



IT Systems Engineering | Universität Potsdam

Natural Language Processing

Machine Learning

Potsdam, 26 April 2012

Saeedeh Momtazi

Information Systems Group

Introduction

2

■ Machine Learning

- Field of study that gives computers the ability to learn without being explicitly programmed.

[Arthur Samuel, 1959]

■ Learning Methods

- Supervised learning
 - Active learning
- Unsupervised learning
- Semi-supervised learning
- Reinforcement learning

Outline

3

- ① Supervised Learning
- ② Semi-Supervised Learning
- ③ Unsupervised Learning

Outline

4

- ① Supervised Learning
- ② Semi-Supervised Learning
- ③ Unsupervised Learning

Supervised Learning

5

Renting budget: 1000 €



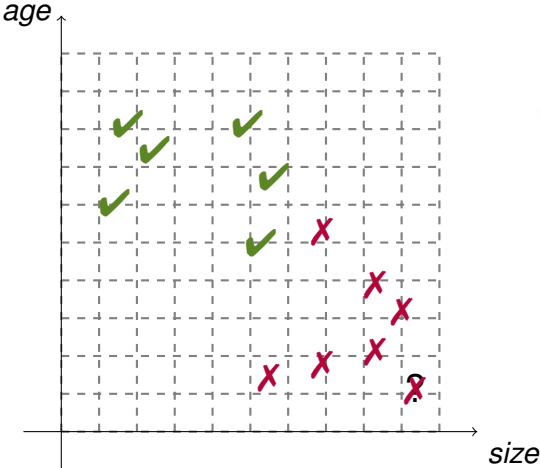
Size: 180 m^2

Age: 2 years



Supervised Learning

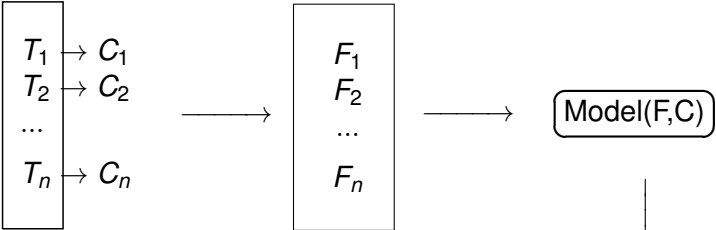
6



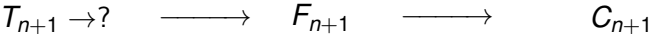
Classification

7

Training



Testing



Applications

8

Problem	Item	Category
POS Tagging	Word	POS
Named Entity Recognition	Word	Named entity
Word Sense Disambiguation	Word	The word's sense
Spam Mail Detection	Document	Spam/Not spam
Language Identification	Document	Language
Text Categorization	Document	Topic
Information Retrieval	Document	Relevant/Not relevant

Part Of Speech Tagging

9

“I saw the man on the roof.”

“I_[PRON] saw_[V] the_[DET] man_[N] on_[PREP] the_[DET] roof_[N]. ”

[PRON]	Pronoun
[PREP]	Preposition
[DET]	Determiner
[V]	Verb
[N]	Noun
...	

Named Entity Recognition

10

“Steven Paul Jobs, co-founder of Apple Inc, was born in California.”

“Steven Paul Jobs, co-founder of Apple Inc, was born in California.”
Person Organization Location

Person
Organization
Location
Date

...

Word Sense Disambiguation

11

“Jim flew his *plane* to Texas.”



“Alice destroys the item with a *plane*.”



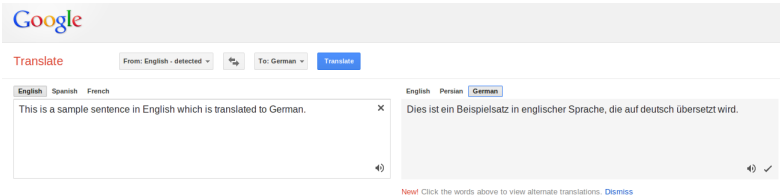
Spam Mail Detection

12



Language Identification

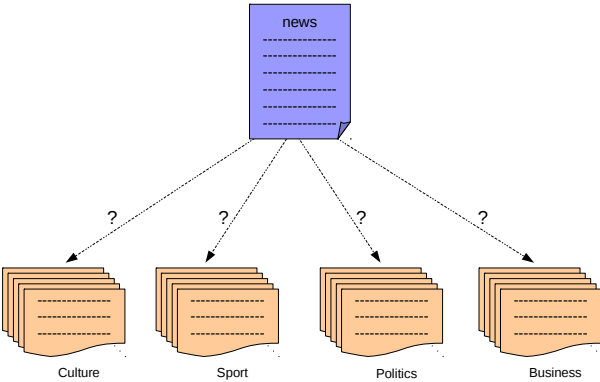
13



The screenshot shows the Google Translate web interface. At the top left is the Google logo. Below it is the 'Translate' section with a 'From: English - detected' dropdown, a language icon, a 'To: German' dropdown, and a blue 'Translate' button. Below the input fields are two text boxes. The left box is labeled 'English' and contains the text 'This is a sample sentence in English which is translated to German.' with a close button (X) and a speaker icon. The right box is labeled 'German' and contains the German translation 'Dies ist ein Beispielsatz in englischer Sprache, die auf deutsch übersetzt wird.' with a speaker icon and a checkmark. Below the right box is a red note: 'New! Click the words above to view alternate translations. Dismiss'.

