# Information Extraction
Taxonomies, classification, detection and linking

## Web Data Mining and Search

# Outline

- Introduction

- From data to information: Taxonomies and classes

- Tasks:
  - Classification
  - Detection
  - Recognition
  - Linking
  - Relation extraction

# Importance of information extraction

P. Jackson and I. Moulinier. 2002. *Natural Language Processing for Online Applications*

- "There is no question concerning the commercial value of being able to classify documents automatically by content. There are myriad potential applications of such a capability for corporate intranets, government departments, and Internet publishers"

- "Understanding the data is one of the keys to successful categorization, yet this is an area in which most categorization tool vendors are extremely weak. Many of the 'one size fits all' tools on the market have not been tested on a wide range of content types."
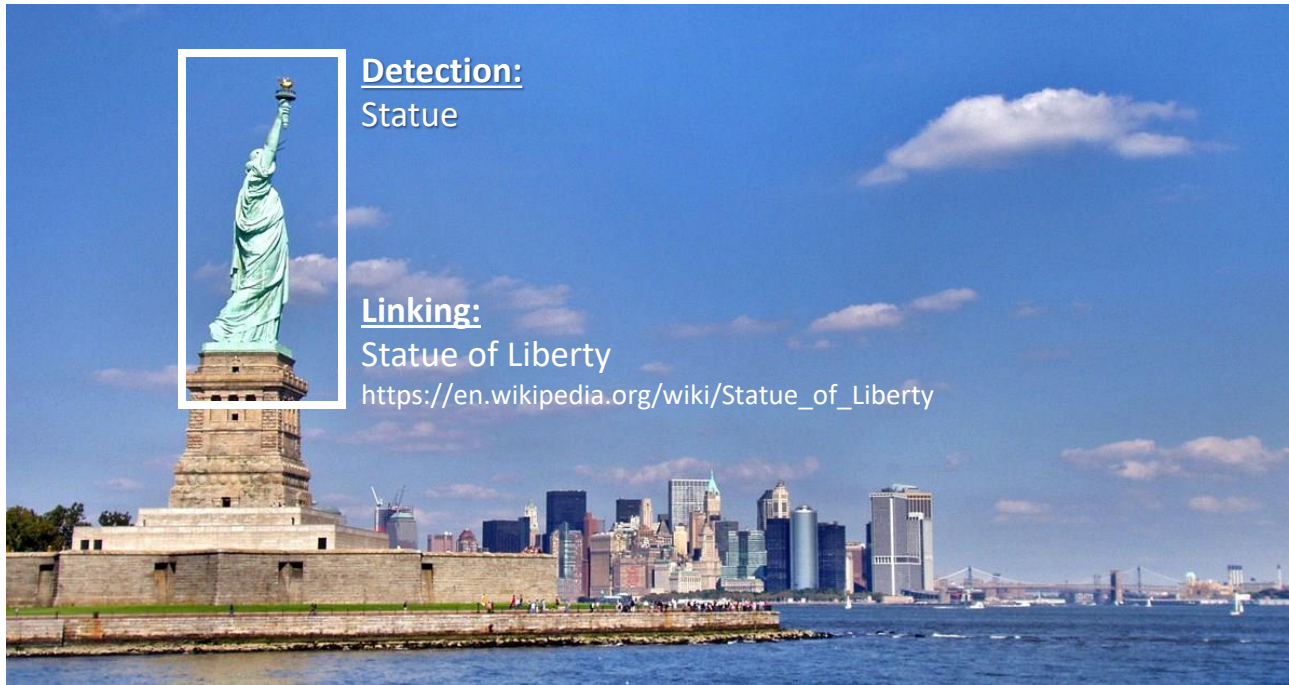
# Real world tasks

- SPAM detection (fake opinions)
- Memes detection (not informational)
- Tampered images
- Sentiment detection (opinions)
- Emergency detection

