# ONLINE MACHINE LEARNING AND DATA MINING
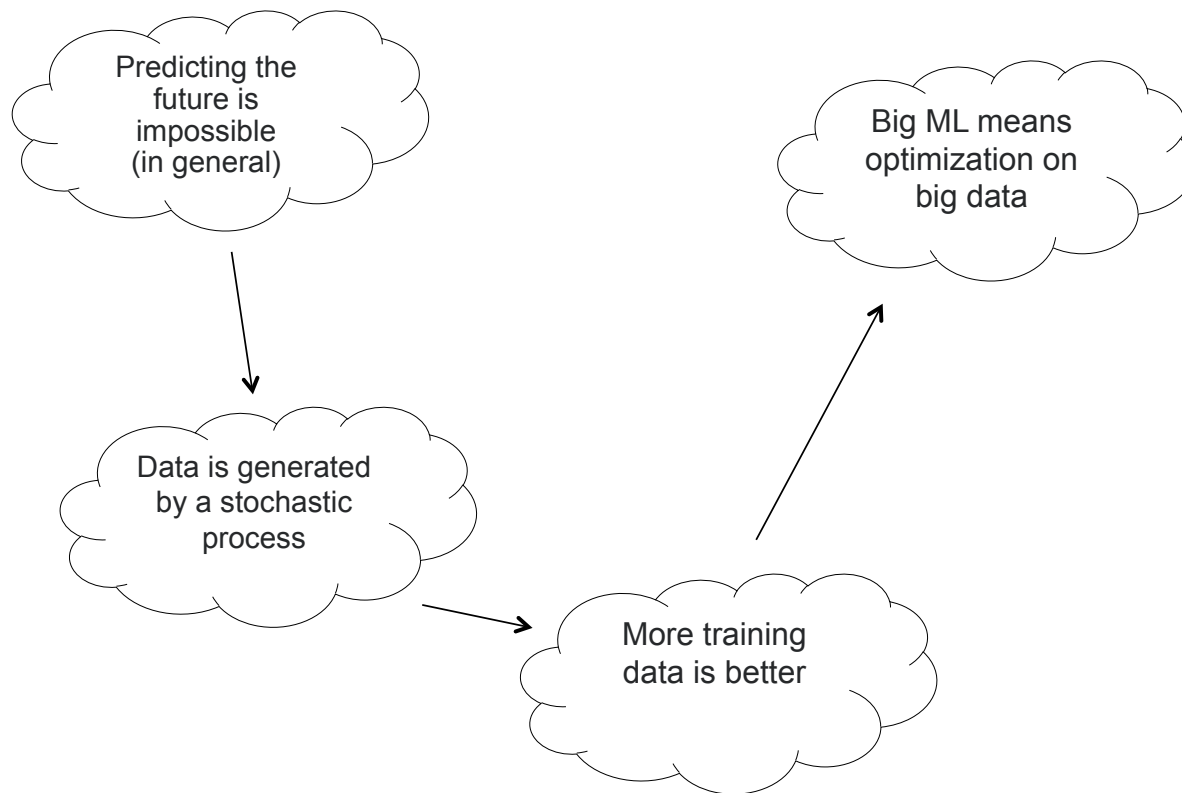
EDO LIBERTY

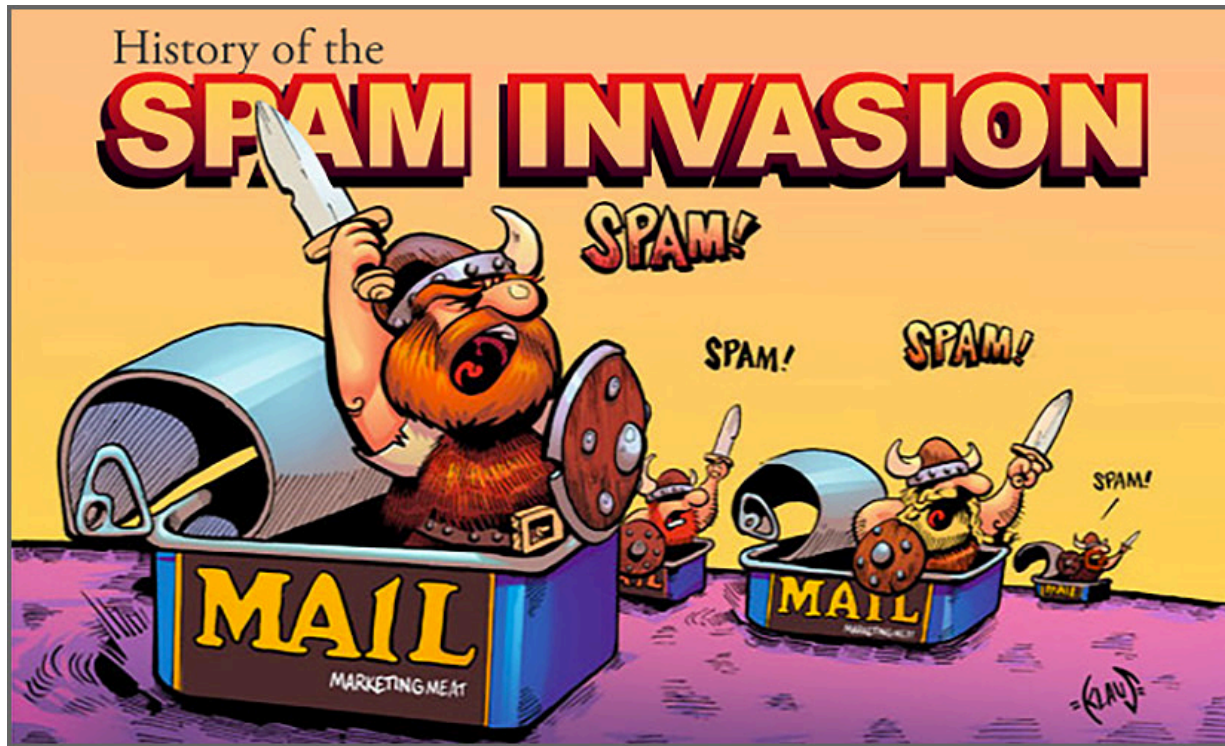# STANDARD MACHINE LEARNING SETTING

= "call"

= "crawl"

= ?

Training Data

Test Data

Train

Model

Apply

"call"

Label

Frequency (40 bands)

time

Convolutional layer

Pooling layer

Fully-connected layers

Softmax layer

YAHOO!

# STANDARD MACHINE LEARNING SETTING

Predicting the future is impossible (in general)

Data is generated by a stochastic process

More training data is better

Big ML means optimization on big data

YAHOO!

## MORE DATA IS OFTEN WORSE (MORE DATA = OLDER DATA)

# OUR ACTIONS HEAVILY INFLUENCE THE DATA

# THE FUTURE IS OFTEN NOT LIKE THE PAST!



Same story line or not?

1) The answer depends on the future
2) We have to decide now...

**YAHOO!**

# HAVING "A MODEL" IS COMPLETELY UNIMPORTANT



Elements of information theory, Cover, 1991
Efficient algorithms for universal portfolios, Kalai, Vempala, 2003
Efficient Algorithms for Online Game Playing and Universal Portfolio Management, Agarwal, Hazan, 2006

YAHOO!

# ONLINE ALGORITHMS
# (DECISION MAKING WITHOUT PREDICTING)

YAHOO!

# THE SKI RENTAL PROBLEM



Rent: x$ /day

Buy: 1000$

YAHOO!

# THE SKI RENTAL PROBLEM

```
      ┌────┐
      │ 70 │