Text Categorization

14



Information Retrieval

15



Google Search

I'm Feeling Lucky

[Information technology - Wikipedia, the free encyclopedia](#)
en.wikipedia.org/wiki/Information_technology

Information Technology (IT) is concerned with technology to treat information. The acquisition, processing, storage and dissemination of vocal, pictorial, textual ...

↳ Information systems - Information history - Category:Information technology

[Information Technology – All About Information Technology - Wh...](#)
jobsearchtech.about.com/od/careersintechology/p/ITDefinition.htm

Information Technology and IT definition. What **information technology** actually means. How **information technology** is different from computer science.

[RIT Information Sciences & Technology](#)

www.ist.rit.edu/

offers bachelors and masters degrees in **information technology**, a masters degree in software development and management, and an advanced certificate in ...

[ScienceDaily: Information Technology News](#)

www.sciencedaily.com/news/computers.../information_technology/

1 day ago – **Information Technology**. Read the latest in IT research from research institutes around the world. Updated daily, full-text, images, free.

[Government of India, Department of Information Technology \(DIT...](#)

www.mit.gov.in/

Developing the **information technology** industry. Includes an organisation chart, subsidiary bodies.

[Information Technology - Everything You Need to Know](#)

informationtechnology.net/

What is **Information Technology**? **Information Technology**, or IT, is the study, design, creation, utilization, support, and management of computer-based ...

[Information Technology](#)

www.ibef.org/industry/informationtechnology.aspx

The Indian **information technology** (IT) industry has played a key role in putting India on the global map and is now envisioned to become a US\$ 225 billion ...

[Information Technology - WetFeet.com](#)

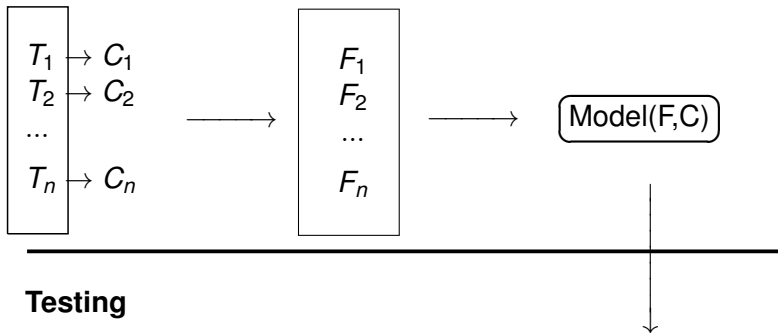
www.wetfeet.com/careers-industries/careers/information-technology

Information Technology. Overview. E-mail, personal computers, and the Internet: These products of the information age have become common currency among ...

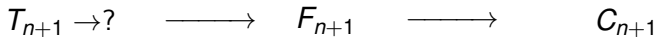
Classification

16

Training



Testing

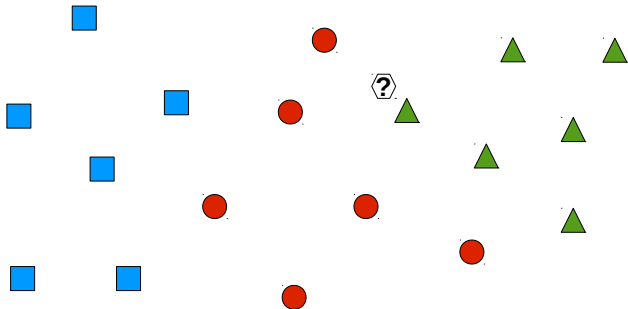


Classification Algorithms

- K Nearest Neighbor
- Support Vector Machines
- Naïve Bayes
- Maximum Entropy
- Linear Regression
- Logistic Regression
- Neural Networks
- Decision Trees
- Boosting
- ...

