

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 (15%)
- Assignment (10%)
- Essay (55%)

Office hours: By appointment

NLP in the Old Days

1950s-1980s: handwritten rules



IBM'S WATSON (right) AND FRIENDS:† For a mathematical wizard . . .

NLP Now

late 1980s until now: statistical learning



Learning

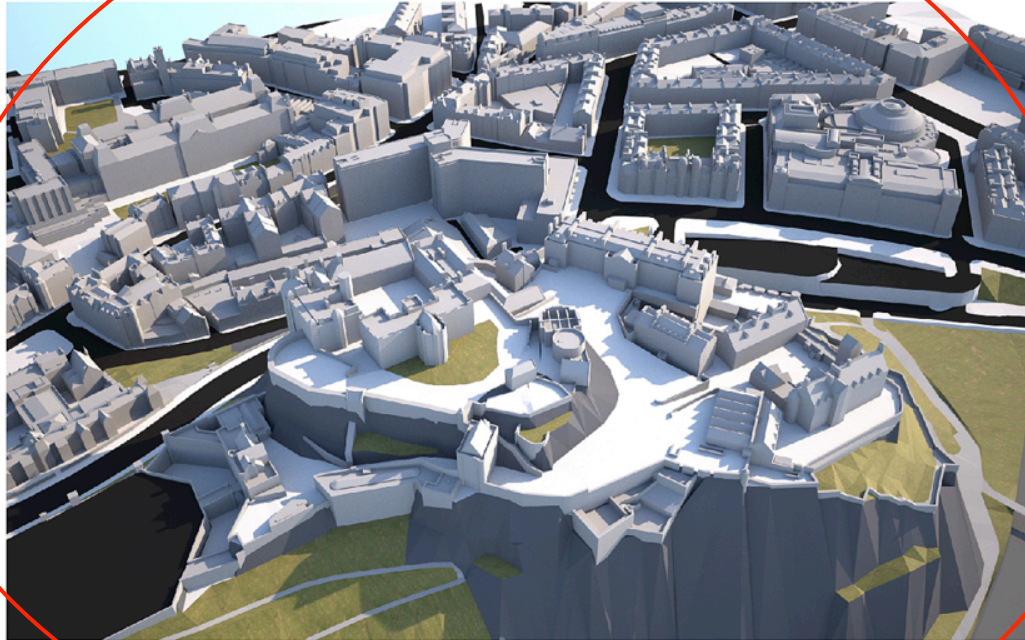
Learning is:

- 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

$$\Omega = \{ \text{the, cat, dog, sit, chase} \}$$

$p: \Omega \rightarrow [0, 1]$ - $p(\omega)$ is the probability attached to ω

$$p(\omega) \geq 0$$

$$\sum_{\omega \in \Omega} p(\omega) = 1$$

$$\int_{\omega} p(\omega) d\omega = 1$$

Random variables

Random variable: $X: \Omega \rightarrow \mathbb{R}$

$$\Omega = \{the, dog, cat, \dots\}$$

$$X_a(\omega) = \text{count of } a\text{'s in } \omega$$

$$X_a(the) = 0 \quad X_a(cat) = 1$$

$$\Omega_2 = \{-ed, -ing, -ion\}$$

$$X(\omega) = \text{suffix of the word}$$

$$X: \Omega \rightarrow \Omega_2$$

Random variables induce probability distribution:

$P(X = \text{ion}) =$ the probability of a word w ending in $-\text{ion}$

$$= \sum_{\substack{w \\ \text{s.t. } w \\ \text{ends in} \\ \text{ion}}} p(w)$$

$\begin{cases} 1 & \text{if } \Pi \text{ is true} \\ 0 & \text{o/w} \end{cases}$

$$= \sum_{w \in \Omega} I(\underbrace{w \text{ ends in } -\text{ion}}_{\Pi}) p(w) =$$

$$= E[I(w \text{ ends in ion})]$$

Model Family

A set of probability distributions (unigram example):

