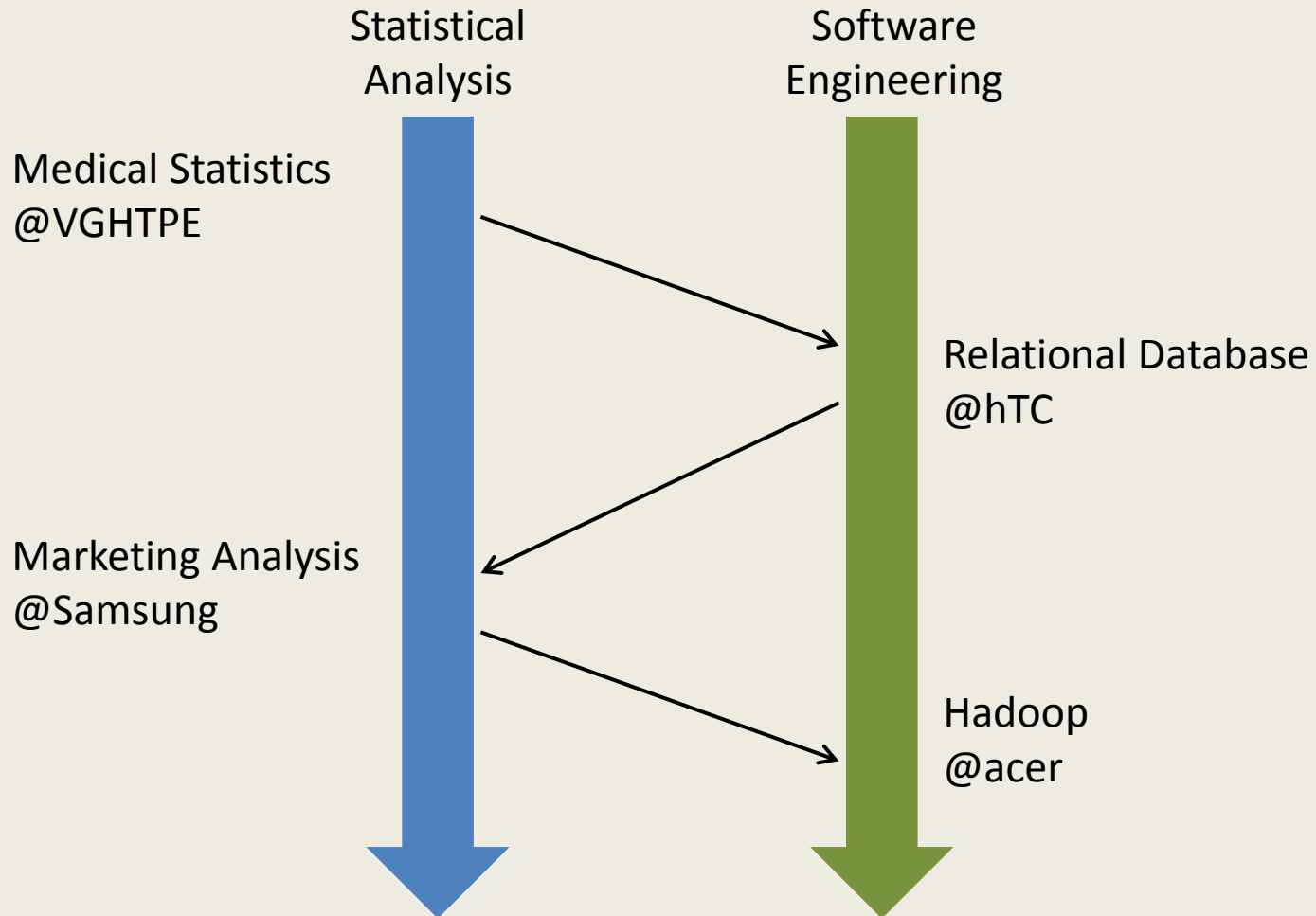


Between Statistics and IT



Definition

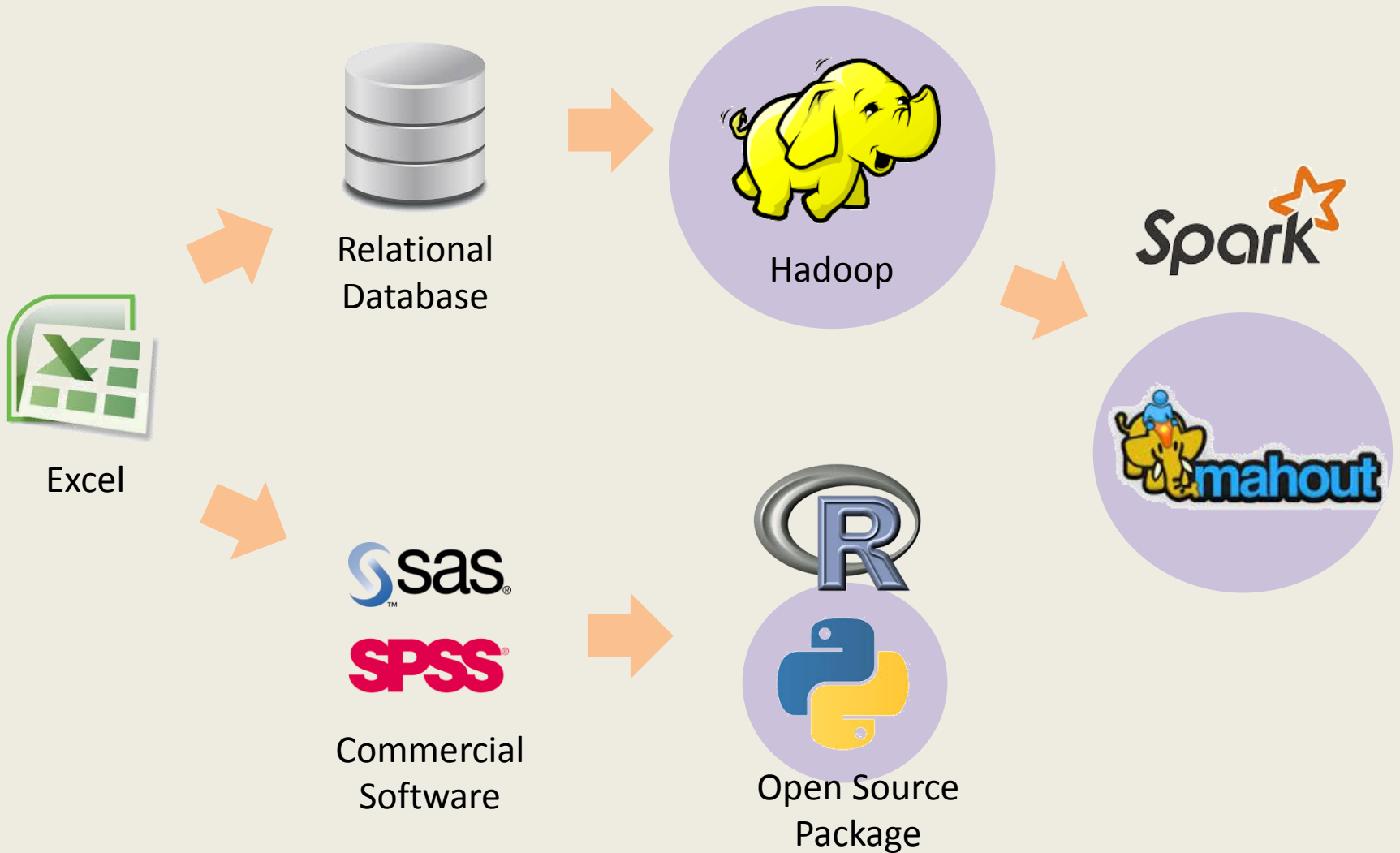
- Data Scientist (n.):

Person who is better at statistics than any software engineer and better at software engineering than any statistician.