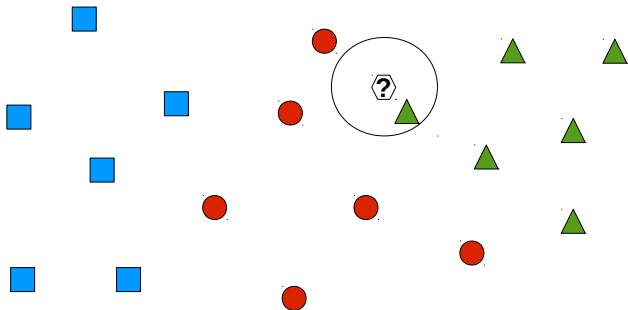
K Nearest Neighbor

18



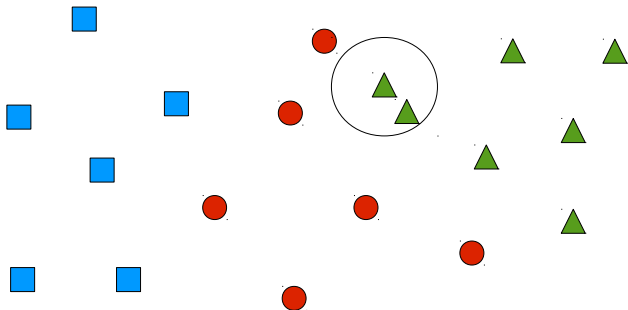
K Nearest Neighbor

19



K Nearest Neighbor

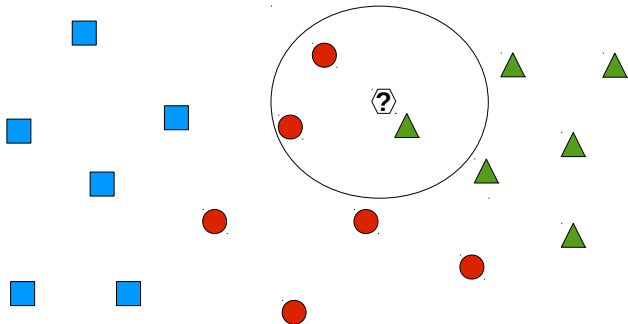
20



- 1 Nearest Neighbor

K Nearest Neighbor

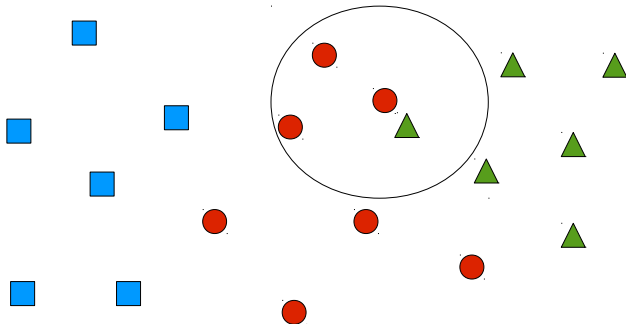
21



- 3 Nearest Neighbor

K Nearest Neighbor

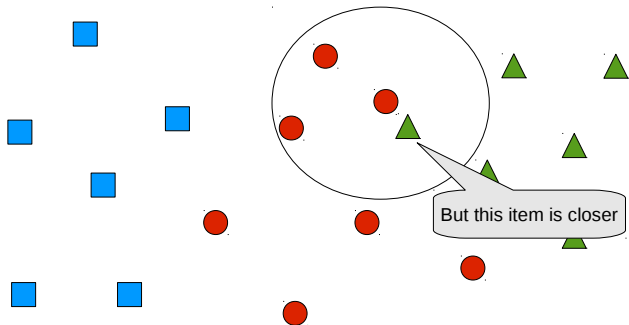
22



- 3 Nearest Neighbor

K Nearest Neighbor

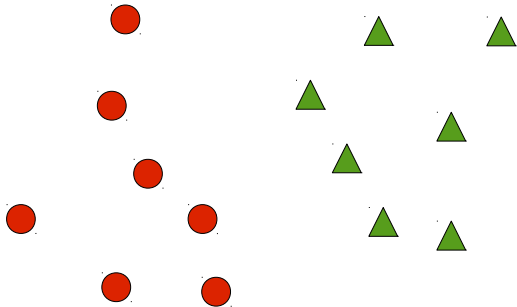
23



- 3 Nearest Neighbor

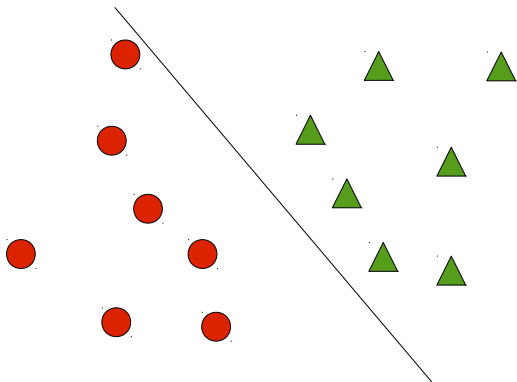
Support Vector Machines

24



Support Vector Machines

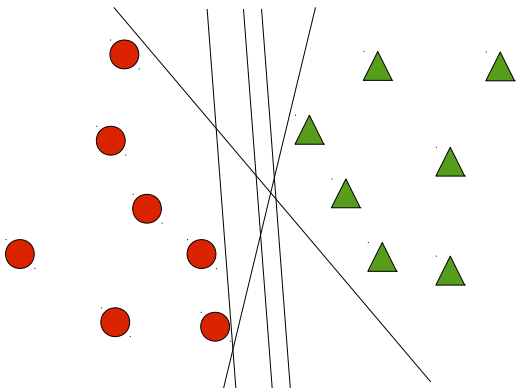
25



- Find a hyperplane in the vector space that separates the items of the two categories

