

Recall our question: how do you tell if two ~~things~~ "things" are different?  
(statistical)

What is ~~language~~ are language models for? For the purposes of this class:  
 → giving a ~~short~~ compact representation of what the language is 'like'.  
 → comparing two language sources. (ex: no country)

We define a language model as a distribution over all strings  $w$  that we wish to consider  
 "ok": "legal".

~~Some distributions aren't model-based per se, e.g. HMMs~~  
 ex: HMMs (for ~~now~~, no a lot's)

[a good counterexample to key in mind: the Poisson over strings of form "dude"  
 - finite # of params, no equivalent PCFs. (see Booth & Thompson '73)]

• a finite non-empty set of states  $q_1, \dots, q_M, q_b, q_e$   
 distinguished "start" begin distinguished "end"

- a finite vocabulary  $V \cup \{\langle \text{begin} \rangle, \langle \text{end} \rangle\}$ , not in  $V$ .

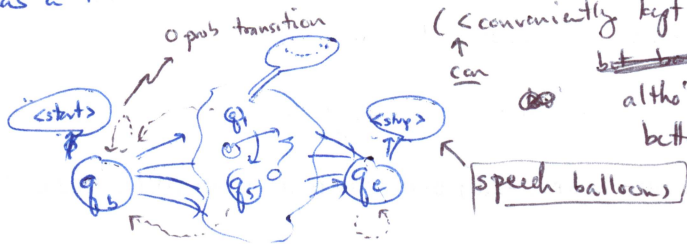
- each state ~~def~~ has an associated emission distribution over  $V^*$  (all sequences of items from  $V$ , including the empty string)

$q_b$  emits only  $\langle \text{begin} \rangle$   
 $q_e$  .. ..  $\langle \text{end} \rangle$   
 $q_i$  does not emit any sequence containing those special items

this might be itself a (state-specific) L.M.

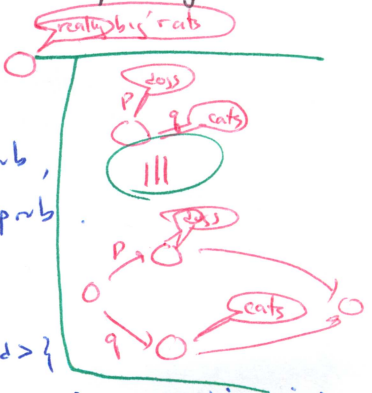
- each state has a transition distribution over all other states, s.t:  
 (conveniently kept in a matrix, all rows sum to 1)  
 can be ~~bad~~ but ~~bad~~ if altho' maybe a sparse rg would be better.

draw 1st  
 - drawing really uncalculated a lot



will further assume every  $q_i$  reachable from  $q_b$  w/ non-zero prob,  
 $q_i$  can reach  $q_e$  w/ non-zero prob  
 → these definitely seem like necessary

~~$p(s) = \dots$~~   
 $p(s) = \text{sum of probs of all paths generating } s, s \in \{\langle \text{begin} \rangle\} V^* \{\langle \text{end} \rangle\}$   
 = product of all the transition probs x prob of corresponding emissions.



nb: should make clear diff. b/w termination and progress (i.e. ~~the~~ and the l.m. itself being proper.

(so you can think of the emissions as also ~~def~~ part of what defines a path.

What does this tell you?

- broken <sup>we have</sup> & distribution over a potentially  $\infty$  # of sequences to sth definable
- by a finite set of parameters (indeed, you have a model  $\theta$ , so you're managing the complexity.

- structure of this model gives you insight into what is "preferred".

- learning algorithms for HMMs, either from labelled or unlabelled data, (altho' you have to set the # of states beforehand)

learnable two training-data paradigms: state-labelled data, unlabelled data

need to specify  $M$  ahead of time. Baum-welch / EM for max-likelihood estimation (tends to find saddle points or local maxima, b/c there's a lot)

- REMEMBER: need constraints (ow., just set all params to  $\infty$ !) - why are n-gram models so much more common than HMMs?

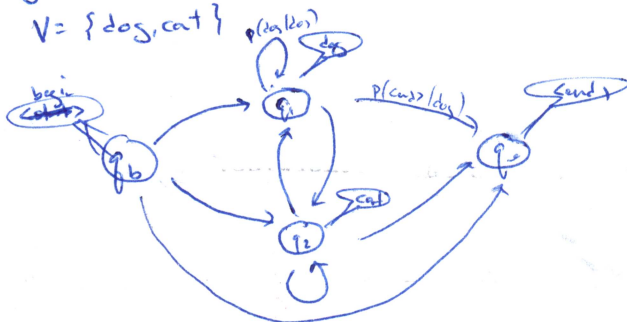
- here, the data basically has state labels! Nothing is really "hidden".

