!