Support Vector Machines

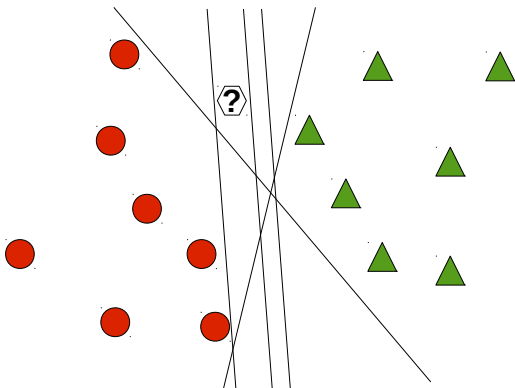
26



- There might be more than one possible separating hyperplane

Support Vector Machines

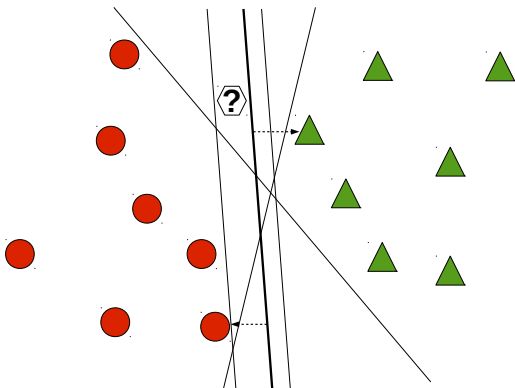
27



- There might be more than one possible separating hyperplane

Support Vector Machines

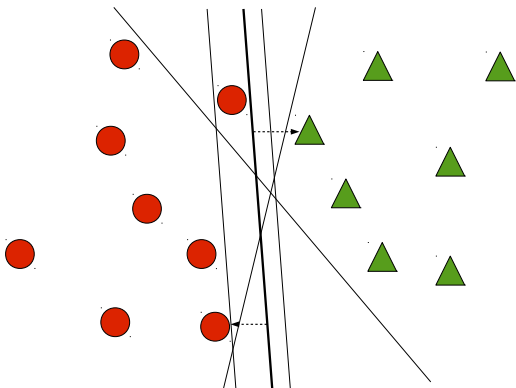
28



- Find the hyperplane with maximum margin
- Vectors at the margins are called support vectors

Support Vector Machines

29



- Find the hyperplane with maximum margin
- Vectors at the margins are called support vectors

Naïve Bayes

30

- Selecting the class with highest probability
⇒ Minimizing the number of items with wrong labels

$$\hat{c} = \operatorname{argmax}_{c_i} P(c_i)$$

- The probability should depend on the to be classified data (d)

$$P(c_i|d)$$

Naïve Bayes

31

$$\hat{c} = \operatorname{argmax}_{c_i} P(c_i)$$

$$\hat{c} = \operatorname{argmax}_{c_i} P(c_i|d)$$

$$\hat{c} = \operatorname{argmax}_{c_i} \frac{P(d|c_i) \cdot P(c_i)}{P(d)}$$

$P(d)$ has no effect

$$\hat{c} = \operatorname{argmax}_{c_i} P(d|c_i) \cdot P(c_i)$$

Naïve Bayes

32

$$\hat{c} = \operatorname{argmax}_{c_j} P(d|c_j) \cdot P(c_j)$$



Likelihood
Probability



Prior
Probability

Maximum Entropy

33

- Assigning a weight λ_j to each feature f_j
 - Positive weight: the feature is likely to be effective
 - Negative weight: the feature is likely to be ineffective
- Picking out a subset of data by each feature
- Voting for each class based on the sum of weighted features

$$\hat{c} = \operatorname{argmax}_{c_i} P(c_i | d, \lambda)$$

Maximum Entropy

34

$$\hat{c} = \operatorname{argmax}_{c_i} P(c_i|d, \lambda)$$

$$P(c_i|d, \lambda) = \frac{\exp \sum_j \lambda_j \cdot f_j(c, d)}{\sum_{c_i} \exp \sum_j \lambda_j \cdot f_j(c_i, d)}$$

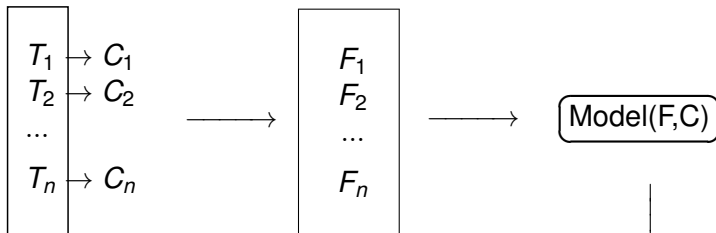
- The expectation of each feature is calculated as follows:

$$E(f_j) = \sum_{(c,d) \in (C,D)} P(c, d) \cdot f_j(c, d)$$

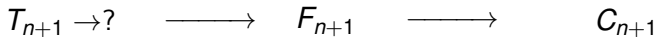
Classification

35

Training

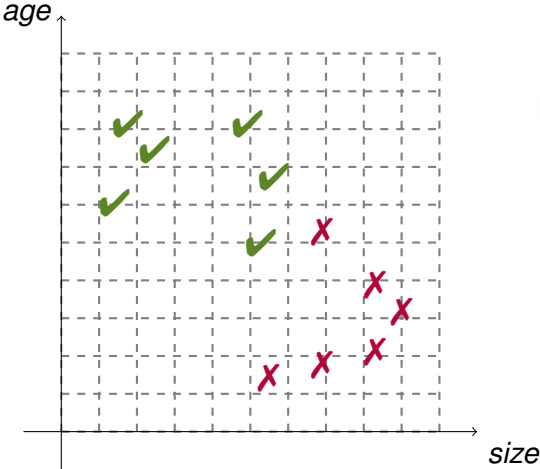


Testing



Feature Selection

36



- Location
- Number of rooms
- Balcony
- ...

Feature Selection

37

- Bag-of-words:
 - Each document can be represented by the set of words that appear in the document
 - Result is a high dimensional feature space
 - The process is computationally expensive

- Solution
 - Using a feature selection method to select informative words

Feature Selection Methods

38

- Information Gain
- Mutual Information
- χ -Square

Information Gain

39

- Measuring the number of bits required for category prediction w.r.t. the presence or absence of a term in the document
- Removing words whose information gain is less than a predefined threshold

$$\begin{aligned}
 IG(w) = & - \sum_{i=1}^K P(c_i) \log P(c_i) \\
 & + P(w) \sum_{i=1}^K P(c_i|w) \log P(c_i|w) \\
 & + P(\bar{w}) \sum_{i=1}^K P(c_i|\bar{w}) \log P(c_i|\bar{w})
 \end{aligned}$$

