

# 4062

## 'Advanced Computational Linguistics'

### Machine Learning Techniques in Machine Translation, Speech Recognition and Topic Modelling

[www.scss.tcd.ie/Martin.Emms/4062](http://www.scss.tcd.ie/Martin.Emms/4062)

Martin Emms

March 31, 2020

# Unsupervised Learning from Data

Machine Translation

# Unsupervised Learning from Data

Machine Translation

Speech Recognition

# Unsupervised Learning from Data

Machine Translation

Speech Recognition

Topic Modelling



# Unsupervised Learning from Data

Machine Translation

Speech Recognition

Topic Modelling

} **Use Unsupervised Machine Learning**

# Unsupervised Learning from Data

Machine Translation

Speech Recognition

Topic Modelling

} **Use Unsupervised Machine Learning**

Possibly other attempts to infer from Big Data } **will use Unsupervised Machine Learning?**

# In Machine Translation

Unsupervised methods are used to go from **sentence** pairs to **word** pairs

# In Machine Translation

Unsupervised methods are used to go from **sentence** pairs to **word** pairs

la maison est grand	the house is big
c'est un haricot vert	its a green bean
je le lui donne	I give it to him
⋮	⋮
⋮	⋮
⋮	⋮

⇒

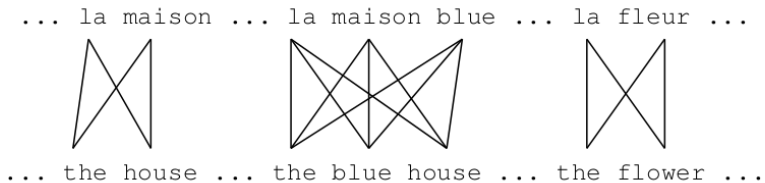
$p(la the)$	0.453
$p(le the)$	0.334
$p(maison house)$	0.876
$p(bleu blue)$	0.563



# Done using EM Algorithm

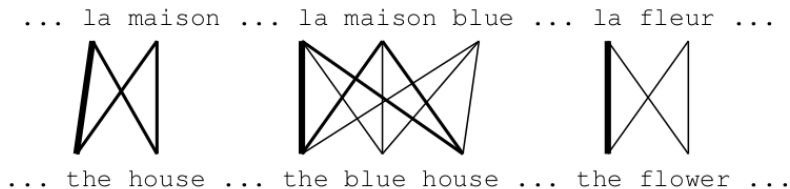
- ▶ Incomplete data
  - ▶ if we had *complete data*, would could estimate *model*
  - ▶ if we had *model*, we could fill in the *gaps in the data*
- ▶ Expectation Maximization (EM) in a nutshell
  1. initialize model parameters (e.g. uniform)
  2. assign probabilities to the missing data
  3. estimate model parameters from completed data
  4. iterate steps 2–3 until convergence

# Learning Word Translations



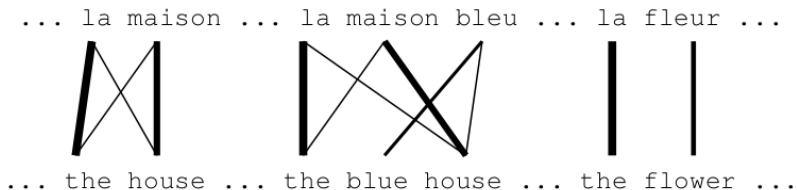
- ▶ Initial step: all alignments equally likely

# Learning Word Translations



- ▶ After one iteration
- ▶ Alignments, e.g., between **la** and **the** are more likely

# Learning Word Translations



- ▶ After another iteration
- ▶ It becomes apparent that alignments, e.g., between **fleur** and **flower** are more likely (pigeon hole principle)

# Learning Word Translations

... la maison ... la maison bleu ... la fleur ...  
/ | | X | |  
... the house ... the blue house ... the flower ...