      └────┘
         │
      ⎰  ⎱
       ↓
 ┌──────────────┐
 │  Computation │
 └──────────────┘
         │
         ↓
      ┌────┐
      │ R  │
      └────┘
```

70
+



_____

70

YAHOO!

# THE SKI RENTAL PROBLEM

| 70 | 90 |

Computation

| R | R |

70
+
90
_____
160

YAHOO!

# THE SKI RENTAL PROBLEM

| 70 | 90 | 80 |

Computation

| R | R | B |

70
+
90
+
1000

_____

1160

YAHOO!

# THE SKI RENTAL PROBLEM

| 70 | 90 | 80 | 70 |
|----|----|----|----|

Computation

| R | R | B | |
|---|---|---|---|



70
+
90
+
1000
+
0
_____

1160

# THE SKI RENTAL PROBLEM

| 70 | 90 | 80 | 70 |
|---|---|---|---|

Computation

| R | R | R | R |
|---|---|---|---|



```
  70
+
  90
+
  80
+
  70
_____
 310
```

You should have rented all along...

# THE SKI RENTAL PROBLEM

| 70 | 90 | 80 | 70 | ·· | | ·· | 90 | 88 | 72 | 79 | Input

Computation

$1000 (green)    $1000 (blue)

| R | R | R | R | ·· | | ·· | B | | | | Output

YAHOO!

# THE SKI RENTAL PROBLEM



ALG <= 2 OPT

Algorithm

Buy

Optimal in hindsight

YAHOO!

# ONLINE LINEAR CLASSIFICATION

YAHOO!

# ONLINE MACHINE LEARNING

Emails

Computation

Spam?   N

**YAHOO!**

## ONLINE MACHINE LEARNING

Emails

Computation

Spam?  N  N

YAHOO!

# ONLINE MACHINE LEARNING

Emails

Computation

Spam?  N  N  Y

YAHOO!