-- Josh Wills @Cloudera



Useful Tools



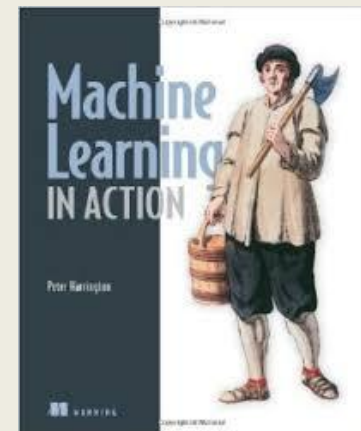
Python Packages

- NumPy – 向量運算
- SciPy – 科學計算
- Matplotlib – 繪圖
- Pandas – 資料分析
- StatsModels – 統計模型
- scikit.learn – 機器學習

- Integrated Package: Anaconda of Python(X,Y)

Machine Learning in Action

- Machine Learning Algorithms with Python
 - Classification: KNN, decision tree, naïve Bayes, logistic regression, SVM, AdaBoost
 - Predictive: regression, tree-based regression
 - Unsupervised: K-mean, Apriori, FP-growth
 - Others: PCA, SVD



Famous Cases

- Walmart - 啤酒與尿布
 - Walmart運用購物籃分析, 發現每週四啤酒與尿布一起銷售的機率非常高, 調查後發現男性常在週四採購周末Party需要的啤酒, 妻子則會順便提醒先生順便買小孩的尿布. 得知此事實後 Walmart調整這兩項商品貨架位置, 使得獲利上升.
-- 1998年哈佛商業評論
- Target - 預測少女懷孕
 - 一位父親怒氣沖沖的到Target百貨向經理控訴: 你們怎麼可以寄送嬰兒用品折價券給未成年高中生女兒? 一個月後, 父親又回去跟經理道歉, 原來女兒真的懷孕了. Target百貨依據消費者個人資料以及購物行爲改變(無香味乳液與葉酸營養品), 成功預測消費者已經懷孕
-- 2012年紐約時報

Hadoop Streaming with Python

- `hadoop jar $HADOOP_HOME/hadoop-streaming.jar \`
`-D mapred.job.name='job001' \`
`-input /data \`
`-output /output \`
`-file knn_map.py \`
`-mapper knn_map.py \`
`-file knn_reduce.py \`
`-reducer knn_reduce.py`

Hive with Python

- ```
from (
 select transform(a)
 using '${env:HOME}/josh/script/s1.py'
 as a1
 from table1
) t1
select transform(a1)
using '${env:HOME}/josh/script/s2.py'
as b1
;
```



# Mahout

- 機器學習函式庫
- MapReduce為主, 目前轉向Spark
- 演算法三大類別:
  - 推薦(recommendation)
  - 群集(clustering)
  - 分類(classification)



# Recommendation (User-Based)

假設有顧客1~5,  
對商品編號101~107  
的評價如下表  
現在要向顧客1號推薦商品

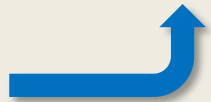
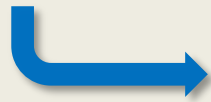
經計算後, 將4, 5號顧客評價高  
分的商品104, 106推薦給顧客1

| item \ user | 1   | 2   | 3   | 4   | 5   |
|-------------|-----|-----|-----|-----|-----|
| 101         | 5   | 2   | 2   | 5   | 4   |
| 102         | 3   | 2.5 |     |     | 3   |
| 103         | 2.5 | 5   |     | 3   | 2   |
| 104         |     | 2   | 4   | 4.5 | 4   |
| 105         |     |     | 4.5 |     | 3.5 |
| 106         |     |     |     | 4   | 4   |
| 107         |     |     | 5   |     |     |