# Classification, detection, linking



**Detection:**
Statue

**Linking:**
Statue of Liberty
https://en.wikipedia.org/wiki/Statue_of_Liberty
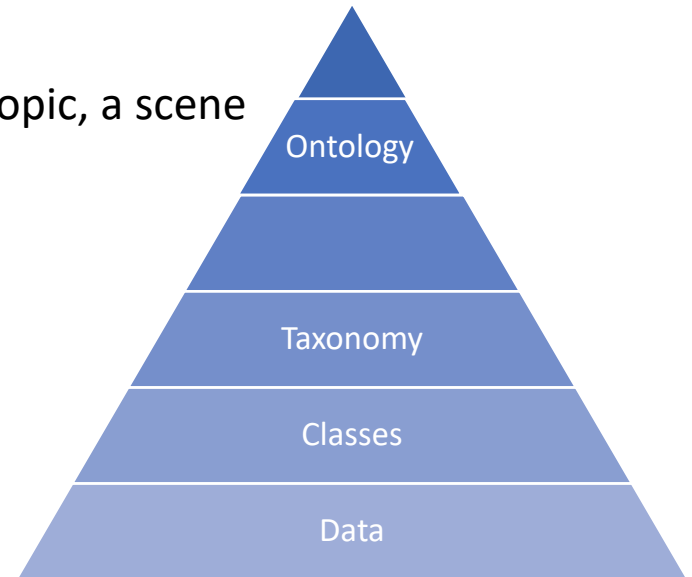
**Classification:**
Sea side
Statue
City
Sky

**Linking:**
New York City
https://en.wikipedia.org/wiki/New_York_City

# From data to information
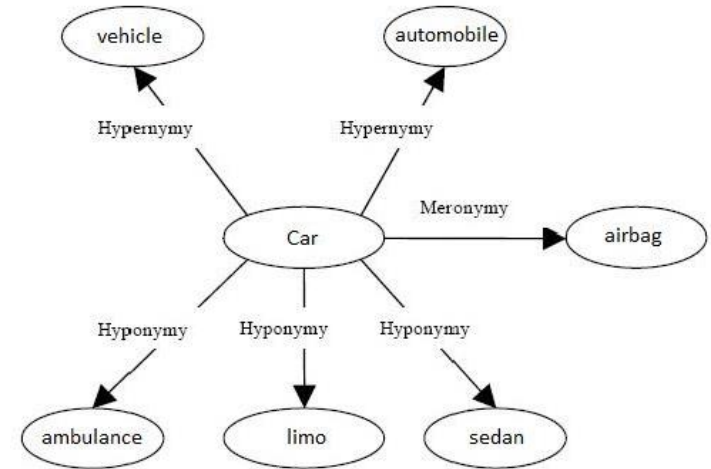
Web Data Mining and Search

# From data to information

- A taxonomy is concerned with classifying and organizing hierarchically concepts of a specific domain.

- It is important to identify the list of items that need to be detected.
  - These items are domain specific, and can be a topic, a scene type, a visual object or a named entity.
  - They are normally associated to a class in a supervised learning task.

Ontology

Taxonomy

Classes

Data

# WordNet: A lexical database

"WordNet® is a large lexical database of English. Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept. Synsets are interlinked by means of conceptual-semantic and lexical relations. "

"WordNet interlinks specific senses of words. As a result, words that are found in close proximity to one another in the network are semantically disambiguated. Second, WordNet labels the semantic relations among words, whereas the groupings of words in a thesaurus does not follow any explicit pattern other than meaning similarity."



https://wordnet.princeton.edu/

# ImageNet: A visual taxonomy

- Selected words of WordNet are illustrated in ImageNet.

- Currently, there are over 14.000 concepts illustrated.

- Roughly 1.000 concepts are used by VOC.

- Great impact in advancing the state of the art.

http://image-net.org/explore.php

# Domain specific taxonomies

- Domain specific terminologies are curated by domain experts and are designed with specific tasks and workflows in mind.