Information Gain

40

$$P(c_i) = \frac{N_i}{N}$$

$$P(w) = \frac{N_w}{N}$$

$$P(c_i|w) = \frac{N_{iw}}{N_i}$$

$$P(\bar{w}) = \frac{N_{\bar{w}}}{N}$$

$$P(c_i|\bar{w}) = \frac{N_{i\bar{w}}}{N_i}$$

N : # docs

N_i : # docs in category c_i

N_w : # docs containing w

$N_{\bar{w}}$: # docs not containing w

N_{iw} : # docs in category c_i containing w

$N_{i\bar{w}}$: # docs in category c_i not containing w

Mutual Information

41

- Measuring the effect of each word in predicting the category
 - How much does its presence or absence in a document contribute to category prediction?

$$MI(w, c_i) = \log \frac{P(w, c_i)}{P(w) \cdot P(c_i)}$$

- Removing words whose mutual information is less than a predefined threshold

$$MI(w) = \max_i MI(w, c_i)$$

$$MI(w) = \sum_i P(c_i) \cdot MI(w, c_i)$$

χ -square

42

- Measuring the dependencies between words and categories

$$\chi^2(w, c_i) = \frac{N \cdot (N_{iw}N_{i\bar{w}} - N_{i\bar{w}}N_{i_w})^2}{(N_{iw} + N_{i\bar{w}}) \cdot (N_{i_w} + N_{i\bar{w}}) \cdot (N_{iw} + N_{i_w}) \cdot (N_{i\bar{w}} + N_{i\bar{w}})}$$

- Ranking words based on their χ -square measure

$$\chi^2(w) = \sum_{i=1}^K P(c_i) \cdot \chi^2(w, c_i)$$

- Selecting the top words as features

Feature Selection

43

- These models perform well for document-level classification
 - Spam Mail Detection
 - Language Identification
 - Text Categorization

- Word-level Classification might need another types of features
 - POS Tagging
 - Named Entity Recognition

(will be discussed later)

Shortcoming

44

- Data annotation is labor intensive

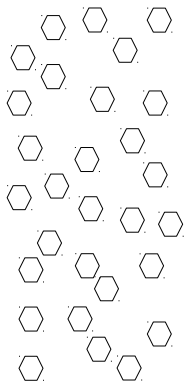


- Solution:
 - Using a minimum amount of annotated data
 - Annotating further data by human, if they are very informative

Active Learning

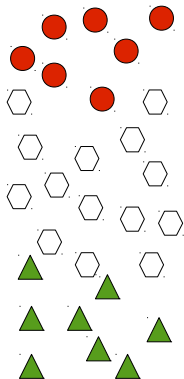
Active Learning

45



Active Learning

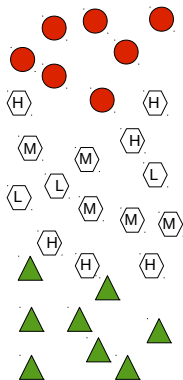
46



- Annotating a small amount of data

Active Learning

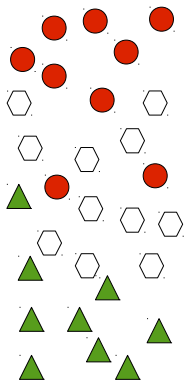
47



- Annotating a small amount of data
- Calculating the confidence score of the classifier on unlabeled data

Active Learning

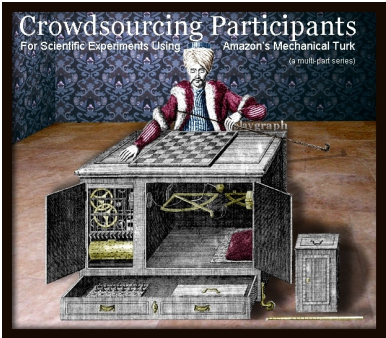
48



- Annotating a small amount of data
- Calculating the confidence score of the classifier on unlabeled data
- Finding the informative unlabeled data (data with lowest confidence)
- Annotating the informative data by the human

Active Learning

Amazon Mechanical Turk



Outline

50

- ① Supervised Learning
- ② Semi-Supervised Learning
- ③ Unsupervised Learning

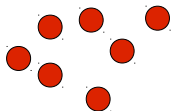
Semi-Supervised Learning

51

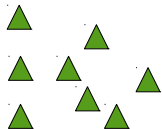
- Problem of data annotation



- Solution:
 - Using minimum amount of annotated data
 - Annotating further data **automatically**, if they are easy to predict

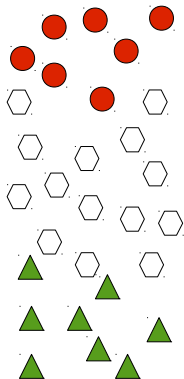


- A small amount of labeled data



Semi-Supervised Learning

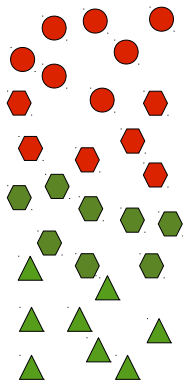
53



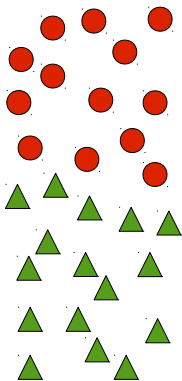
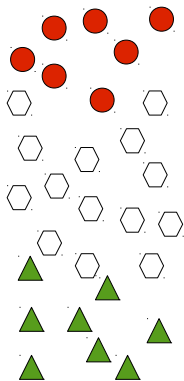
- A small amount of labeled data
- A large amount of unlabeled data

