

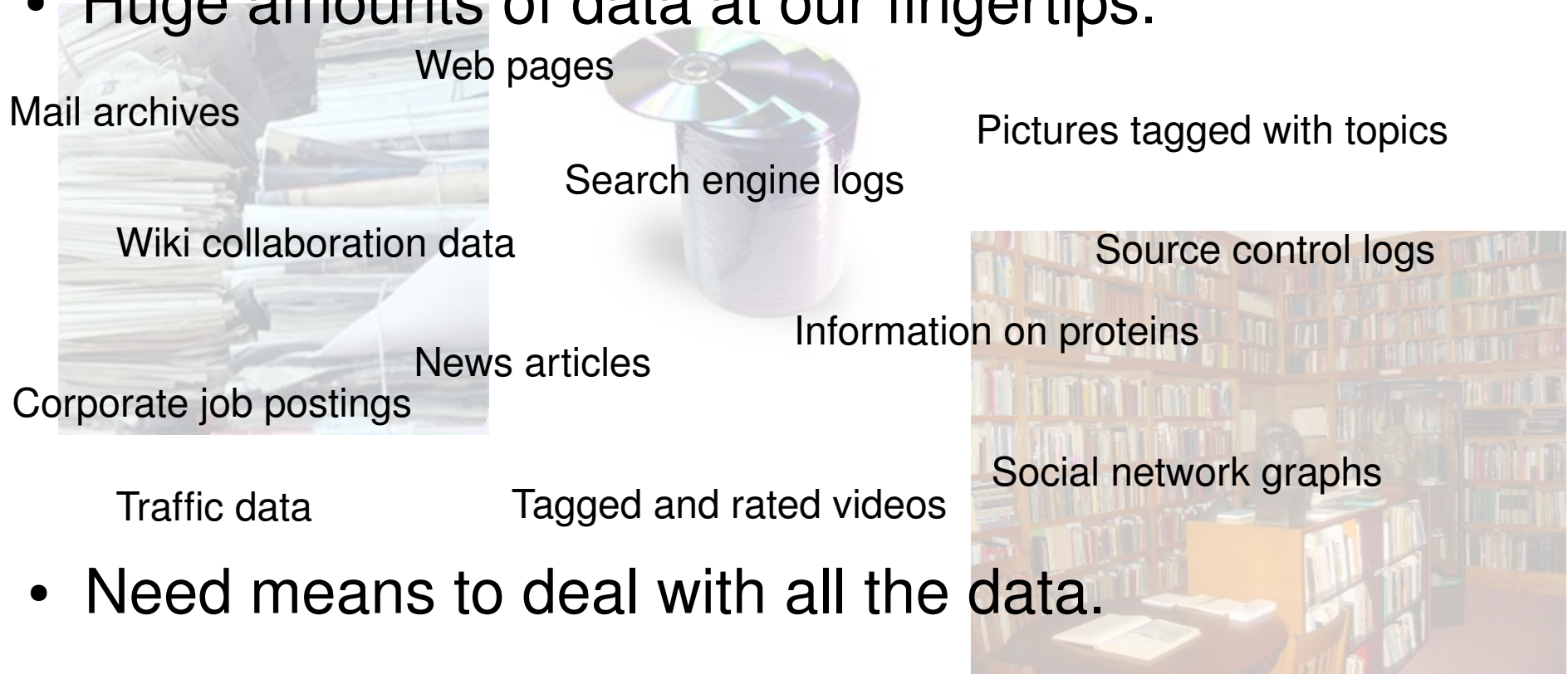
# Apache Mahout

Bringing Machine Learning to Industrial Strength

# Problem Setting

---

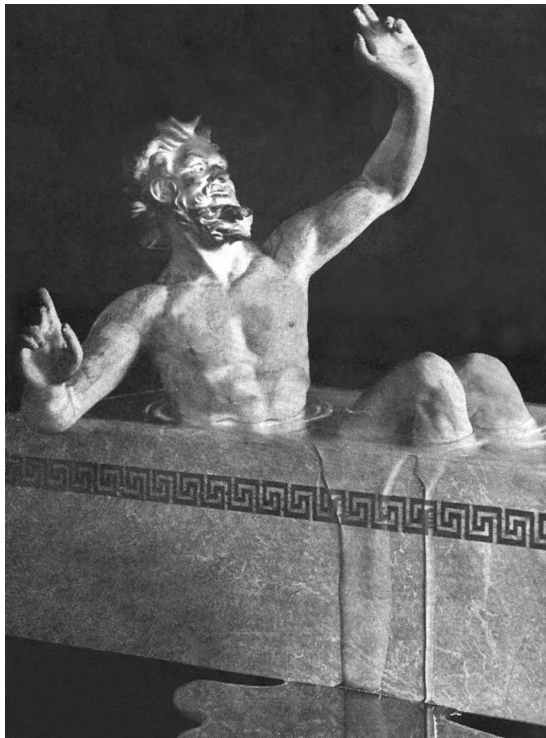
- Huge amounts of data at our fingertips.



- Need means to deal with all the data.
-

# Problem Setting

- Nature generates data.