# ONLINE MACHINE LEARNING



Number of mistakes is compared to the best classifier in hindsight!

Variants of SGD have this property

Computation

N N Y N · · · · Y N N Y

Prediction, Learning, and Games, Cesa-Bianchi, Lugosi, 2006

YAHOO!

# ONLINE PRINCIPAL COMPONENT ANALYSIS

Online Principal Components Analysis, Boutsidis, Garber, Karnin, Liberty 2014
Online PCA with Spectral Bounds, Karnin, Liberty, 2015

YAHOO!

$x_i$

YAHOO!

$$\Phi\Phi^T x$$

$$x_i$$

$\Phi\Phi^T x$

$x_i$

$\|x_i - \Phi^T\Phi x_i\|$

Eigenpets: https://bioramble.wordpress.com/2015/09/01/

YAHOO!

# ONLINE PRINCIPAL COMPONENT ANALYSIS

---
**Algorithm 1** Fixed Error: Conceptual Algorithm

---
**input:** $X, \Delta$
$U \leftarrow$ all zeros matrix
**for** $x_t \in X$ **do**
    **if** $\|(I - UU^T)X_{1:t}\|^2 \geq \Delta$
        Add the top left singular vector of $(I - UU^T)X_{1:t}$ to $U$
    **yield** $y_t = U^T x_t$
**end for**

---

Online PCA with Spectral Bounds, Karnin, Liberty, 2015

YAHOO!

# Online PCA with Spectral Bounds

# ONLINE K-MEANS CLUSTERING

An Algorithm for Online K-Means Clustering, Liberty, Sriharsha, Sviridenko, 2014

YAHOO!

# K-MEANS CLUSTERING

http://research.ics.aalto.fi/mi/software/ne/

YAHOO!

# K-MEANS CLUSTERING

- Roughly 20,000 documents
- 20 topics:
    - Graphics
    - PC hardware
    - Baseball
    - For-sale
    - Politics
    - …

http://qwone.com/~jason/20Newsgroups/

http://research.ics.aalto.fi/mi/software/ne/

YAHOO!

# K-MEANS CLUSTERING

1) One can cluster
   points fully online

2) Create only slightly
   more than k centers

3) Be competitive with the best
   offline clustering to k clusters

---

**Algorithm 2** Online $k$-means algorithm

---

**input:** $V, k$

$C \leftarrow$ first $k + 1$ distinct vectors in $V$; and $n = k + 1$

(For each of these **yield** itself as its center)

$w^* \leftarrow \min_{v,v' \in C} \|v - v'\|^2 / 2$

$r \leftarrow 1; q_1 \leftarrow 0; f_1 = w^*/k$

**for** $v \in$ the remainder of $V$ **do**

    $n \leftarrow n + 1$

    **with probability** $p = \min(D^2(v, C)/f_r, 1)$

        $C \leftarrow C \cup \{v\}; q_r \leftarrow q_r + 1$

    **if** $q_r \geq 3k(1 + \log(n))$ **then**

        $r \leftarrow r + 1; q_r \leftarrow 0; f_r \leftarrow 2 \cdot f_{r-1}$

    **end if**

    **yield:** $c = \arg\min_{c \in C} \|v - c\|^2$

**end for**

---

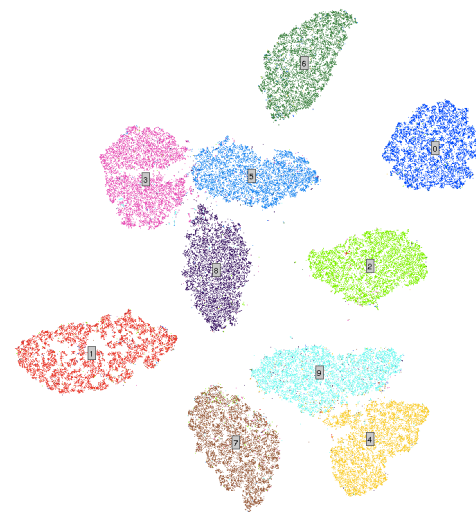An Algorithm for Online K-Means Clustering, Liberty, Sriharsha, Sviridenko 2015

YAHOO!

# ONLINE K-MEANS CLUSTERING



An Algorithm for Online K-Means Clustering, Liberty, Sriharsha, Sviridenko 2015
k-means++: the advantages of careful seeding, Arthur, Vassilvitskii, 2006

# STREAMING ALGORITHMS

# OPEN SOURCE FROM YAHOO
EDO LIBERTY

YAHOO!

# DATASKETCHES.GITHUB.IO



Sketches Library from YAHOO!

A Java software library of *stochastic* *streaming algorithms*

Overview Download GitHub Comments

YAHOO!

YAHOO!

# DISTRIBUTED STORAGE

YAHOO!