$$P(\langle \text{cat} \rangle \rightarrow \text{dog, dog} \langle \text{stop} \rangle)$$

$$P(\rightarrow) \cdot 1 \cdot P(\rightarrow) \dots$$

ex: a "bigram lm" as an HMM:

$V = \{\text{dog, cat}\}$



$q_i = \text{"just said } v_i \text{"}$

vs. unigram model:



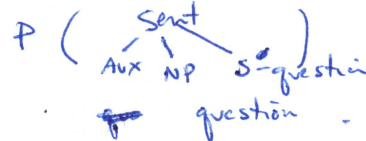
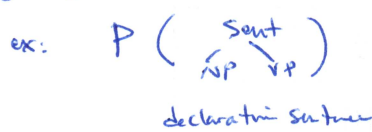
~~linear dependence of what is to be next generated~~

~~we can build more co~~

~~but it's not just finite state~~ - we aren't just restricted to finite state models

ex: probabilistic context-free grammars (PCFGs) - analogous defns for dependency grammars. aka stochastic CFGs (SCFGs)

- each category type has a distribution over its decompositions, a finite set:



~~what~~

$P(s) = \sum_{\theta} \text{sum of probs of each trees generating } s$

↳ product of each decomposition used in a tree

$$P \left( \begin{array}{c} \text{sent} \\ \swarrow \quad \downarrow \quad \searrow \\ \text{NP} \quad \text{VP} \\ \text{they} \quad \text{left} \end{array} \right) = P \left( \begin{array}{c} \text{sent} \\ \swarrow \quad \searrow \\ \text{NP} \quad \text{VP} \\ \text{they} \quad \text{left} \end{array} \right) P(\text{NP} \mid \text{they}) P(\text{VP} \mid \text{left})$$

q: ~~does this~~ does every choice of valid 'decomposition' distribution 'induce a proper prob on all strings? Turns out it's tricky.

- [Booth; Thompson '73] conditions under which you would get a proper distribution over all strings.

[Chi & Gema '98]



For this  $\theta = (p)$  <that's all we need>

$$P_0(\text{"police"}) = q$$

$$P_0(\text{"police police"}) = p \cdot q \cdot q = p q^2$$

$$P_0(\text{"police police police"}) = [p \cdot p \cdot q^3] + [p \cdot q \cdot p \cdot q^2] = 2 p^2 q^3$$

$$\sum_s P_0(s) = \sum_{i=1}^{\infty} \# \text{ of 'bracketings' of length-} i \text{ sequence} \cdot q^i \cdot p^{i-1}$$

generate  $i$  leaves  
got  $i$  "parents" of those leaves.

Let  $X_d$  = prob of all trees w/ depth ~~height~~  $\leq d$  : fringe contains only words

$$X_{d=1} = q$$

$$X_2 = q + p q^2$$

$$X_i = q + p X_{i-1}^2$$

$$X_3 = q + p q^2 + p (q + p q^2)^2 (q + p q^2)^2$$

apparently converges to  $\min(1, \frac{q}{p})$

$$= \sum_{i=1}^{\infty} (\# \text{ bracketings}(i)) \cdot q^i (q p)^{i-1}$$

$$\sum_{i=1}^{\infty} \frac{1}{i} \binom{2(i-1)}{i-1} q (q p)^{i-1}$$

<presumably the same geometric sum?>

$$\approx q \sum_{i=1}^{\infty} \frac{1}{i} \frac{4^{i-1}}{\sqrt{\pi(i-1)}} (q p)^{i-1}$$

abandon this

if  $p > 1/2$ , this will not be proper (too much probability on the "expansion"  $p$  as opposed to the "generation of words").

ex: content model: [Burgin; Lee '04]

- code available (also Alexander Passos)

→ (learn) how topics relate to each other w/in docs.

- etc. state

→ potentially relevant to projects:

<probably have fixed set of topics>

(etc. state)

talked about LM on just top-most frequent words

stylistic

what do you do w/ the 'other' words to make them not dominate?

entropy → fixed-max-length strings.

2nd attempt, using Chi & Gema's approach:

let  $X_d$  = prob of all trees w/ depth  $\leq d$  ( : fringe contains only words)



$$\rightarrow p \cdot X_d \cdot X_d$$

and then  $\frac{S}{\text{police}} = q$

$$X_{d+1} = q + p X_d^2$$

guess that  $X_{d+1} \geq X_d \geq 0$  ~~so  $X_{d+1} \leq q + p X_{d+1}^2$~~

If we assume convergence, then we solve for  $X = q + p X^2$ , or  $0 = p X^2 - X + q$ .

By quadratic formula, roots are:

$$\Rightarrow 0 = (p x - (1-p))(x-1)$$

$$= (p x - q)(x-1)$$

solutions are  $x = q/p$  or  $x = 1$ .

(thanks to Chanhoo Tan for idea of assuming convergence)