- In the medical domain, the SNOMED-CT is intended to describe medical conditions, procedures, admin, etc.
    - http://browser.ihtsdotools.org/

- In the computer science domain the ACM Computing Classification Scheme is widely used to classify published articles.
    - https://dl.acm.org/ccs/ccs.cfm

# Wikipedia as a database

- Wikipedia contains large amounts of information largely unstructured but structured as a taxonomy.

- **DBPedia** aims to create a rigorous database out of Wikipedia.

https://en.wikipedia.org/wiki/Portal:Contents

- A key application is to link data to Wikipedia entries.

# Which and how many are detectable?

- An important question to ask is which and how many items of the taxonomy are detectable in data?

- A few (well separated ones)?          -> Easy!

- A zillion closely related ones?         -> Not so easy…
  - Think: Yahoo! Directory, Library of Congress classification, legal applications
  - Quickly gets difficult!
    - Classifier combination is always a useful technique
      - Voting, bagging, or boosting multiple classifiers
    - Much literature on hierarchical classification
      - Definitely helps for scalability, even if not in accuracy
    - May need a hybrid automatic/manual solution

# Taxonomies and classification

- In practice, only a few elements of the taxonomy should be used as classes for classification
  - Only the ones offering a stable document class representation.

- The ultimate goal is to link information to an entry on a taxonomy capturing the target domain.

- Ultimately more complete domain representation should be used, e.g. an ontology.

# Classification

Web Data Mining and Search

# Document classification



| Colour | Texture | Region | Semantic | Metadata |
|---|---|---|---|---|
| | | | windmill, sky, sea,buildings | Date: 7 Dec 06 Author: Joao, Place: Portugal |

Input features  **Classification**  class

# Classification task

- For new unseen documents, we wish to classify documents with one of the known classes.

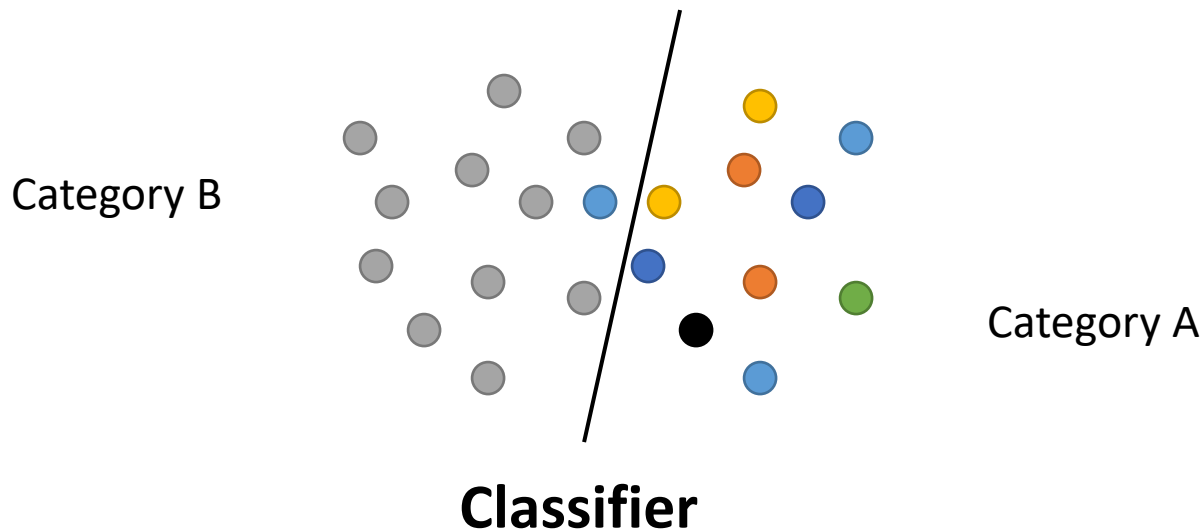- New documents are represented in some feature space and then a machine learning algorithm classifies the new documents.