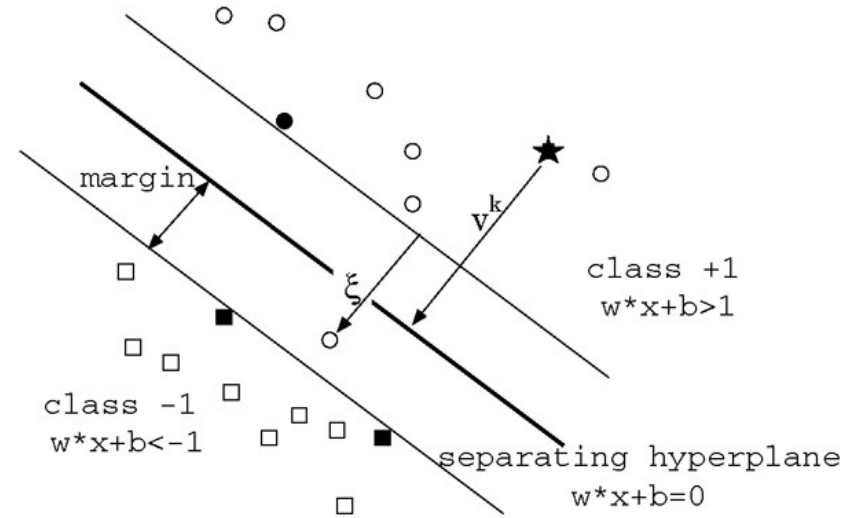
- Archimedes generates model.

$$\frac{\text{Density of Object}}{\text{Density of Fluid}} = .$$

$$\frac{\text{Weight}}{\text{Weight} - \text{Apparent immersed weight}}$$

# Problem Setting

- Nature generates data.
- ML generates models.

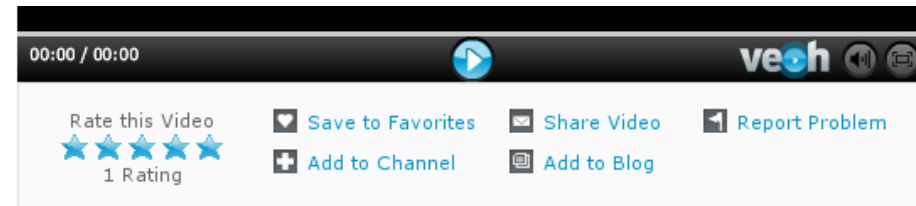


# Where is ML used already?

- Search result clustering.
- “Did you mean” feature.
- Auto completion.
- Language detection.
- Analysis of tags.



Did you mean: [university](#)



# Once upon a time

---



- How it all began:
    - Summer 2007: Crazy developers needed scalable ML.
    - Mailing list and wiki followed quickly.
  - Contacted people
    - from research.
    - from related Apache projects.
  - Rather large community even before project start.
  - 25.01.2008: Project Mahout launched.
-

# Who we are



Dawid Weiss  
Carrot2



Karl Wettin  
Lucene



Grant Ingersoll  
Lucene PMC



Ted Dunning  
The Veoh guy



Jeff Eastman  
Welcome!

Otis Gospodetnic – Lucene  
Erik Hatcher – Lucene (among others)



Isabel Drost  
(that would be myself)

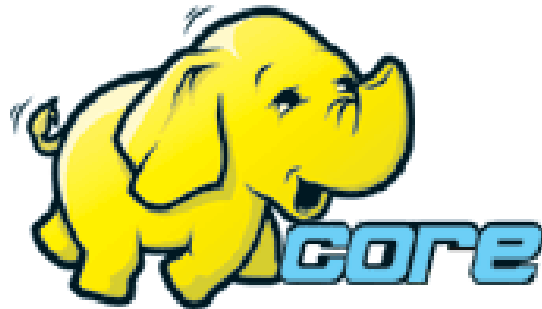
# Our Mission

---

- Build learning algorithms that are scalable.
- Context:



Hama – matrix support



Hadoop – parallelization



Lucene – provides the use cases

---



# Initial contributions

---

- k-Means implementation.
    - Started with non-parallel version.
    - Ported to Hadoop already.
  - Matrix computation package.
    - Building block of many machine learning algorithms.
    - Together with Hama towards parallel matrices.
- 
-

# Initial Contributions

---



- Work on Naïve Bayes, Perceptron, PLSI/EM.
- Integrate Taste
  - Collaborative filtering project at sourceforge.
  - Item based recommendation.
  - User based recommendation.



# GSoC @ Mahout

---



- *“Implementing Logistic Regression in Mahout”*
- *“Codename Mahout.GA for mahout-machine-learning”*
- *“The Implementation of Support Vector Machine Algorithm at Hadoop Platform”*
- *“Application to participate in Mahout”*
- *“Mahout application (Neural Networks)”*
- *“DeCoDe - A smart code search engine based on lucene to show how Mahout work.”*
- *“Applying for mahout machine learning (Neural Networks)”*
- *“Implementation of the Principal Components Analysis algorithm for Mahout”*
- *“MAHOUT Naive Bayes implementation”*

# Conclusions

---



- This is just the beginning: Even the logo is a draft :)
  - High demand for scalable machine learning.
  - We need you – in case you have:
    - A good deal of enthusiasm.
    - Solid mathematical knowledge to understand ML papers.
    - Either proficient in or willing to learn about Hadoop.
    - Or: A lot of data and want to know what to learn from it.
  - [mahout-dev@apache.org](mailto:mahout-dev@apache.org) [mahout-user@apache.org](mailto:mahout-user@apache.org)
-