依據共同評價過的商品  
計算顧客之間的相似度  
發現顧客1號與4, 5號正相關  
與2號負相關  
3號則無法計算(樣本不夠)

| user \ user | 1     | 2     | 3     | 4     | 5     |
|-------------|-------|-------|-------|-------|-------|
| 1           | 1.00  | -0.76 | -     | 1.00  | 0.94  |
| 2           | -0.76 | 1.00  | -     | -0.97 | -0.94 |
| 3           | -     | -     | 1.00  | -1.00 | -0.65 |
| 4           | 1.00  | -0.97 | -1.00 | 1.00  | 0.88  |
| 5           | 0.94  | -0.94 | -0.65 | 0.88  | 1.00  |

| item \ user | 1   | 2   | 3   | 4   | 5   |
|-------------|-----|-----|-----|-----|-----|
| 101         | 5   | 2   | 2   | 5   | 4   |
| 102         | 3   | 2.5 |     |     | 3   |
| 103         | 2.5 | 5   |     | 3   | 2   |
| 104         | 5   | 2   | 4   | 4.5 | 4   |
| 105         |     |     | 4.5 |     | 3.5 |
| 106         | 4   |     |     | 4   | 4   |
| 107         |     |     | 5   |     |     |



# Recommendation Options

- 推薦方法
  - User-based, Item-based, Slope-One, SVD, Knn Item-Based, Tree Cluster
- 相似度估計
  - Pearson Correlation, Euclidean Distance, Consine, Spearman Correlation, Tanimoto Coefficient, LogLokelihood
- 鄰居選擇方法
  - Nearest-N, Threshold
- 其他因素
  - Boolean? Weight?

# Run MapReduce

```
hadoop jar $MAHOUT_HOME/mahout-core-0.9-job.jar \
org.apache.mahout.cf.taste.hadoop.item.RecommenderJob \
-Dmapred.input.dir=/data \
-Dmapred.output.dir=/output \
-s SIMILARITY_PEARSON_CORRELATION
```

# Clustering

- Google 新聞群組

## 焦點新聞

### 立院裡裡外外全是反服貿抗議學生

中時電子報 - 15分鐘之前    

... tw\_chinatimes. 「反黑箱服貿」的抗議學生19日持續占領立法院青島東路側門停車場及議場，黑壓壓都是反黑箱服貿的抗議學生19日持續占領立法院青島東路側門停車場及議場，黑壓壓都是反...




佔領國會》服貿延燒佔據立院學生表達三訴求 自由時報

群眾佔領議場國會史上首次 中央廣播電台

長篇：反服貿夜襲學生攻立院霸主席台200學生奇襲打破玻璃門衝進議場 台視新聞

[閱讀即時報導](#) »

### 馬航副機師道晚安前早設定調頭

NOWnews - 5分鐘之前    

馬航MH370航班家屬們在周一投票要求直接與馬來西亞政府官員對話。國際中心／綜合報導. 美國媒體報導，馬來西亞失聯航國家廣播公司新聞網 (NBC News) 報導，消息來源透露，...

NBC：馬航轉向12分鐘後才最後通話 中時電子報

馬航道晚安前早設定掉頭 中央廣播電台

長篇：馬航班機失聯前早設定掉頭 中央通訊社

[閱讀即時報導](#) »

### 超強傳美可秘錄某國所有通話

中央廣播電台 - 7分鐘之前    

位於馬里蘭州的美國國安局(NSA)，因為監視全球網路祕密計畫曝光，引起軒然大波，圖為國安局總部的空照圖。(AFP). 「華某個國家境內所有電話，並可回放重聽30天內的通話紀錄。

# Twitter Text Mining

- The Acer engineers are smoking crack again ?



