

schedule checking appts

- have diff comparison:

comparing ~~comparing~~ language models. [A] - what can you tell from them?
[B] how can you compare them?

[A] language models

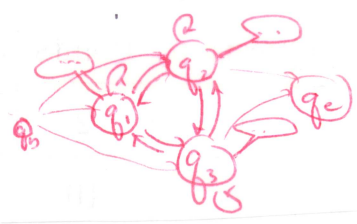
→ which we spent some time on last time:
remember we were broadly constraining these as:

(^{prob.} distributions over strings $\langle \text{begin} \rangle V^* \langle \text{end} \rangle$ for fixed non-empty vocab V)

• ~~HMM-based lang. models~~
hidden-markov-model based language models - "ngram" models as a special case.
probabilistic context-free-grammar based language models

typically only implicitly noted.
Put here as a reminder to key the generation model in mind.

"structured" language model like a hidden Markov Model:



can be too big to look @.

[could, of course, just start w/ fewer states; but how satisfying is that, since you built the one w/ a larger state set.]

should:
if " $P(x)$ "
 $= P(x)P(x|x)$
 $= P(x)P(x)$
then $P(x) + P(x) = 2$ for this x
so don't confuse "base distribution" w/ distribution over all strings.

(1) Fix length l , compute top- k highest prob paths of length l . [product of its transition & emission probs]
"what does this HMM 'like to do'?"

- dynamic programming.

[of course, you might want to vary l]

(2) Fix sequence of interest s , compute top- k paths that assign highest probs to s .

reminder: can have more than one path generate the same sequence.

<why this instead of (1), sometimes?>

- ~~might not care ab~~ your LM might spend a lot of prob on sequences you don't care about, i.e., nonsense
- you might be particular interested in possible analyses of a test seq.

<ex: Polymath discussion>

[of course, you might want to vary s]

(3) look @ ~~distribution~~ transition or emission ~~prob~~ distribution for a particular state.

~~note: this generalises to n-gram mod~~

→ $p(x)$ over ~~some~~ some space X .

For simplicity, assume X ~~finite~~ finite, and write θ_i for the prob of the i th elt.

Notation:

(b/c we will want these p 's to be variables, soon).

constraints: $\theta_i \geq 0, \sum \theta_i = 1$.
let θ be the vector of θ_i 's.

- note: n-gram models fit this: $p(x|y)$: prob given word x that next word is y .
 θ over some finite vocabs - each x induces its own distribn.

ex: $X = \{\text{the}, \text{dog}\}$

$$P(\text{dog} | \text{the}) > P(\text{the} | \text{the})$$

$$\theta_{\text{dog}} > \theta_{\text{the}}$$

- what can we say about the θ_i 's, besides just eyeballing them.

- measures of "diversity" or "spread"

- are lots of things equally likely, or are there only a few things that are highly likely?

~~are important ones~~: the entropy, Gini-Simpson index (others possible, like L_1 or L_2 norms)

entropy $H(\theta) = \sum_i \theta_i \log\left(\frac{1}{\theta_i}\right)$

↳ ~~big~~ small when θ_i big.
"surprise" in seeing event i .

"back-integral" as $p(x_i)$

= "expected surprise".

ex: ~~$\theta_i = 0$~~ $\theta_i = \begin{cases} 1, & i=1 \\ 0 & \text{o.w.} \end{cases}$

$$H(\theta^{\text{shock}}) = 0 + \sum_{i \neq 1} 0 = 0$$

you are never surprised, b/c you will always see x_1 .

~~$H(\theta) = -\sum_i \theta_i \log(\theta_i) = -\log(1) = 0$~~

So we know what has 0 entropy.
But could this be negative?

And what ~~could be the~~ kind of distribution has the highest entropy?

Find ~~the~~ ~~max~~ ~~of~~ ~~the~~ ~~entropy~~ ~~function~~ $H(\vec{\theta}) - \lambda (\sum_i \theta_i - 1)$
constant opt.

(doing this exercise even if you already know the answer b/c useful later).

$$\frac{\partial}{\partial \theta_j} \left(\sum_i \theta_i \log \left(\frac{1}{\theta_i} \right) - \lambda \left(\sum_i \theta_i - 1 \right) \right)$$

when $\theta_i \neq \theta_j$, the whole thing is a constant wrt θ_j .

$$= \frac{\partial}{\partial \theta_j} \left[\theta_j \log \left(\frac{1}{\theta_j} \right) - \lambda \theta_j \right]$$

~~θ_j~~

$$= -\theta_j \cdot \frac{1}{\theta_j} - \log \theta_j - \lambda \quad \text{set to } 0$$

$$-1 - \log \theta_j - \lambda = 0$$

~~$\log \theta_j =$~~

$\log \theta_j = -(1-\lambda)$ constant. So all the θ_j 's are the same.

