

Topics in Natural Language Processing

Shay Cohen

Institute for Language, Cognition and Computation

University of Edinburgh

Lecture 1

Topics in NLP

- We will cover the basic methodology in NLP
- There will be a focus on statistical learning
- Even more so, *structured prediction*

Topics in NLP

Prerequisites:

- Some familiarity with machine learning and probability
- If something is unclear, ask!

Things to Do:

- Student presentations (20%)
- Brief paper responses (25%)
- Essay (55%)

Office hours: By appointment

NLP Now

late 1980s until now: statistical learning



Learning

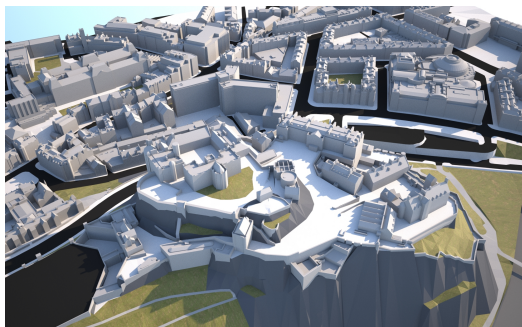
Learning is:

- Experience translated into expertise/knowledge
- Memorisation with generalisation

Machine learning and NLP:

- Experience = Training data
- Knowledge = Decoder or Prediction Model
- Used to either **mimic** humans or **transcend** their abilities

What is a Model?



From Merriam-Webster:

- a usually small copy of something
- a set of ideas and numbers that describe the past, present, or future state of something (such as an economy or a business)

When is a model a good model?

What is a Statistical Model?

Predict the future. Probabilistically.



Probability and Statistics: Reminder

Probability distribution? Example: unigram model

$\Omega = \{ \text{the, cat, chased, on, a, ...} \}$ sample space

$$p(\omega) \geq 0 \quad \sum_{\omega} p(\omega) = 1$$

Random variables

Random variable:

$$X(\omega) = \Omega_2$$

$$X(\text{jumping}) = \text{ing}$$

$$X(\text{jumped}) = \text{ed}$$

$X(\omega)$ returns the last two letters

$$p(X(w) = ed) = \sum_{\substack{w \text{ ends} \\ \text{in } ed}} p(w)$$

$$p(X(w) = 1) = \sum_{\substack{X(w) = 1 \\ w \text{ is related} \\ \text{to a movie}}} p(w) = \sum_w p(w) X(w) = E[X]$$

$X(w)$ is 1

if the word is related to movies

$X(w)$ is 0 o/w

Model Family

Produced with a Trial Version of PDF Annotator - www.PDFAnno

A set of probability distributions (unigram example):

$$\{ p_1(\omega), p_2(\omega), p_3(\omega), \dots \}$$

Parameters

A set of parameters:

$$\Theta$$

$p(w|\theta)$ = probability that depends on θ

$$\{ p(w|\theta) \mid \theta \in \Theta \}$$

$$p_1(w) \dots p(w|\theta_1)$$

$$p_2(w) \dots p(w|\theta_2) \dots$$

$$\theta \in \mathbb{R}^{|S|}$$

$$\theta_i \in (0, 1)$$

$$\sum_i \theta_i = 1$$

Another Parametrisation

Produced with a Trial Version of PDF Annotator - www.PDFAnno

Rely on properties of the words:

We can build paramterisations
that treat words as "groups"

Estimation

What is training data?

$$w_1, w_2, w_3, \dots, w_n$$

Estimation

What is the fit of the data to the model?