# Twitter Text Mining

## Aspire R7 的四種使用模式



筆電模式

觸控模式

分享模式

平板模式

# Twitter Text Mining

- 文字分析發現Meme已經開始傳播並衍生出三種版本
- 找出病毒傳播的網站並消毒澄清
- 將“acer engineers”設為關鍵字並持續追蹤

|              |                                                   |            |               |
|--------------|---------------------------------------------------|------------|---------------|
| <b>Total</b> |                                                   | <b>495</b> |               |
|              | <b>the acer engineers are smoking crack again</b> | <b>381</b> |               |
|              | i.imgur.com                                       | 240        |               |
|              | 9gag.com                                          | 52         |               |
|              | _invalid                                          | 12         |               |
|              | bestof4chan.org                                   | 9          |               |
|              | www.bestof9gag.com                                | 7          |               |
|              | www.reddit.com                                    | 7          |               |
|              | zumlerr.com                                       | 4          |               |
|              | adf.ly                                            | 2          |               |
|              | imgur.com                                         | 2          |               |
|              | limk.com                                          | 2          |               |
|              | q.gs                                              | 2          |               |
|              | www.facebook.com                                  | 2          |               |
|              | anotherboringblog.com                             | 1          |               |
|              | geekstumbles.com                                  | 1          |               |
|              | hanyokasha.tumblr.com                             | 1          |               |
|              | kartelnometry.tumblr.com                          | 1          |               |
|              | latestfunnystuff.com                              | 1          |               |
|              | latestvcnews.wordpress.com                        | 1          |               |
|              | lazeblog.com                                      | 1          |               |
|              | m.vk.com                                          | 1          |               |
|              | omg-pictures.tumblr.com                           | 1          |               |
|              | tumblokami.tumblr.com                             | 1          |               |
|              | uberhumor.com                                     | 1          |               |
|              | www.awwomg.com                                    | 1          |               |
|              | www.dailymao.net                                  | 1          |               |
|              | <b>the acer engineers are doing drugs again</b>   | <b>103</b> |               |
|              | twitter.com                                       | 103        | @BestProHumor |
|              | <b>Acer engineers want to see the world burn</b>  | <b>11</b>  |               |
|              | www.quiterly.com                                  | 1          |               |



# Classification

- Gmail垃圾郵件分類

The screenshot shows the Gmail interface with the 'Spam (9731)' folder highlighted in the left sidebar. The main content area displays a list of spam messages, including various types of phishing and malware attempts.

philipp.lenssen@

Google Mail BETA

Search Mail Search the W

Compose Mail

Inbox (22)

Starred ☆

Sent Mail

Drafts (18)

