# Apache Mahout

## Scaling Machine Learning

Presented by:
Isabel Drost

# Agenda

- Motivation.

- Machine learning?

- Introducing Mahout.

- How can you help?

# Some motivation.

# Follow news stories



September 10, 2008 by Alex Barth
http://www.flickr.com/photos/a-barth/2846621384



Search through papers. Automatic topic tracker.

# Movie recommendation



March 22, 2008 by Crystian Cruz
http://www.flickr.com/photos/crystiancruz/2353895708



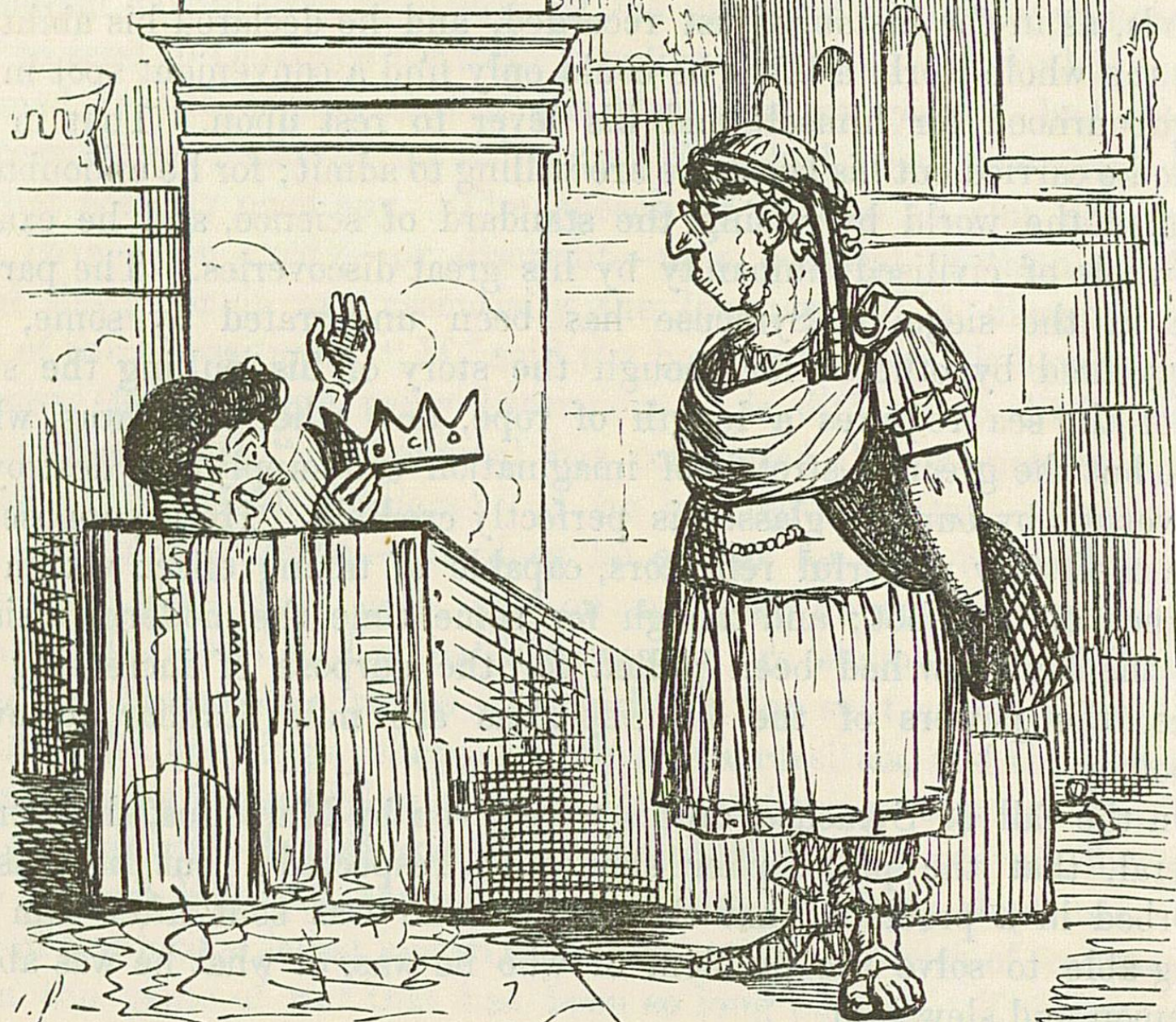IMDB + movie reviews. Aggregate reviews from IMDB, twitter, ...