Semi-Supervised Learning

54



- A small amount of labeled data
- A large amount of unlabeled data
- Solution
 - Finding the similarity between the labeled and unlabeled data
 - Predicting the labels of the unlabeled data

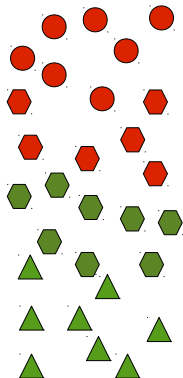


- Training the classifier using
 - The labeled data
 - Predicted labels of the unlabeled data

Shortcoming

56

- Introducing a lot of noisy data to the system



Shortcoming

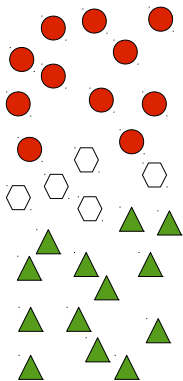
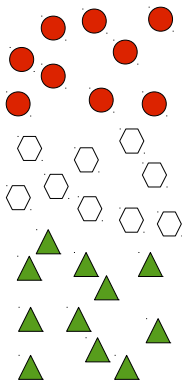
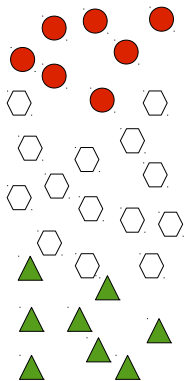
57

- Introducing a lot of noisy data to the system

- Solution
 - Adding unlabeled data to the training set, if the predicted label has a high confidence

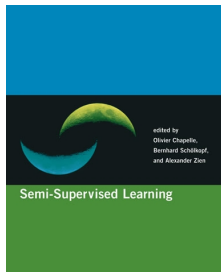
Semi-Supervised Learning

58



Related Books

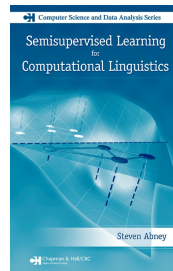
59



Semi-Supervised Learning

by O. Chapelle, B. Schölkopf, A. Zien
MIT Press

2006



Semisupervised Learning for
Computational Linguistics

by S. Abney
Chapman & Hall

2007

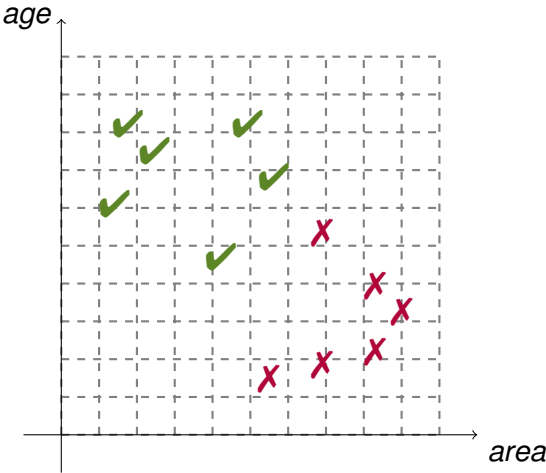
Outline

60

- ① Supervised Learning
- ② Semi-Supervised Learning
- ③ Unsupervised Learning

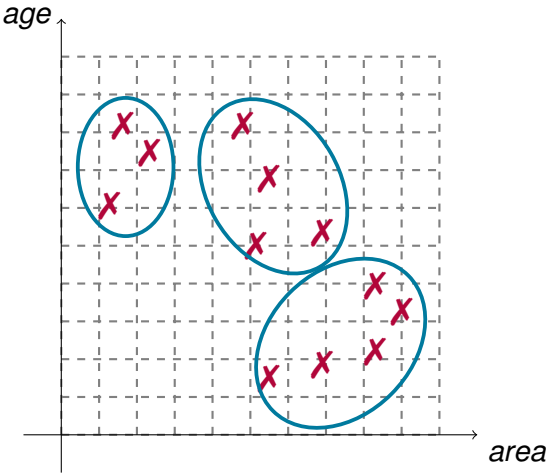
Unsupervised Learning

61



Unsupervised Learning

62



Clustering

- Working based on the similarities between the data items
- Assigning the similar data items to the same cluster

Applications

64

- Word Clustering
 - Speech Recognition
 - Machine Translation
 - Named Entity Recognition
 - Information Retrieval
 - ...

- Document Clustering
 - Information Retrieval
 - ...

Speech recognition

65



“Computers can recognize speech.”

“Computers can wreck a nice peach.”

Machine Translation

66

“The cat eats ...” \Rightarrow *“Die Katze frisst ...”*
“Die Katze isst ...”

Language Modeling

67

Corpus Texts:

“I have a meeting on Monday evening”

“You should work on Wednesday afternoon”

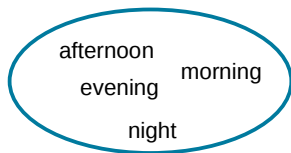
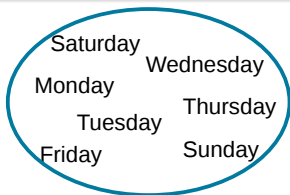
“The next session of the NLP lecture in on Thursday morning”

no observation in the corpus.



*“The talk is on **Monday morning**.”*

“The talk is on Monday molding.”