All Mail

**Spam (9731)**

Trash

Contacts

Labels

Invite a friend

Delete Forever Not Spam More Actions ... Refresh

Select: All, None, Read, Unread, Starred, Unstarred

(messages that have been in Spam more than 30 days wi

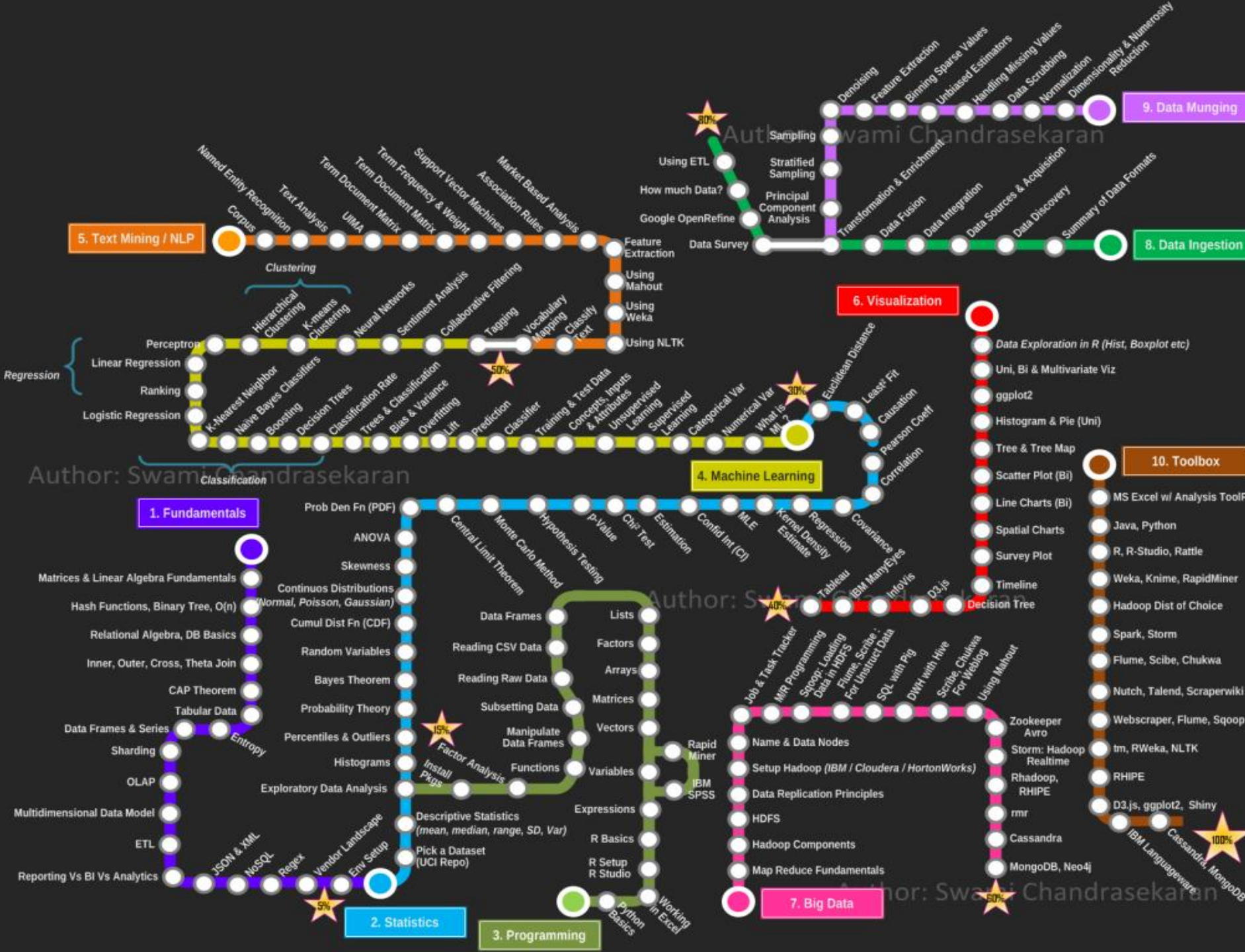
|                          |                          |                           |
|--------------------------|--------------------------|---------------------------|
| <input type="checkbox"/> | ☆ info                   | ←\$BBg?M\$N=P2q\$!y7hl    |
| <input type="checkbox"/> | ☆ info                   | ←\$B:FEY\$N3NG'%Z!<%8←    |
| <input type="checkbox"/> | ☆ 校友朱晓梅                  | philipp.lenssen你好吗,我      |
| <input type="checkbox"/> | ☆ MAILER-DAEMON          | ВНИМАНИЕ: Сообщение       |
| <input type="checkbox"/> | ☆ MAILER-DAEMON          | 郵件傳輸失敗！                   |
| <input type="checkbox"/> | ☆ MAILER-DAEMON          | **Message you sent bloc   |
| <input type="checkbox"/> | ☆ postmaster             | [ERR] "Brianna" expecti   |
| <input type="checkbox"/> | ☆ 瘋狂集殺                   | “父親節”最佳贈禮~~瘋狂             |
| <input type="checkbox"/> | ☆ postmaster             | Undeliverable: Interesti  |
| <input type="checkbox"/> | ☆ Internet Mail Deliverv | Deliverv Notification: De |