- Lots and lots of data.

- Structured and unstructured.

# Mission

Provide scalable data mining algorithms.

# Machine Learning?
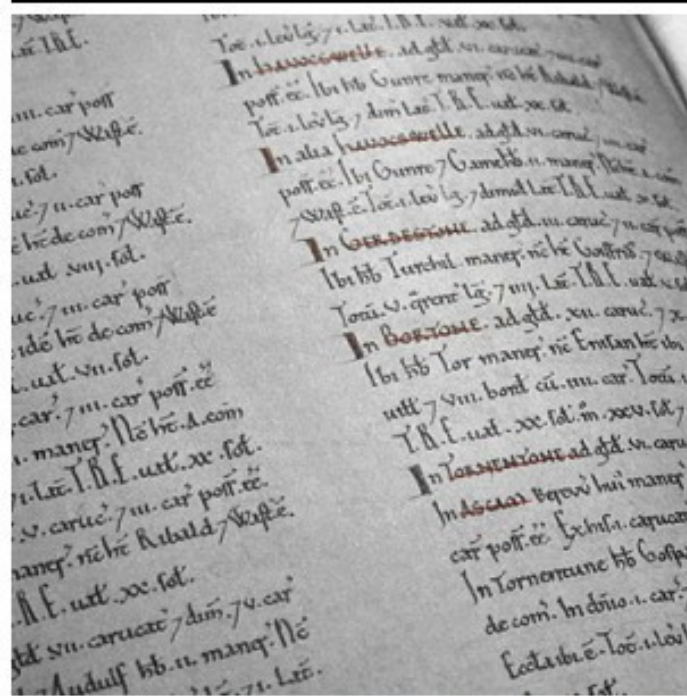
εὕρηκα

# Archimedes generates model:
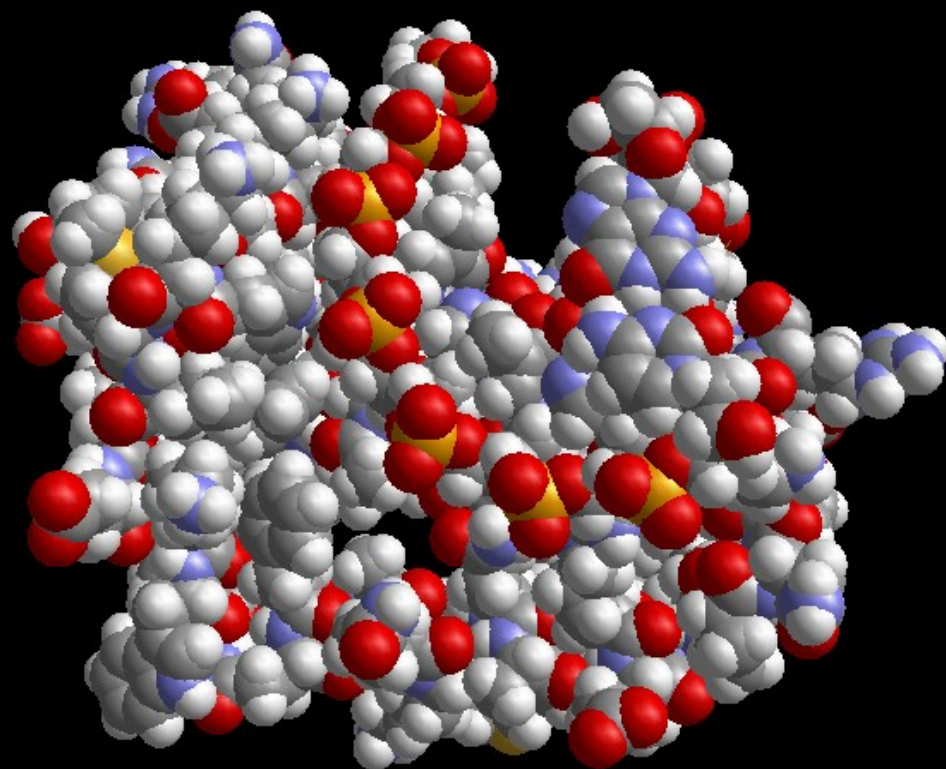
$$\frac{Density\ of\ Object}{Density\ of\ Fluid} = .$$