- ▶ eventually converges, inherent hidden structure found by EM
- ▶ **Learns** translation probabilities

$p(la the)$	0.453
$p(le the)$	0.334
$p(maison house)$	0.876
$p(bleu blue)$	0.563

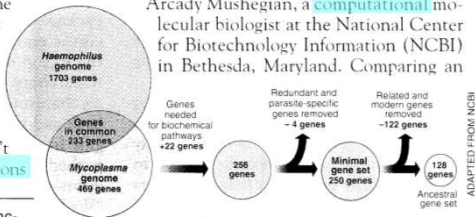
# Topics in Documents

## Seeking Life's Bare (Genetic) Necessities

COLD SPRING HARBOR, NEW YORK—How many genes does an organism need to survive? Last week at the genome meeting here,\* two genome researchers with radically different approaches presented complementary views of the basic genes needed for life. One research team, using computer analyses to compare known genomes, concluded that today's organisms can be sustained with just 250 genes, and that the earliest life forms required a mere 128 genes. The other researcher mapped genes in a simple parasite and estimated that for this organism, 800 genes are plenty to do the job—but that anything short of 100 wouldn't be enough.

Although the numbers don't match precisely, those predictions

“are not all that far apart,” especially in comparison to the 75,000 genes in the human genome, notes Siv Andersson of Uppsala University in Sweden, who arrived at the 800 number. But coming up with a consensus answer may be more than just a genetic numbers game, particularly as more and more genomes are completely mapped and sequenced. “It may be a way of organizing any newly sequenced genome,” explains Arcady Mushegian, a computational molecular biologist at the National Center for Biotechnology Information (NCBI) in Bethesda, Maryland. Comparing an



\* Genome Mapping and Sequencing, Cold Spring Harbor, New York, May 8 to 12.

**Stripping down.** Computer analysis yields an estimate of the minimum modern and ancient genomes.

Documents exhibit multiple topics **biology** **genetics** **computation** ...

# Topics in Documents

## Topics

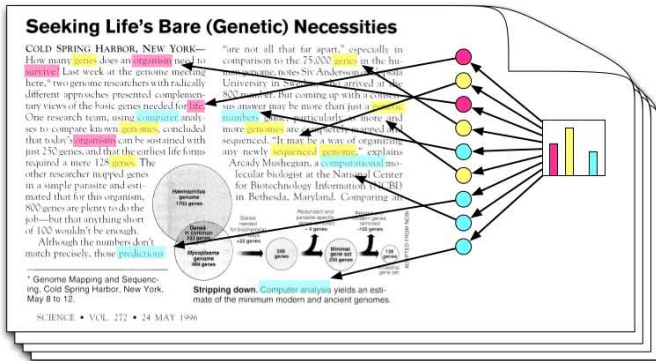
gene 0.04  
dna 0.02  
genetic 0.01  
...

life 0.02  
evolve 0.01  
organism 0.01  
...

brain 0.04  
neuron 0.02  
nerve 0.01  
...

data 0.02  
number 0.02  
computer 0.01  
...

## Documents

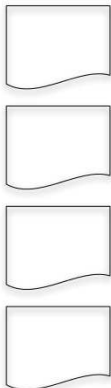


## Topic proportions and assignments

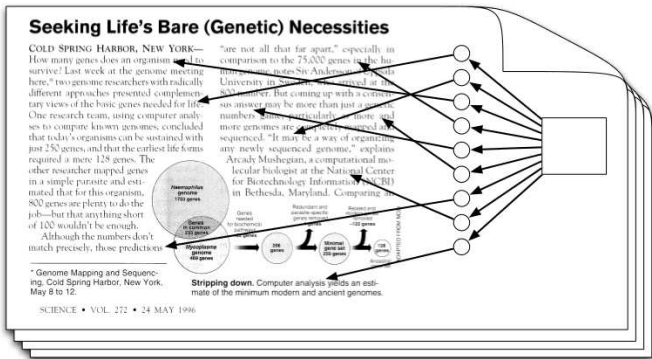
- ▶ **topic** = distribution over words
- ▶ **document** = distribution over topics