# Classification Algorithm

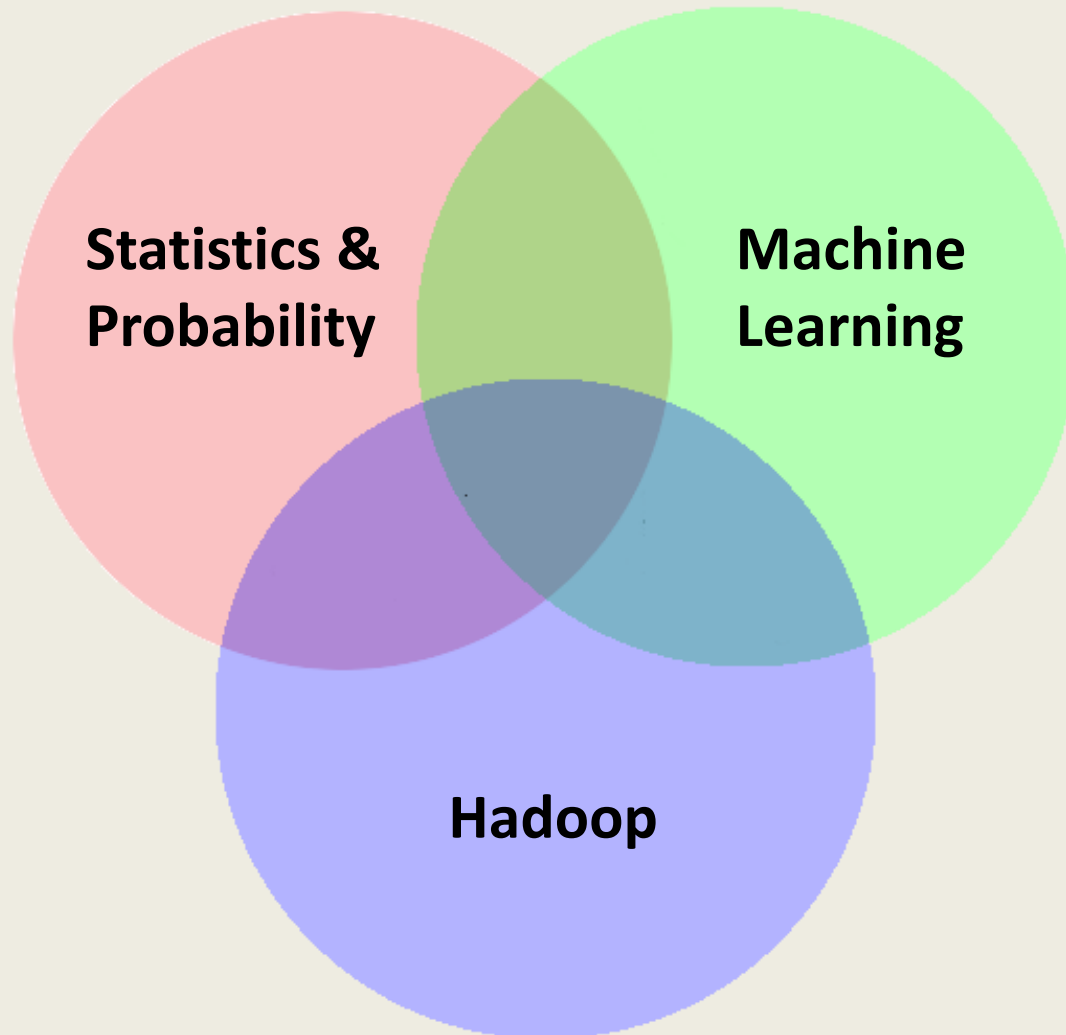
- Logistic Regression
- Naïve Bayesian
- Random Forests
- Hidden Markov Models

# Cloudera Certification

- Hadoop Developer CCDH
- Hadoop Admin CCAH
- HBase Specialist CCSHB
- CCP: Data Scientist
  - Data Science Essentials (DS-200)
  - Data Science Challenge



# Data Science Essentials (DS-200)



# Statistics & Probability

- Statistics Concepts
- Descriptive Statistics
- Probability Distributions
- Bayes Probability
- Likelihood Function
- Bagging and Boosting
- Other Statistics Concepts

# Statistics & Probability

- 參考資料: 各種機率統計書籍
- 注重機率論, 統計模擬及相關線性代數
- 不重一般統計學: 假設檢定, 參數估計, ...