Category B

Category A

**Classifier**

# Perceptron

- All sample vectors $x^{(j)}$ have their corresponding label $y^{(j)} = \{+1, -1\}$

- **The perceptron performs a binary prediction $\hat{y}$ based on the observed data $x$ :**

$$\hat{y} = f(x) = \begin{cases} +1 & , if\ x_2 \geq m \cdot x_1 + b \\ -1 & , if\ x_2 < m \cdot x_1 + b \end{cases}$$

# Model error

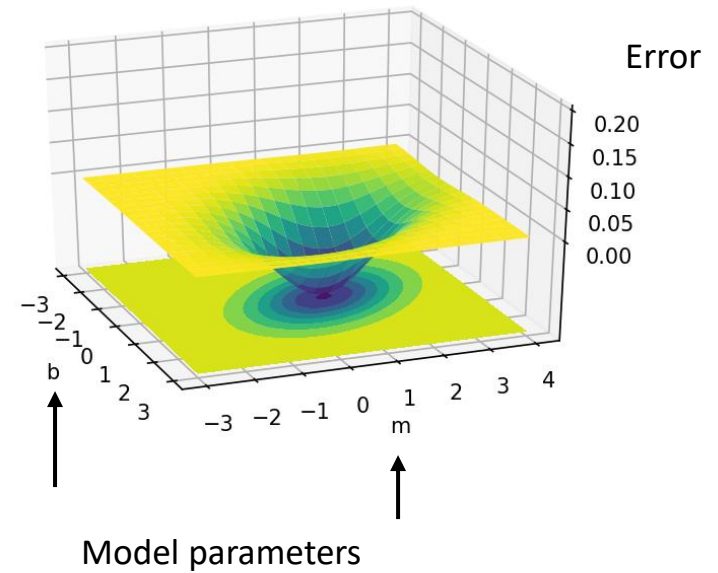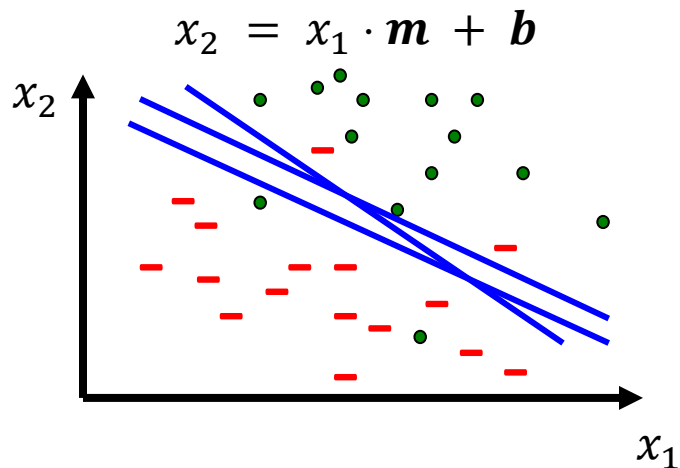- The Mean Square Error (MSE) measures the error between the true labels and the predicted labels

$$MSE = \frac{1}{TotalSamples} \sum_{i}^{TotalSamples} (label_i - predictedLabel_i)^2$$

$$x_2 - m \cdot x_1 - b = 0$$

# Minimizing the error

$$MeanSquareError = \frac{1}{TotalSamples} \sum_{i}^{TotalSamples} (label_i - predictedLabel_i)^2$$

$$x_2 = x_1 \cdot \boldsymbol{m} + \boldsymbol{b}$$

Error

Model parameters

# Learning to minimize the model error

- Initialize the model with random weights
- Compute the model predictions
- Compute the error of each prediction
- Update the model with the <u>samples incorrectly classified</u>.