- max "surprise" (least concentration) when you have "all possibilities equally likely".

Same

~~What about Gini S~~

Gini-Simpson index $GS(\vec{\theta}) = 1 - \sum_i \theta_i^2$: prob that 2 samples will not be the same.

$= 1 - \sum_i \theta_i \theta_i$
 " \rightarrow θ_i when θ_i is θ_i
 " \rightarrow expectation of θ_i
 ~~$P(x_i)$~~ (cf. entropy).

\bullet $GS(\vec{\theta}^{sharp}) = 1 - 1 - \sum_{i=1} 0 = 0 \checkmark$

max? $\frac{\partial}{\partial \theta_j} \left(1 - \sum_i \theta_i^2 - \lambda (\sum_i \theta_i - 1) \right)$

$$= \frac{\partial}{\partial \theta_j} \left(\theta_j^2 - \lambda \theta_j \right) = 2\theta_j - \lambda$$

set to 0, again, all θ_j 's the same.

[B] comparing two LMs.

could use L_2, L_1 , etc.

Today, KL divergence: $D(\vec{\theta} \parallel \vec{\varphi}) = \sum \theta_i \log \frac{\theta_i}{\varphi_i}$ \rightarrow some terms can be negative.

$$= -\sum \theta_i \log \varphi_i - H(\vec{\theta}_i)$$

if $\vec{\varphi} = \vec{\theta}$, $D(\vec{\theta} \parallel \vec{\theta}) = -\sum \theta_i \log \theta_i - H(\theta_i) = H(\vec{\theta}_i) - H(\vec{\theta}) = 0$.

$$\frac{\partial}{\partial \varphi_j} \left(-\sum_i \theta_i \log \varphi_i - H(\vec{\theta}) - \lambda (\sum_i \varphi_i - 1) \right)$$

$$= -\theta_j \cdot \frac{1}{\varphi_j} - \lambda \quad \text{set to 0,} \quad -\frac{\theta_j}{\varphi_j} = \lambda \quad \text{or} \quad -\frac{\theta_j}{\lambda} = \varphi_j$$

Since θ_i 's already sum to 1,

$$\lambda = -1$$

so, $\vec{\theta} = \vec{\varphi}$ is at least a critical point.

2nd partials: $\frac{\partial^2}{\partial \varphi_i \partial \varphi_i} = 1/\varphi_i^2 > 0$, at least one is > 0 . So, Hessian is positive definite, and we have a minimum.

And note that if $\exists j$ s.t. $\theta_j > 0, \varphi_j = 0$, $D(\vec{\theta} \parallel \vec{\varphi}) = \infty \dots$

is that weird. Two 'finite' things having 'infinite' distance?

but it makes sense: 2 distributions when one says something is possible that another says impossible, they are irreconcilable!

~~al~~ ... altho' zeroes in ~~estimate~~ ^{data} \rightarrow you think those things are necessarily impossible...

- solution: (a) smooth your distributions < if had time, would have done Kneser-Ney >

(b) use Jensen-Shannon divergence: $\pi_i \neq 0$ if θ_i or $\varphi_i \neq 0$.

$$\vec{\pi} = \frac{1}{2} (\vec{\theta} + \vec{\varphi}) \quad \text{-- never } \vec{0}$$

$$JSD(\vec{\theta}, \vec{\varphi}) = \frac{1}{2} (D(\vec{\theta} \parallel \vec{\pi}) + D(\vec{\varphi} \parallel \vec{\pi}))$$

- or the skew divergence:

$$D_\alpha(\vec{\theta}, \vec{\varphi}) = D(\vec{\theta} \parallel \alpha \vec{\varphi} + (1-\alpha) \vec{\theta})$$

cross-entropy: $-\sum \theta_i \log \varphi_i$ (expected'n vrt $\vec{\theta}$ of surprise according to $\vec{\varphi}$)

for empirical unigram model:

$$-\sum_i \frac{\#(x_i)}{\text{data size}} \log \varphi_i = -\frac{1}{\text{data size}} \log \prod \varphi_i^{\#(x_i)}$$

which looks like the 'prob' assigned to the sample.

Ex: in no country for old members, empirical = empirical lang. sample, $\vec{\varphi}$ = user's l.u.

In general:

for large sample, do $\frac{1}{\text{sample size}} P_{\vec{\theta}}(\text{sample})$

then are ~~using~~ better arguments using Jensen's inequality or other log properties, but I wanted to stick w/ using the same techniques throughout the lecture.