$f(w_1, \dots, w_n, \theta)$ tells us whether θ
predicts well w_1, \dots, w_n

$$\theta^* = \underset{\theta}{\operatorname{arg\,max}} f(w_1, \dots, w_n, \theta)$$

our final model

NLP Problem Example: Document Classification

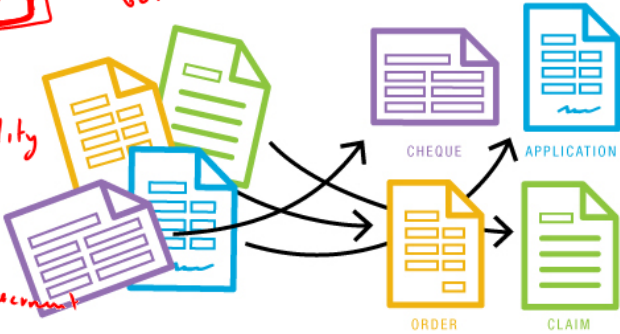
sentiment analysis, document topic, ...

$$p(c|d) = \frac{p(d,c)}{p(d)}$$

$p(d,c)$ ← our basic component

$p(c|d)$
↓
probability of class c for document d

$$p(d) = \sum_c p(d,c) = \sum_c p(c|d)$$



$$\Omega = D \times C$$

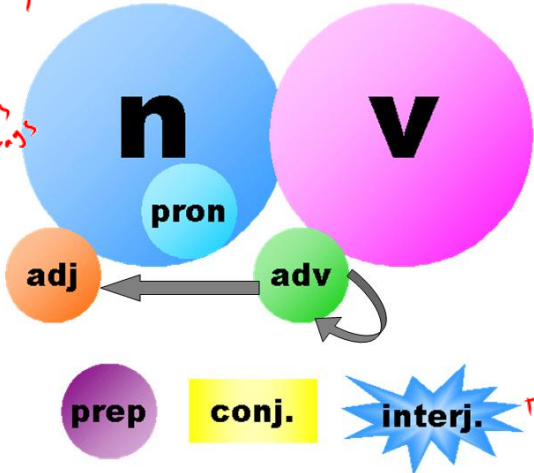
$(d,c) \in \Omega$ d is the document c - the class

NLP Problem Example: POS Tagging

map words to their part-of-speech tags

$$\Omega = V^* \times T^*$$

(v, t)
 \downarrow
 sentence \rightarrow pos tags

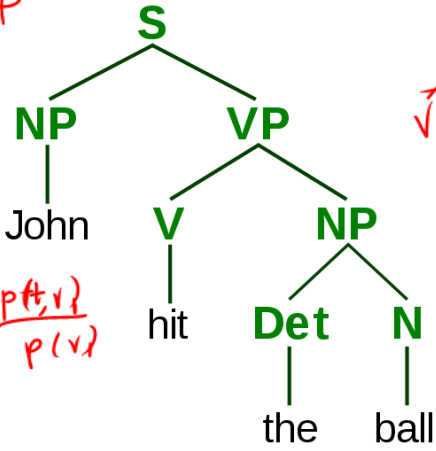
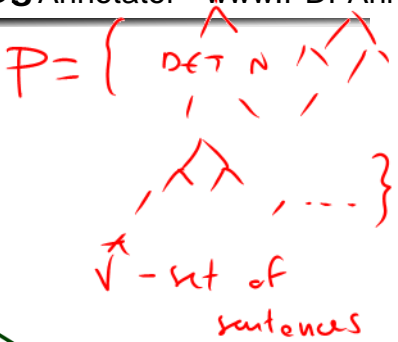


$\{Noun, VERB, PREP, ADJ, ADV\}$
 V - set of words
 $\Omega = T^* \times V$
 $T = \{t_1, t_2, \dots\}$

$$T^* = \{t_1, t_1, t_2, t_3, t_2, t_1, \dots\}$$

map sentences to their syntax

$$\Omega = V^* \times P$$



$$\frac{p(t|v)}{p(v,t)}$$

$$\frac{p(t,v)}{p(v)}$$



Back to Modelling

What if the space to model is complex? Modelling documents.

Modelling a Problem

- Define a sample space
- Define the structure of the sample space
- Decide on a parametrisation

Then one can proceed with data collection and learning

Modelling - Tradeoffs

- “Exact copy”, detailed
- Not too many parameters
- Efficient to work with

Next class

Paradigms in statistical learning

- Frequentist approaches
- Bayesian approaches
- “Computer science approaches?”