$$\underline{\underline{\mathcal{M}}} = \{ p_1, p_2, \dots \}$$

model (pointing to \mathcal{M})
model (pointing to p_1)

$$p_i: \Omega \rightarrow [0, 1]$$

Parameters

A set of parameters: Θ for $\theta \in \Theta$ $p(w|\theta)$

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

$$\Omega = \{ \text{the, dog, ...} \} \quad |\Omega| = V$$

$p(w)$ = prob. of word w
unigram

find a θ that characterizes what we'd expect
from $p(w)$

$$\Theta = \{ \theta \in \mathbb{R}^{V-1} \text{ s.t. } 1 \geq \theta_i \geq 0 \quad \sum_{i=1}^V \theta_i = 1 \}$$

Another Parametrisation

Rely on properties of the words:

θ will be just a vector of length 26

$\theta_a \theta_b \theta_c, \dots$

$$p(w) = \prod_{i=1}^{|w|} p(w_i)$$

Estimation

What is training data?

$$w^{(1)}, w^{(2)}, w^{(3)} \dots \in \Omega$$

Estimation

What is the fit of the data to the model?

$f(w^{(1)}, \dots, w^{(n)}, \theta)$ - tells what is the fit of θ to $w^{(1)}, \dots, w^{(n)}$

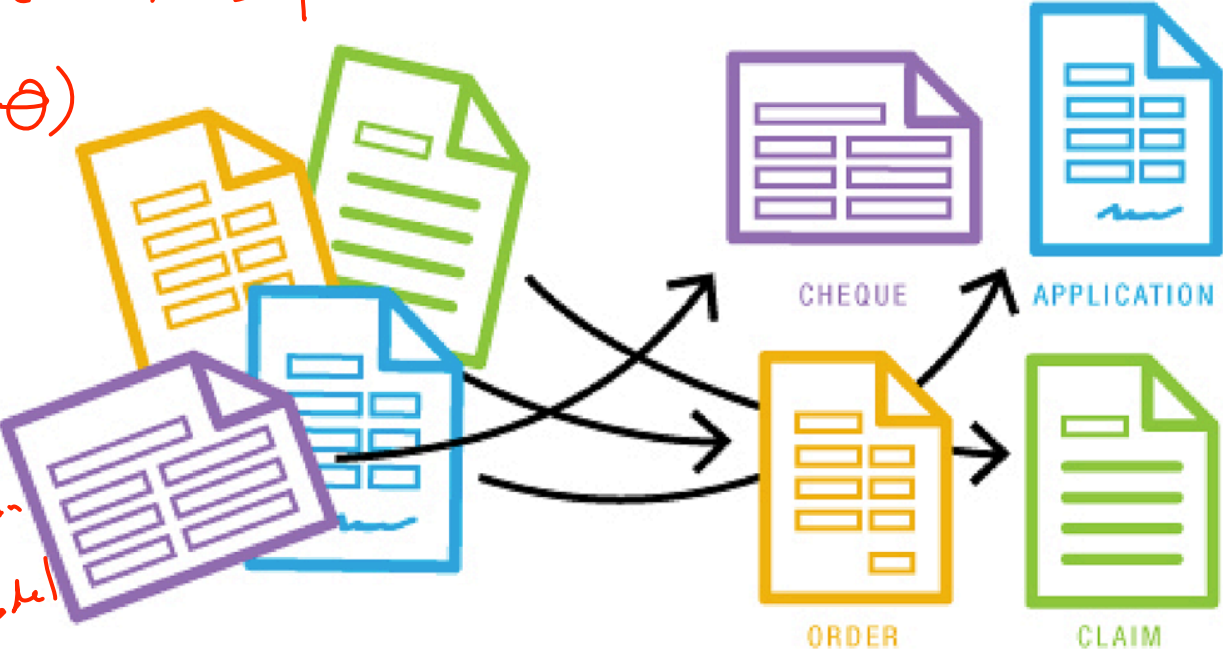
$$\theta^* = \arg \max_{\theta} f(w^{(1)}, \dots, w^{(n)}, \theta)$$

