# Unsupervised Learning on an Approximate Corpus

Jason Smith, Jason Eisner

## Learning from n-grams

Sentence	es:				C	Counts
	time	flies	like	an	arrow	20
	fruit	flies	like	an	orange	3
your	plane	flies	like	an	ostrich	2
n-grams						
11 S' anno 1	time	flies	like			20
	fruit	flies	like			3
		flies	like	an		25
			like	an	arrow	20
your	plane	flies				2

...

## Contributions

- Learning from a finite-state distribution over sentences
  - e.g. an n-gram language model over sentences, instead of individual sentences
- Why?
  - Original corpus unavailable
  - Speed (learning from compressed data)
  - (Fundamental question about weighted grammars)
- Exact and approximate solutions

## Task

- HMM POS tagging (Merialdo 94)
- Many approaches build off of EM

## **Previous Work**



• (slide taken from Lin et al., 2009)

## **Previous Work**



• (slide taken from Lin et al., 2009)

## **Previous Work**



(slide taken from Lin et al., 2009)

#### Full context:

Ν	V	Adv	Det	Ν
time	flies	like	an	arrow
Adj	Ν	V	Det	Ν
fruit	flies	like	an	orange

Full context:

	Ν	V	Adv	Det	Ν
	time	flies	like	an	arrow
	Adj	Ν	V	Det	N
	fruit	flies	like	an	orange
ocal	n-gram	context:			
	?	V? N?	Adv? V?	Det	?
	?	flies	like	an	?

Local n-gram context:

?	V? N?	Adv? V?	Det	?
?	flies	like	an	?

#### Overlapping n-grams:

time	flies	like	20
fruit	flies	like	3
plane	flies	like	2

Counts

Local n-gram context:

?	V? N?	Adv? V?	Det	?
time 80% fruit 12% plane 8%	flies	like	an	?

#### Overlapping n-grams:

time	flies	like	20
fruit	flies	like	3
plane	flies	like	2

Counts

Local n-gram context:

N A	di	V N	Adv v	Det	?
time fruit plane	80% 12% e 8%	flies	like	an	?
verlap	ping n-g	grams:	C	Counts	
tin	ne	flies	like	20	
fru	uit	flies	like	3	
pla	ne	flies	like	2	

## Exploit Overlapping n-grams

Counts

	time	flies	like			20
		flies	like	an		25
			like	an	arrow	20
	fruit	flies	like			3
			like	an	orange	3
your	plane	flies				2
	plane	flies	like			2
			like	an	ostrich	2

# Exploit Overlapping n-grams n-gram language model!







## N-gram language models







## Task

- HMM POS tagging (Merialdo 94)
- Many approaches build off of EM

$$\max \sum_{w \in \text{corpus}} \frac{1}{n} \log \sum_{t} p(t, w)$$
$$\max \sum_{w \in \text{corpus}} c(w) \log \sum_{t} p(t, w)$$

Sentence: time flies like an arrow

#### Sentence: time flies like an arrow

## p(Tag|Last Tag)

	Det	Z	V	•••
Det	0.1	0.1	0.5	
Ν	0.8	0.3	0.4	•••
V	0.1	0.6	0.1	

...

...

#### Sentence: time flies like an arrow

23

## p(Tag|Last Tag)

	Det	Ζ	V	•••
Det	0.1	0.1	0.5	
Ν	0.8	0.3	0.4	•••
V	0.1	0.6	0.1	
				-



	Det	Ζ	V	•••
time	0.01	0.3	0.1	
flies	0.01	0.2	0.2	•••
an	0.33	0.01	0.01	
		•••		

#### Sentence: time flies like an arrow



 $n_c$  : c(w)'s word context window  $n_p$  : p(t,w)'s tag context window

## Supervised learning: HMM

NVAdvDetNtimeflieslikeanarrow

## Supervised learning: HMM

transition counts: estimating p(Tag|Last Tag)



## Supervised learning: HMM

emission counts: estimating p(Word|Tag)





### What if someone tagged our n-grams?



### What if someone tagged our n-grams?

c(w) c(t,w):