| Observation | Prediction | Error | Update |
|:---:|:---:|:---:|:---:|
| -1 | -1 | 0 | 0 |
| -1 | +1 | -1 | -1*x |
| +1 | -1 | +1 | +1*x |
| +1 | +1 | 0 | 0 |

# Learning algorithm

```
[ ]:  b=0
      m=0
      model = [m,b]

      max_iters = 30
      mean_square_error = []
      for iter in range(0,max_iters):

          # Compute the model predictions
          predicted_labels = ((observations_x2 - m*observations_x1 - b ) >= 0)*2-1

          # Compute the model error
          error_of_all_samples = (true_labels-predicted_labels)/2

          # Update the model parameters
          update_m = np.mean(error_of_all_samples*observations_x1)
          update_b = np.mean(error_of_all_samples)

          m = m - update_m*0.1
          b = b - update_b*0.1
```

$$\hat{y} = f(x) = \begin{cases} +1 & ,if\ x_2 - m \cdot x_1 - b \geq 0 \\ -1 & ,if\ x_2 - m \cdot x_1 - b < 0 \end{cases}$$

$$error = (y - \hat{y})/2 = \begin{cases} +1 \\ 0 \\ -1 \end{cases}$$

$$update_m = error \cdot x_1$$

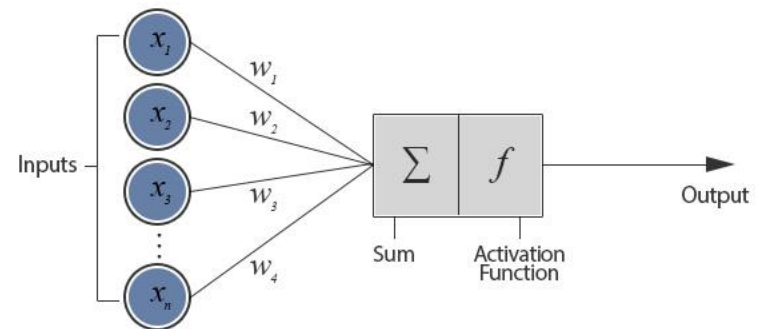$$m = m - update_m \cdot learning_{rate}$$

# Perceptron: general formulation

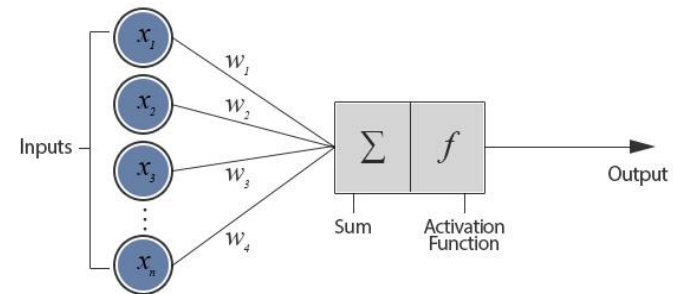- **Binary classification:**

$$z = w_0 + w_1 x_1 + \ldots + w_n x_n$$

$$\hat{y} = f(z) = \begin{cases} +1 & , if\ z \geq 0 \\ -1 & , if\ z < 0 \end{cases}$$



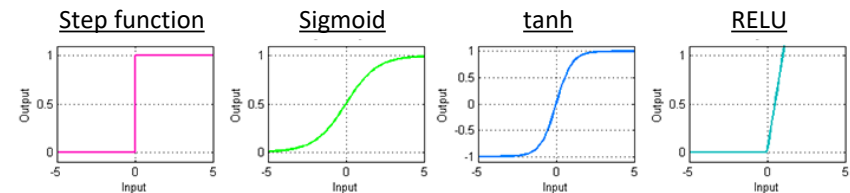- **Input:** Vectors $x^{(j)}$ and labels $y^{(j)}$
  - Vectors $x^{(j)}$ are real valued where $\|x\|_2 = 1$
- **Goal:** Find vector $w = (w_1, w_2, \ldots, w_d)$
  - Each $w_i$ is a real number

# Activation functions

- The perceptron was initially proposed with the step function.