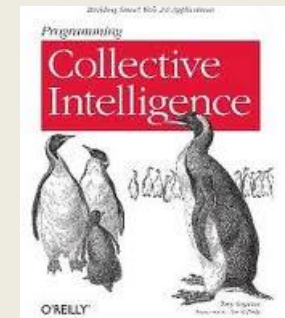
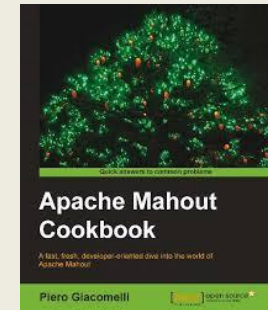
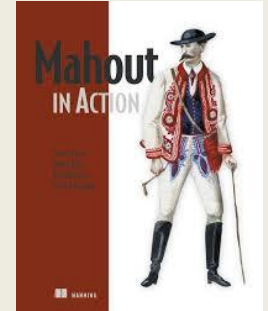
# Machine Learning

- Algorithms Concept
  - Naïve Bayes
  - SVM
  - K-Mean
  - Collaborative Filtering
- Over/Under Fitting
- VC dimension
- Stochastic Gradient Descent



# Machine Learning

- References:
  - Machine Learning in Action
  - Mahout in Action
  - Programming Collective Intelligence
- Coursera:
  - Machine Learning, Andrew Ng
  - 機器學習基石, 林軒田

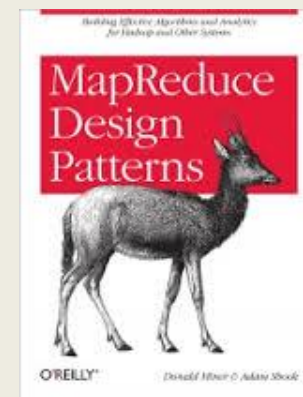
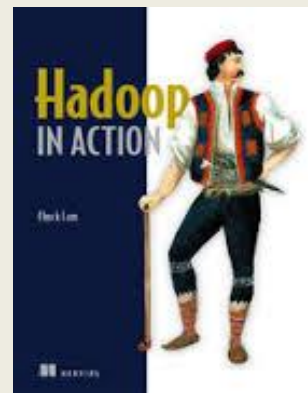
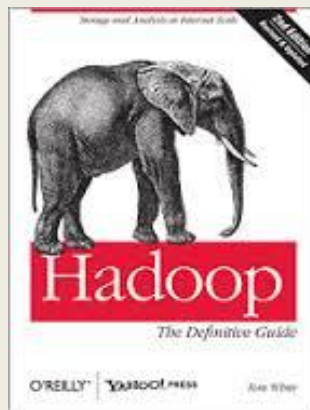


# Hadoop

- MapReduce Concepts
- Combiner and Partitioner
- Distributed Cache
- Pseudo Code Analysis
- Hive, Pig , Sqoop and Flume Use Case

# Hadoop

- Reference:
  - Hadoop Definitive Guide
  - Hadoop In Action
  - MapReduce Design Pattern



# Data Science Challenge

- Web Analytics Challenge: Classification, Clustering, and Collaborative Filtering
  - Explore and Summary Data
  - Clean Data
  - Classify Adult or Kid
  - Cluster Users
  - Recommend File to User

# Solution Kit

## [CCP: Data Scientist Solution Kit](#)

- Hadoop Streaming + Python to Clean and Summary JSON data
- Python + bash shell to classify with SimRank Algorithm
- Cluster users with K-Mean Algorithm using Cloudera KM
- Recommend Film with Mahout

# The End