# DISTRIBUTED MODEL (MAP/REDUCE, MESSAGE PASSING, ...)

# DISTRIBUTED MODEL (INDEXES, TABLES, DATABASES, ...)

YAHOO!

# BIG-DATA *META* INFOGRAPHIC

# THE STREAMING COMPUTATIONAL MODEL

YAHOO!

# THE STREAMING COMPUTATIONAL MODEL

| 1 | 7 | 8 | 1 | $\cdot\,\cdot$ | | $\cdot\,\cdot$ | 0 | 1 | 7 | 7 |

$O(n)$ Items

Iterator

Computation

$O(\text{polylog}(n))$ Space

Sketch

Query

YAHOO!

# THE DISTRIBUTED STREAMING COMPUTATIONAL MODEL



The World

Sketch   Sketch   Sketch   Sketch

Merge → Sketch

YAHOO!

# Number of users (easy)



data          Map          Reduce
(count)       (sum)

YAHOO!

# Web Site Logs

| Time | User ID | Site | Time Spent Sec | Items Viewed |
|------|---------|------|----------------|--------------|
| 9:00 | U1 | Apps | 59 | 5 |
| 9:30 | U2 | Apps | 179 | 15 |
| 10:00 | U3 | Music | 29 | 3 |
| 1:00 | U1 | Music | 89 | 10 |
| ... | ... | ... | ... | ... |

# Financial Transactions System Log

| Time | User ID | Site | Purchased | Revenue |
|------|---------|------|-----------|---------|
| 9:00 | U1 | Apps | FaceTune | $3.99 |
| 9:30 | U2 | Apps | Minecraft | $6.99 |
| 10:00 | U3 | Music | Purple Rain | $1.29 |
| 10:05 | U3 | Apps | Minecraft | $6.99 |
| ... | ... | ... | ... | ... |

## Unique User Queries
- Unique users viewing Apps since 9:45...?
- Unique users visiting Apps site AND Music site?
- Unique users visiting Apps site AND NOT Music site?

## Quantile Queries
- The median and 95%ile Time Spent seconds by ...?
- A Frequency Histogram of Time Spent by Split-Points specified at query time?

## Frequency Queries
- The numbers of times each app was purchased

## Join Queries
- For all users that purchased Apps, what is the average / median time spent?

YAHOO!

# Number of <u>unique</u> users (hard)



data     Map (key=user)     Reduce (return 1)     Reduce (sum)

YAHOO!

# Number of <u>unique</u> users (made easy)



data          Map
              (sketch)          Reduce
                                (merge)

YAHOO!

# Current Sketch Implementations

## Count Unique Sketches

– Both Theta Sketches* and HLL Sketches

– **Estimating Cardinality** of a stream of identifiers with duplicates

– **Set Operations** (e.g., Union, Intersection, and Difference)

– Can be extended to produce approximate Joins

## Quantiles Sketches

– Normal or Inverse PMF's, CDF's of streams of numeric values,

using after-the-fact queries.

## Frequent Item Sketches

– Identify the Heavy Hitters of arbitrary objects from a stream of objects

– Estimate the frequency of any item from the stream

YAHOO!

# DataSketches.GitHub.io    Open Source Library

- Dedicated to production quality Sketch implementations.
  - These are not toy algorithms!
  - Heavily used within Yahoo

- Common Attributes
  - True streaming. Single pass, "one-touch" algorithms for either *real-time* or *batch*
  - All Sketches are Mergeable, which makes them highly parallelizable.
  - Designed for multiple large-scale computing environments:
    - Core of library is coded in Java with no external dependencies
    - Easy integration into virtually any system environment
    - Adaptors for Hadoop/Pig and Hadoop/Hive environments
    - Standard library promotes sharing across platforms and organizations
  - Maven deployable and registered with *Maven Central Repository*
    - *http://search.maven.org/#search|ga|1|datasketches*
  - Comprehensive unit tests and testing tools are provided
  - Extensive documentation with Systems Developers in mind
  - All algorithms are backed by published mathematical theory

**YAHOO!**

# Counting distinct elements example

```
$ less emails.csv | wc -l
 10000000
```
→ 10M sender domains from inbound emails

```
$ head -n 5 emails.csv
facebookmail.com
jobsdbalert.co.id
facebookmail.com
twitter.com
bonsplansdujour.net
```
→ There are duplicates

```
$ cat emails.csv | sort | uniq | wc -l
^C

$ cat emails.csv | sort -u -S 100% | wc -l
^C
```
→ Roughly 200Mb and several minutes of CPU (~25 seconds for numbers)

```
$ cat emails.csv | sketch uniq
47618   40772   55589

$ cat emails.csv | sketch uniq 0.01
53782   53351   54216
```
→ < 10Kb of memory and 1.5 Seconds!

YAHOO!