### What if someone tagged our n-grams?



 $n_c : c(w)$ 's word context window  $n_P : p(t,w)$ 's tag context window  $n_q : c(t,w)$ 's tag context window













## Unsupervised learning

Sentence: time flies like an arrow

## p(Tag|Last Tag)

	Det	Ζ	V	•••
Det	0.1	0.1	0.5	
Ν	0.8	0.3	0.4	•••
V	0.1	0.6	0.1	
				-



	Det	Ζ	V	•••
time	0.01	0.3	0.1	
flies	0.01	0.2	0.2	•••
an	0.33	0.01	0.01	
		•••		

38

#### **HMM** Tagging Trellis $p(t,w) \circ c(w):$ Adv V flies like Det V V Adv Ν V Ν an time like flies Ν arrow Ν Det <S> Adj V Adv Det time flies like an Adj V Ν Ν V flies like flies like time an arrow





### Let's tag our own n-grams (EM) $n_q = 2$ c(w) c(t,w) c(w)q(t|w):



# Let's tag our own n-grams (EM) $n_q = 2$ $\frac{n_q = 2}{r_q}$



 $n_c : c(w)$ 's word context window  $n_P : p(t,w)$ 's tag context window  $n_q : q(t|w)$ 's tag context window

$$\log \sum_{t} p_{\theta}(t, w) = \log \sum_{t} q(t|w) \left(\frac{p_{\theta}(t, w)}{q(t|w)}\right)$$
$$\geq \sum_{t} q(t|w) \log\left(\frac{p_{\theta}(t, w)}{q(t|w)}\right)$$
$$= \mathbb{E}_{q(t|w)} [\log p_{\theta}(t, w) - \log q(t|w)]$$

$$\begin{split} \log \sum_{t} p_{\theta}(t, w) &= \log \sum_{t} q(t|w) (\frac{p_{\theta}(t, w)}{q(t|w)}) \\ \text{ensen's inequality} \end{split} \geq \sum_{t} q(t|w) \log(\frac{p_{\theta}(t, w)}{q(t|w)}) \\ &= \mathbb{E}_{q(t|w)} [\log p_{\theta}(t, w) - \log q(t|w)] \end{split}$$

$$\mathbb{E}_{c}(w) \log \sum_{t} p_{\theta}(t, w) = \mathbb{E}_{c}(w) \log \sum_{t} q(t|w) \left(\frac{p_{\theta}(t, w)}{q(t|w)}\right)$$
$$\geq \frac{E_{c}(w)}{t} \sum_{t} q(t|w) \log\left(\frac{p_{\theta}(t, w)}{q(t|w)}\right)$$
$$= \mathbb{E}_{c}(w) q(t|w) \left[\log p_{\theta}(t, w) - \log q(t|w)\right]$$



## How to maximize this bound

 $\mathbb{E}_{c(w)q(t|w)}[\log p_{\theta}(t,w) - \log q(t|w)]$ 

- Updating p(t,w) (M-step): shown earlier
- Updating q(t|w) (E-step): more complex, but has a dynamic programming solution which makes use of finite-state machines
  - Expectation semirings (Eisner 2002), details in paper

## Experiments: EM vs. n-gram EM

- How does EM on a full corpus compare to ngram EM on an approximate corpus?
  - POS tagging accuracy and likelihood
- Standard setup for unsupervised POS tagging with a dictionary
- Reduced tag set (17 tags)
- Limited tag dictionary from WSJ (words must appear 5 times, otherwise all tags are possible)

## Experiments: EM vs. n-gram EM

- n-gram EM parameter choices:
  - $n_c=5 c(w)$  uses up to 5-grams
  - $n_p=2 p(t,w)$  is a bigram HMM
  - n<sub>q</sub>=I q(t|w) conditions tag only on ngram word context (approximate, but saves space)

## Results:WSJ







## Conclusions

- **New problem**: train on an infinite corpus (distribution over sentences)
- **New algorithms**: exact and approximate likelihood maximization
- New results: faster (sublinear) training by compressing corpus into n-gram model