NLP Problem Example: Document Classification

sentiment analysis, document topic, ...

$$\Omega = \{ (d, c) \mid c \text{ is a class, } d \text{ is a document} \}$$

$$p(d, c \mid \theta)$$



Prediction of the model

$$\arg \max_c p(c \mid d, \theta)$$

assignment

$$\theta^* = \arg \max_{\theta} f((d^{(1)}, c^{(1)}), \dots, (d^{(n)}, c^{(n)}), \theta)$$

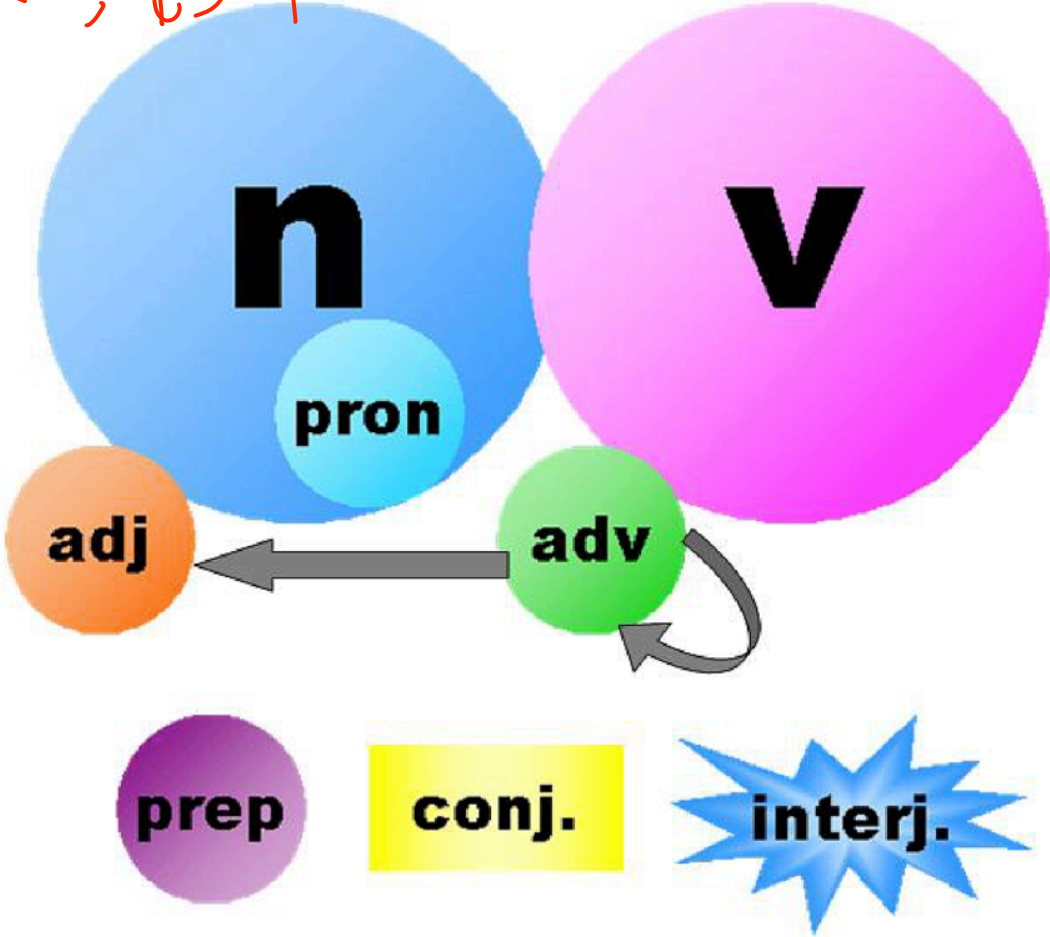
NLP Problem Example: POS Tagging

map words to their part-of-speech tags

$$\Omega = \{ (s, t) \mid s \text{ is a sentence, } t \text{ is a tag} \}$$