# Topics in Documents

Topics



Documents



Topic proportions and assignments

- ▶ can only directly see words
- ▶ but can learn the rest by **unsupervised methods**



# Example learned words-for-topics

## GENETICS

human  
genome  
dna  
genetic  
genes  
sequence  
gene  
molecular  
sequencing  
map  
information  
genetics  
mapping  
project  
sequences

## EVOL. BIOL.

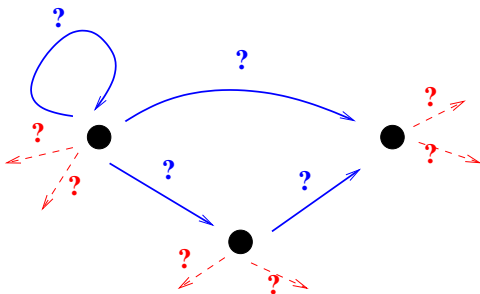
evolution  
evolutionary  
species  
organisms  
life  
origin  
biology  
groups  
phylogenetic  
living  
diversity  
group  
new  
two  
common

## COMPUTING

computer  
models  
information  
data  
computers  
system  
network  
systems  
model  
parallel  
methods  
networks  
software  
new  
simulations

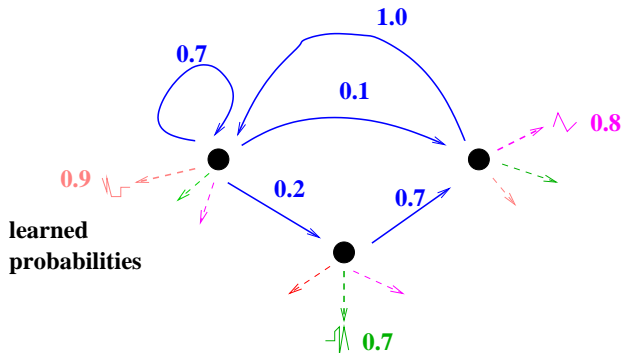
# Hidden Markov Models for Speech or Actions

- ▶ words/action visible as an *evidence* sequence
- ▶ can be modeled with a HMM  
with **transitions between states** and **visible evidence for a state**



# Learning HMM from Data

**data:** 



- ▶ Possible to use **unsupervised learning** to find probabilities concerning *hidden* variables from data with just *visible* evidence
- ▶ Used in Speech Recognizers, Activity Recognizers

## despite the name . . .

- ▶ though we look at techniques which can be applied to language, these techniques/ideas are applicable much more widely than that
- ▶ the material will be just as accessible to ICS students as to CSL students: the number of non-CSL students taking the course last year exceeded the number of CSL students

- ▶ so, how does it differ from other 'ML' modules?

- ▶ so, how does it differ from other 'ML' modules?
  1. pervasive **hidden variables**
  2. data unbounded size, not fixed-size grid of features/values
  
- ▶ still, *why* do this?

- ▶ so, how does it differ from other 'ML' modules?
  1. pervasive **hidden variables**
  2. data unbounded size, not fixed-size grid of features/values
  
- ▶ still, *why* do this?
  1. the existence of algorithms for **unsupervised** is properly surprising
  2. ingenuity in getting from exponential to feasible cost is properly impressive

## Coursework/Labs

- ▶ some 'pencil-and-paper' working things out course work



## Coursework/Labs

- ▶ some 'pencil-and-paper' working things out course work
- ▶ some coding from scratch course work

## Coursework/Labs

- ▶ some 'pencil-and-paper' working things out course work
- ▶ some coding from scratch course work
- ▶ some labs and coursework involving particular toolkits/libraries

## Further questions?

please feel free to get in touch by email to seek further advice