- Historically, other activation functions have been studied.

- It can be shown that the perceptron with the sigmoid activation function corresponds to the logistic regression model.



Inputs $x_1$, $x_2$, $x_3$, ..., $x_n$ with weights $w_1$, $w_2$, $w_3$, $w_4$ → $\Sigma$ (Sum) | $f$ (Activation Function) → Output

**Activation functions**



Step function · Sigmoid · tanh · RELU

# Note regarding model training

- Robustly training a model for Web data is a complex task.

- In most of the cases, we will use pre-trained models.

- These models were trained on large-scale data.

- These pre-trained models are robust and reliable.

# Per-class evaluation measures

| | | Ground-truth | |
|---|---|---|---|
| | | True | False |
| **Method** | True | True positive | False positive |
| | False | False negative | True negative |

- **Recall**: Fraction of docs in class i classified correctly:

$$Recall = \frac{truePos}{truePos + falseNeg}$$

- **Precision**: Fraction of docs assigned class i that are actually about class i:

$$Precision = \frac{truePos}{truePos + falsePos}$$

- **Accuracy**: Fraction of docs classified correctly:

$$Accuracy = \frac{truePos + trueNeg}{truePos + falsePos + trueNeg + falseNeg}$$

\* abragência, precisão e exatidão.

# Micro- vs. Macro-Averaging

- If we have more than one class, how do we combine multiple performance measures into one quantity?

- **Macroaveraging**: Compute performance for each class, then average.

- **Microaveraging**: Collect decisions for all classes, compute contingency table, evaluate.

# Micro- vs. Macro-Averaging: Example

| Class 1 | | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 10 | 10 |
| Classifier: no | 10 | 970 |

| Class 2 | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 90 | 10 |
| Classifier: no | 10 | 890 |

| Micro Ave. Table | | |
|---|---|---|
| | Truth: yes | Truth: no |
| Classifier: yes | 100 | 20 |
| Classifier: no | 20 | 1860 |

- Macroaveraged precision: $(0.5 + 0.9)/2 = 0.7$
- Microaveraged precision: $100/120 = .83$

- Microaveraged score is dominated by score on common classes

# Good practice: Make a confusion matrix

- This (i, j) entry means 53 of the docs actually in class i were put in class j by the classifier.



- In a perfect classification, only the diagonal has non-zero entries
- Look at common confusions and how they might be addressed

# Detection and recognition

Web Data Mining and Search
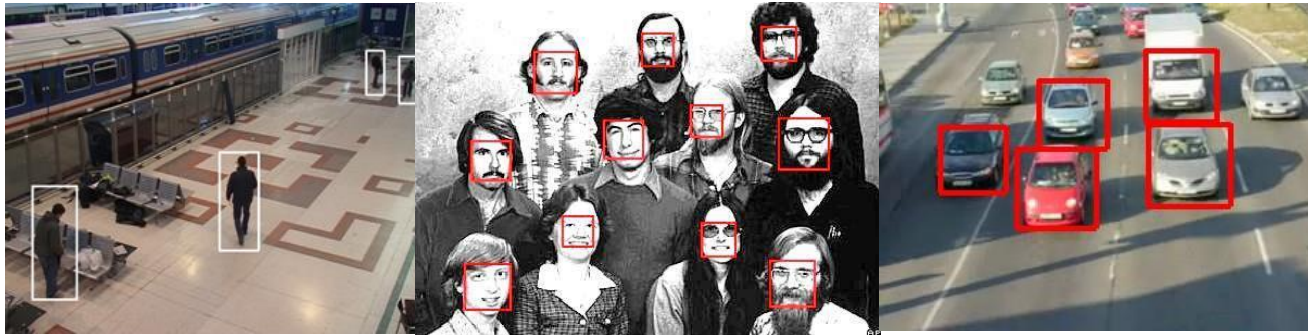
# Detection and recognition

- Detecting and recognizing "things" in natural language and images is a necessary first step in many more complex tasks.