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

# Machine learning generates model

# Machine learning pipeline

| Gather data. (and meta data). | → | Identify characteristics. | → | Chose right algorithm. |
| --- | --- | --- | --- | --- |

| Keep model in sync when nature changes. | ← | Train on the gathered data. | ← | Tune parameters of your algorithm. |
| --- | --- | --- | --- | --- |

"a million Letters"
D Sharon Pruitt
95 Ent Road
Hanscom AFB MA, 01731

SERVE WITH PRIDE

JURY DUTY
SERVE WITH PRIDE
USA 41

USA FIRST-CLASS FOREVER

Post Card

Hugs
de Loto -

e touch
e spot,

Hello
S. 3.1

Hello
You!

# Machine learning pipeline

| Gather data. (and meta data). | → | Identify characteristics. | → | Chose right algorithm. |
|---|---|---|---|---|

| Keep model in sync when nature changes. | ← | Train on the gathered data. | ← | Tune parameters of your algorithm. |
|---|---|---|---|---|

E-Bay

Auction
status?

Phishing
Spam?

Different
topic

Requested
password?

password

Apache

One of your mails: $\begin{pmatrix} 0 \\ 1 \\ 1 \\ 0 \\ 0 \\ 0 \\ 0 \\ 1 \\ 1 \end{pmatrix}$

. . .

Hadoop

London

Lucene

London

# Machine learning pipeline

| | | |
|---|---|---|
| Gather data. (and meta data). | Identify characteristics. | Chose right algorithm. |
| Keep model in sync when nature changes. | Train on the gathered data. | Tune parameters of your algorithm. |

margin

class +1
w*x+b>1

class -1
w*x+b<-1

$v^k$

$\xi$

separating hyperplane
w*x+b=0

# Machine learning pipeline

| Gather data. (and meta data). | → | Identify characteristics. | → | Chose right algorithm. |

Chose right algorithm. → Tune parameters of your algorithm.

| Keep model in sync when nature changes. | ← | Train on the gathered data. | ← | Tune parameters of your algorithm. |

# Parameter tuning

- Penalty for mistakes.