Language Modeling

68

Corpus Texts:

*“I have a meeting on **Monday evening**”*

*“You should work on **Wednesday afternoon**”*

*“The next session of the NLP lecture in on **Thursday morning**”*

⇒ [Week-day] [day-time]

Class-based Language Model

Information Retrieval

69

Who invented the automobile



*“The first **car** was invented by Karl Benz.”*

“Thomas Edison invented the first commercially practical light.”

“Alexander Graham Bell invented the first practical telephone.”

car

automobile

vehicle

Information Retrieval

70

About 144,000,000 results (0.14 seconds)

Camel

www.camel.com/

R.J. Reynolds Tobacco Company only markets its tobacco products to tobacco consumers who are 21 years of age or older. In order to be eligible to receive ...
You've visited this page 2 times. Last visit: 4/16/12

Camel - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Camel

A **camel** is an even-toed ungulate within the genus *Camelus*, bearing distinctive fatty deposits known as humps on its back. There are two species of **camels**: the ...
↳ Bactrian camel - Dromedary - Australian feral camel - Camel (disambiguation)

Camel (band) - Wikipedia, the free encyclopedia

[en.wikipedia.org/wiki/Camel_\(band\)](http://en.wikipedia.org/wiki/Camel_(band))

Camel are an English progressive rock band formed in 1971. Whilst they didn't achieve the large-scale fame of some of their '70s contemporaries (Pink Floyd, ...

Apache Camel, Index

camel.apache.org/

Apache **Camel** provides support for Bean Binding and seamless integration with popular frameworks such as Spring, Blueprint and Guice. **Camel** also has ...

Welcome to the Official Camel Website

www.camelproductions.com/

Official site with news, tour information, timeline, merchandise and jukebox. Home site of founder Andy Latimer.

Camel Pictures and Facts

fahs.net/camel-pictures-facts/

A comprehensive look at **camels** and their vital role in history. Take a fun quiz, and see how much you learned! Many of the **camel** pictures are also desktop ...

San Diego Zoo's Animal Bytes: Camel

www.sandiegozoo.org/animals/bytes/camel.html

Get accurate animal information about **camels** in an easy-to-read style from the San Diego Zoo's Animal Bytes. Buy tickets online and plan a visit to the Zoo or ...

Camel - Free listening, videos, concerts, stats, & pictures at Last.fm

www.last.fm/music/Camel

Watch videos & listen free to **Camel**: Freefall, Supersteiner & more, plus 58 pictures. **CAMEL** is a progressive rock group from Guildford, Surrey, England.

CAMEL - NU RMX UP!!!! | Free Music, Tour Dates, Photos, Videos

www.myspace.com/camel@nuco

CAMEL - NU RMX UP!!!!'s official profile including the latest music, albums, songs, music videos and more updates.

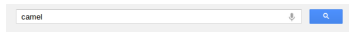
Programming Perl, 3rd Edition - O'Reilly Media

prog.oreilly.com/product/0789596002/1.do

Camels are large ruminant mammals, weighing between 1000 and 1600 ... All this having been said, the **Camel** Book is getting a new edition, slated (as of this ...

Information Retrieval

71



About 144,000,000 results (0.14 seconds)

Camel

www.camel.com/

R.J. Reynolds Tobacco Company only markets its tobacco products to tobacco consumers who are 21 years of age or older. In order to be eligible to receive ...
You've visited this page 2 times. Last visit: 4/16/12

Camel - Wikipedia, the free encyclopedia

en.wikipedia.org/wiki/Camel

A **camel** is an even-toed ungulate within the genus *Camelus*, bearing distinctive fatty deposits known as humps on its back. There are two species of **camels**: the ...
↳ Bactrian camel - Dromedary - Australian feral camel - Camel (disambiguation)

Camel (band) - Wikipedia, the free encyclopedia

[en.wikipedia.org/wiki/Camel_\(band\)](http://en.wikipedia.org/wiki/Camel_(band))

Camel are an English progressive rock band formed in 1971. Whilst they didn't achieve the large-scale fame of some of their '70s contemporaries (Pink Floyd, ...

Apache Camel - Index

camel.apache.org/

Apache Camel provides support for Bean Binding and seamless integration with popular frameworks such as Spring, Blueprint and Guice. **Camel** also has ...

Welcome to the Official Camel Website

www.camelproductions.com/

Official site with news, tour information, timeline, merchandise and jukebox. Home site of founder Andy Latimer.

Camel Pictures and Facts

fahs.net/camel-pictures-facts/

A comprehensive look at **camels** and their vital role in history. Take a fun quiz, and see how much you learned! Many of the **camel** pictures are also desktop ...

San Diego Zoo's Animal Bytes: Camel

www.sandiegozoo.org/animalbytes/animal-camel.html

Get accurate animal information about **camels** in an easy-to-read style from the San Diego Zoo's Animal Bytes. Buy tickets online and plan a visit to the Zoo or ...

Camel - Free listening, videos, concerts, stats, & pictures at Last.fm

www.last.fm/music/Camel

Watch videos & listen free to **Camel**. Freefall, Supersteiner & more, plus 58 pictures. **CAMEL** is a progressive rock group from Gulkiford, Surrey, England.

CAMEL - NU RMX UP!!!! | Free Music, Tour Dates, Photos, Videos

www.myspace.com/camelbandco

CAMEL - NU RMX UP!!!!'s official profile including the latest music, albums, songs, music videos and more updates.

Programming Perl, 3rd Edition - O'Reilly Media

shop.oreilly.com/product/078959600271.do

Camels are large ruminant mammals, weighing between 1000 and 1600 ... All this having been said, the **Camel** Book is getting a new edition, slated (as of this ...