- The "things" that can be detected/recognized include:

| Type | Tag | Sample Categories | Example sentences |
|------|-----|-------------------|-------------------|
| People | PER | people, characters | **Turing** is a giant of computer science. |
| Organization | ORG | companies, sports teams | The **IPCC** warned about the cyclone. |
| Location | LOC | regions, mountains, seas | The **Mt. Sanitas** loop is in **Sunshine Canyon**. |
| Geo-Political Entity | GPE | countries, states, provinces | **Palo Alto** is raising the fees for parking. |
| Facility | FAC | bridges, buildings, airports | Consider the **Golden Gate Bridge**. |
| Vehicles | VEH | planes, trains, automobiles | It was a classic **Ford Falcon**. |

**Figure 18.1**   A list of generic named entity types with the kinds of entities they refer to.

# Object detection

- How to detect a <u>face</u>, a <u>person</u> or a <u>car</u> in a picture?

- How to find pictures of <u>BigBen</u> or the <u>Eiffel Tower</u>?

# Object detection

- To solve the problem of finding objects at multiple scales:
  - The image is scalled multiple times
  - At each scale the entire image is scanned for faces on each possible position.

- This is the convolution operation that we will study later in the course.

# Named entity recognition

- Recognizing a named entity is an important task to extract the meaning of a sentence or a natural language document.

- Ambiguity can exist in the form of polysemy and synonyms.

| Name | Possible Categories |
|---|---|
| *Washington* | Person, Location, Political Entity, Organization, Vehicle |
| *Downing St.* | Location, Organization |
| *IRA* | Person, Organization, Monetary Instrument |
| *Louis Vuitton* | Person, Organization, Commercial Product |

**Figure 18.2** Common categorical ambiguities associated with various proper names.

[PER Washington] was born into slavery on the farm of James Burroughs.
[ORG Washington] went up 2 games to 1 in the four-game series.
Blair arrived in [LOC Washington] for what may well be his last state visit.
In June, [GPE Washington] passed a primary seatbelt law.
The [VEH Washington] had proved to be a leaky ship, every passage I made...

**Figure 18.3** Examples of type ambiguities in the use of the name *Washington*.

# Named entity recognition

- To detect or recognize a named entity, one needs to run a classifier over the sequence of tokens of the natural language sentence.



**Figure 18.7**   Named entity recognition as sequence labeling. The features available to the classifier during training and classification are those in the boxed area.

# Face detection



- As done previously for classification...
  - Positive and negative examples must be gathered
  - Features must be computed for each image
  - A classifier must be estimated

- Positive examples should cover a wide variation of poses, illuminations, instances, etc.

- Challenges:
  - Where is the face?
  - What's the face size?
  - Given an image patch, how to classify the patch as a face or not?

# Linking

Web Data Mining and Search

# Entity linking

- Entity linking concerns the task of linking a mention or an image region to the unique identifier of the entity represented in that piece of data.

- "The task of entity linking is to associate an occurrence in text with the representation of some real-world entity in an ontology, a list of entities in the world, like a gazeteer."

- "Perhaps the most common ontology used for this task is Wikipedia, in which each Wikipedia page acts as the unique id for a particular entity."

# Person recognition

- Once a face is detected, the goal is to link the face image to the person.
  - Ideally linking the face to some named entity in a taxonomy.

- The image face needs to be classified into one of the existing classes, i.e. one of the known persons.

# Summary

- Information extracton tasks
  - https://web.stanford.edu/~jurafsky/slp3/18.pdf

- From data to information: Taxonomies and classes
  - WordNet, ImageNet, SNOMED-CT

- Linear classifier:
  - http://d2l.ai/chapter_linear-networks/index.html