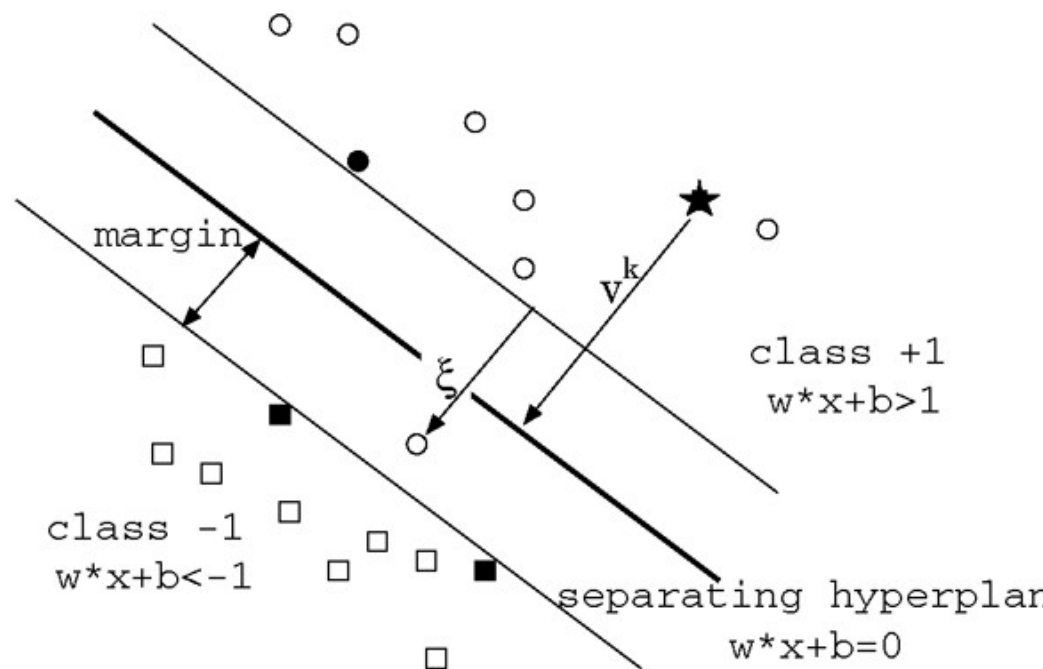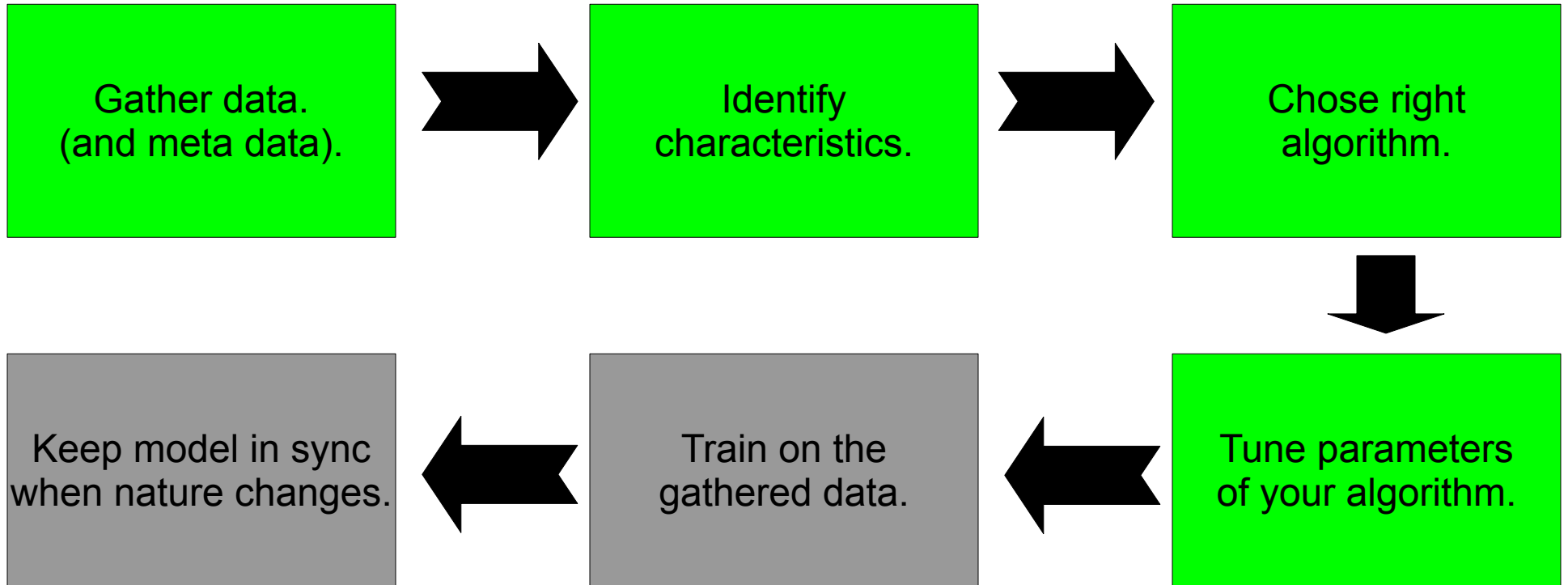- Kernel type for data transformation.

- Tune kernel parameters.

# Machine learning pipeline

| | | |
|---|---|---|
| Gather data. (and meta data). | Identify characteristics. | Chose right algorithm. |
| Keep model in sync when nature changes. | Train on the gathered data. | Tune parameters of your algorithm. |

# Training

- Build model from data.