Clustering documents based on their similarities

Clustering Algorithms

72

- Flat
 - K-means

- Hierarchical
 - Top-Down (Divisive)
 - Bottom-Up (Agglomerative)
 - Single-link
 - Complete-link
 - Average-link

K-means

73

- The best known clustering algorithm
- Works well for many cases
- Used as default / baseline for clustering documents
- Algorithm
 - Defining each cluster center as the mean or centroid of the items in the cluster

$$\vec{\mu} = \frac{1}{|c|} \sum_{\vec{x} \in c} \vec{x}$$

- Minimizing the average squared Euclidean distance of the items from their cluster centers

K-means

74

```
Initialization: Randomly choose  $k$  items as initial centroids  
while stopping criterion has not been met do  
  for each item do  
    Find the nearest centroid  
    Assign the item to the cluster associated with the nearest centroid  
  end for  
  for each cluster do  
    Update the centroid of the cluster based on the average of all items in the cluster  
  end for  
end while
```

■ Iterating two steps:

□ Re-assignment

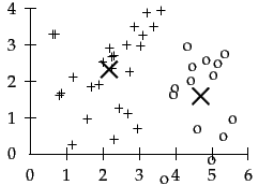
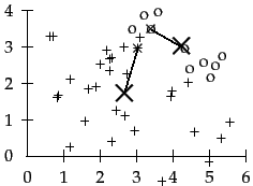
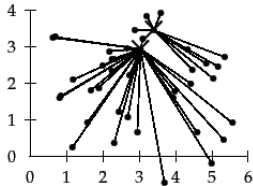
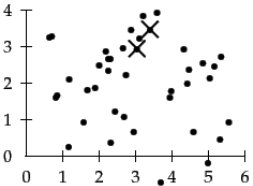
- Assigning each vector to its closest centroid

□ Re-computation

- Computing each centroid as the average of the vectors that were assigned to it in re-assignment

K-means

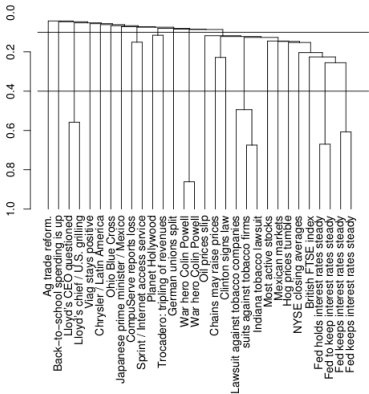
75



http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletKM.html

Hierarchical Agglomerative Clustering (HAC)

- Creating a hierarchy in the form of a binary tree



Hierarchical Agglomerative Clustering (HAC)

77

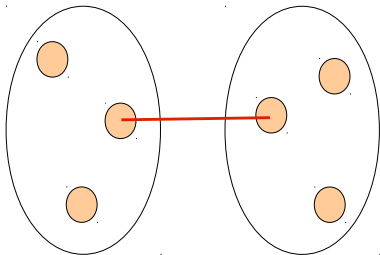
Initial Mapping: Put a single item in each cluster
while reaching the predefined number of clusters **do**
 for each pair of clusters **do**
 Measure the similarity of two clusters
 end for
 Merge the two clusters that are most similar
end while

- Measuring the similarity in three ways:
 - Single-link
 - Complete-link
 - Average-link

Hierarchical Agglomerative Clustering (HAC)

78

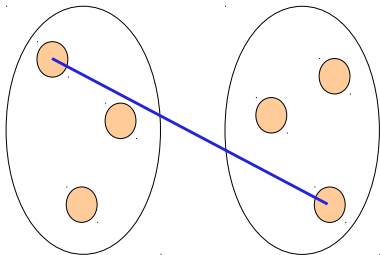
- Single-link / single-linkage clustering
 - Based on the similarity of the most **similar** members



Hierarchical Agglomerative Clustering (HAC)

79

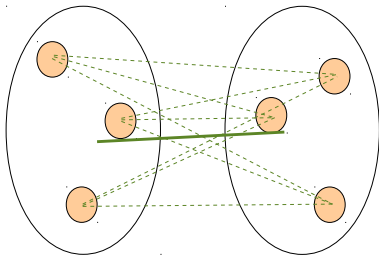
- Complete-link / complete-linkage clustering
 - Based on the similarity of the most **dissimilar** members



Hierarchical Agglomerative Clustering (HAC)

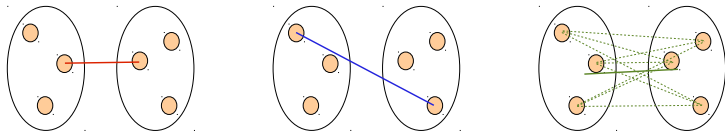
80

- Average-link / average-linkage clustering
 - Based on the average of all similarities between the members



Hierarchical Agglomerative Clustering (HAC)

81



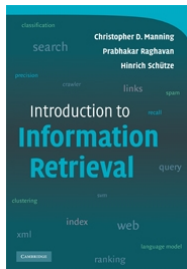
http://home.dei.polimi.it/matteucc/Clustering/tutorial_html/AppletH.html

Further Reading

82

Introduction to Information Retrieval

C.D. Manning, P. Raghavan, H. Schütze Cambridge
University Press 2008



[http://nlp.stanford.edu/IR-book/html/
htmledition/irbook.html](http://nlp.stanford.edu/IR-book/html/htmledition/irbook.html)

Chapters 13,14,15,16,17