$$p(s, t | \theta)$$

$$\underset{t \text{ max}}{p(t | s, \theta)}$$

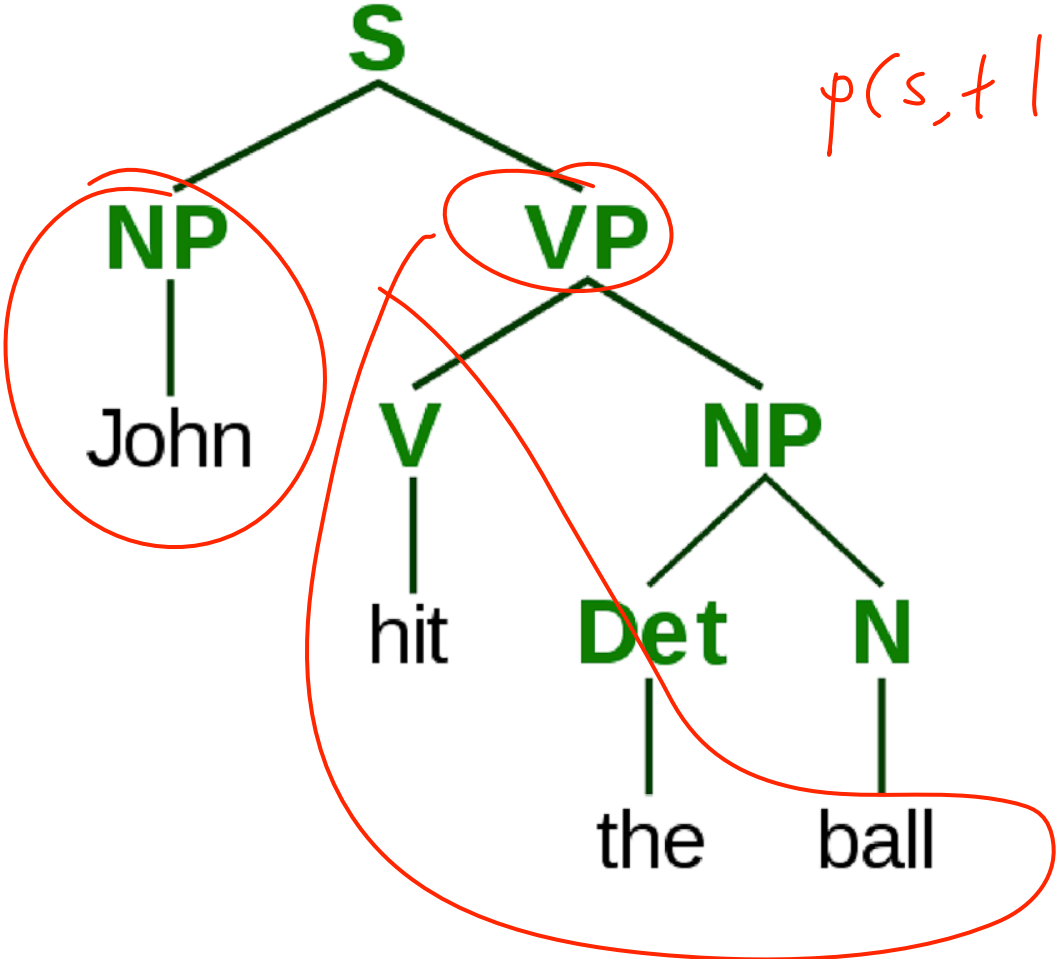


NLP Problem Example: Parsing

map sentences to their syntax

$$\Omega = \{ (s, t) \mid s \text{ is sentence, } t \text{ is a tree} \}$$

$$p(s, t \mid \theta)$$



NLP Problem Example: FrameNet Parsing

find predicate-argument structure

James **has** a **university degree** in astronomy .

 POSSESSION LOCALE_BY_USE QUANTITY

Locale **Quantity**

Owner **Possession**

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?”