# Machine learning pipeline

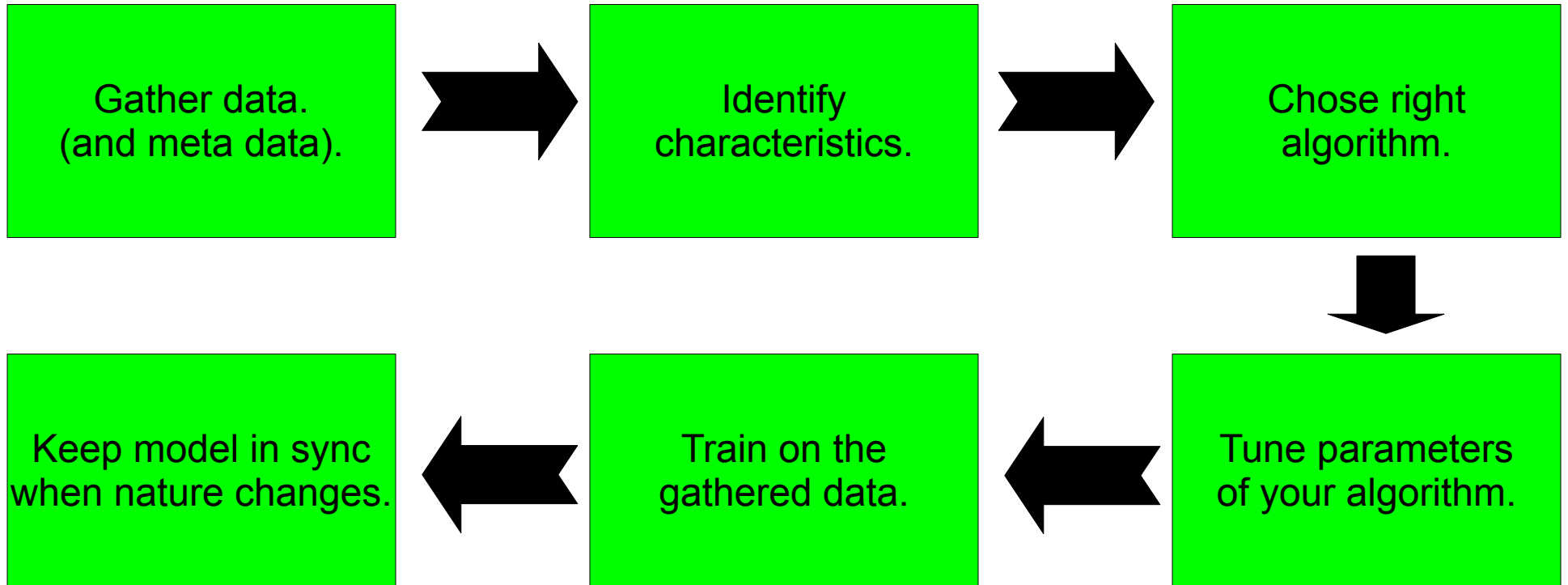| | | |
|---|---|---|
| Gather data. (and meta data). | → Identify characteristics. | → Chose right algorithm. |
| Keep model in sync when nature changes. | ← Train on the gathered data. | ← Tune parameters of your algorithm. |

# Nature changes?

- Spammers adapt to spam filters.
- Users write mails in different styles.
- Expand to new languages.
- ...

# Machine learning pipeline

| Gather data. (and meta data). | → | Identify characteristics. | → | Chose right algorithm. |
|---|---|---|---|---|

↓

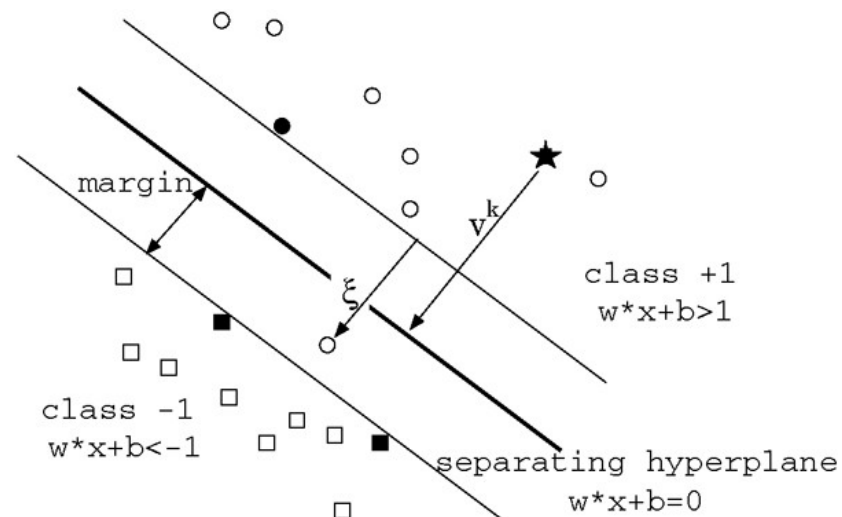| Keep model in sync when nature changes. | ← | Train on the gathered data. | ← | Tune parameters of your algorithm. |
|---|---|---|---|---|

# Introducing Mahout

# Classification

- Categorize data.

- Examples:

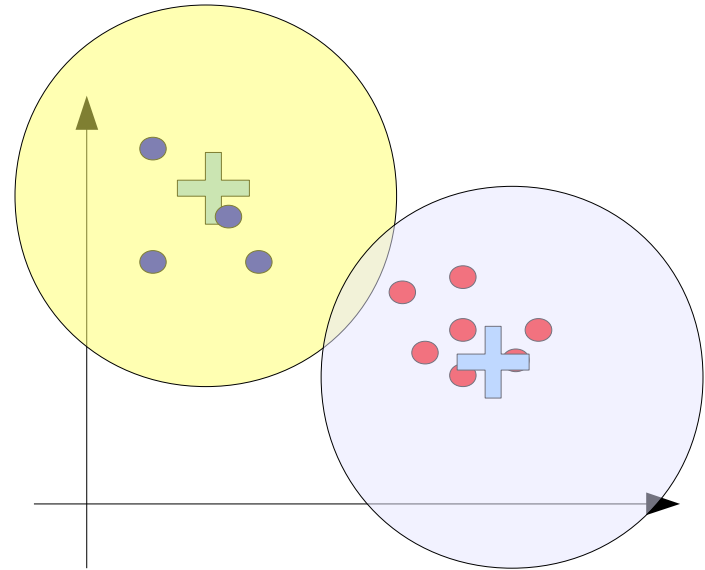  - Identify spam mails.

  - Classify movies as "Action", "Comedy" ...

# Classification

- Naive bayes.

- Complementary naive bayes.

- Winnow/Perceptron

- Others upcoming.

# Discovering groups of data



- Group data by similarity.

- Examples:
  - News articles by topic.
  - Developers by favorite modules.

# Discovering groups of data

- Canopy.

- K-Means.

- Dirichlet based.

- PLSI.

- Others upcoming.

# Recommendation mining

- Recommend items.

- Examples:

    - Find books a user my like.

    - Identify movies a user likes.

# Upcoming

- More algorithms.

- More examples.

# What Mahout can do for you

"Why should I participate?"

Jumpstart your project with proven code.

Discuss with researchers and engineers.

Become a community member.

**http://.../pub/mirrors/apache/lucene/mahout/0.1/**

Thank you to all those
making this possible.

mahout-dev@apache.org

mahout-user@apache.org

- We need You:
  - Enthusiasm.
  - Mathematical knowledge.
  - Proficiency in Hadoop.
  - Interest in understanding data.

July 9, 2006 by trackrecord
http://www.flickr.com/photos/trackrecord/185514449

I WANT YOU

# Some advertising

Berlin - June* at 5p.m.

newthinking store Berlin

Tucholskystr. 48

Hadoop** User/Developer Meeting Germany

\* Exact date is set by speaker – that is you!

\*\* Lucene, Tika, Solr, UIMA, Mahout, katta, ... people welcome.