## (1) intro

• some terminology:

```
                    S
         ┌──────────┴──────┐
        NP                 VP
         │          ┌──────┴────┐
        DT         VBZ         NP
                          ┌─────┼─────┐
                         DT    JJ     NN
```

} syntax

→ part of speech

this    is    a    simple    sentence    → words
        be          SIMPLE1   SENTENCE1   → morphology
        3sg         having    string of words  } semantics
        present     few parts satisfying the
                              grammatical rules
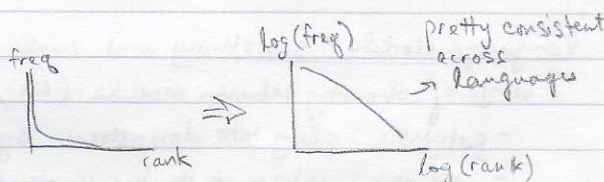                              of a language

→ discourse

## why is NLP hard?

• NLP is hard b/c of ambiguity    But it is an instructive one
  - word senses: bank (finance or river?)
  - part of speech: chair (noun or verb?)
  - syntactic structure: I saw a man with a telescope
  - quantifier scope: Every child loves some movie
  - multiple: I saw her duck.

• data is sparse due to Zipf's law: freq vs. rank

$$f \cdot r \approx k$$ where $f$ = freq of word, $r$ = rank (by $f$), $k$ = constant

$$\Rightarrow \log f = \log k - \log r \qquad y = b + mx$$



freq vs rank; log(freq) vs log(rank) — pretty consistent across languages

• variation: think WSJ vs. social media
• expressivity: can express same thing many different ways
• context dependence and unknown representation: need world knowledge (sometimes/usually), don't know how to represent "meaning"

## (2) text corpora

• good if naturally-occuring, includes metadata, or has linguistic annotations
• for sentiment analysis, good to have some sort of ratings/labels
• problems w/ counting positive/negative words: sense ambiguity, sarcasm/irony, (other) context
• tokenization: adding logical boundaries between separate word/punctuation tokens not already separated by spaces
  - only one part of preprocessing/normalization
• "domain" of corpus includes mode of communication (e.g., speech, writing), topic (e.g., politics, physics), genre (e.g., news, tweet), audience (formality, politeness, complexity)
• domain adaption corrects for when training data and test data don't come from the same source

## (3) n-gram language models

• this lecture talks about sentence probabilities
• language model gives approx. for true prob. of an arbitrary sequence of words
• uses of a language model: spelling correction, automatic speech recognition, machine translation
  - in each case, get possibilities from some model (e.g., error model, acoustic model, translation model) and use language model to choose among alternatives
• unrealistic to do MLE on complete sentences; instead, multiply probabilities of sub-parts
• n-gram model:

$$P(w_1, \ldots, w_n) = \prod_{i=1}^{n} P(w_i \mid w_1, w_2, \ldots, w_{i-1}) \approx P(w_i \mid w_{i-2}, w_{i-1}) \text{ for trigram model}$$

$$\approx P(w_i \mid w_{i-1}) \text{ for bigram}$$

$$\approx P(w_i) \text{ for unigram}$$

• estimation w/ trigram model:

$$P_{MLE}(w_i \mid w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i)}{C(w_{i-2}, w_{i-1})}$$

1

n-gram language models (cont'd)
- practical implementation of tri-grams:
$$P(\vec{w}) = P(w_1) P(w_2|w_1) \prod_{i=3}^{n} P(w_i|w_{i-2}, w_{i-1})$$

   - but this doesn't capture the fact that some word strings are more/less likely to start/end a sentence
   - can assume $w_{-1} = w_0 = \langle s \rangle$ and $w_{n+1} = \langle /s \rangle$;

$$P(\vec{w}) = \prod_{i=1}^{n+1} P(w_i|w_{i-2}, w_{i-1})$$

   - typically use negative log probs since probs can get small

(4)
language models: smoothing and evaluation
- ways of choosing between models (e.g., bigram and trigram)
   - extrinsic: plug into downstream system and (somehow) evaluate that system
   - intrinsic: inherent to the current task; usually quicker/easier
- one intrinsic measure: entropy:
$$H(x) = \sum_x -P(x) \log_2 P(x) = \mathbb{E}[-\log_2 P(X)]$$

   - this measures uncertainty/disorder/surprise



$H(X) = 0 \qquad H(X) = 1 \qquad H(X) = 2 \qquad H(X) = 1.4 \qquad H(X) = 0.2$  → or, # yes/no questions

   - intuition: average number of bits needed to encode $X \approx$ entropy of $X$
- a good model assigns high probability to sequences that actually have high probability; put another way, our model should have low uncertainty (entropy) about which word comes next
   $\Rightarrow$ this can be measured using cross-entropy
   - note: cross-entropy $\geq$ entropy
   - computing per-word cross entropy:    average across words    $H_M(w_1, ..., w_n) = -\frac{1}{n} \log_2 P_M(w_1, ..., w_n)$
- perplexity: $2^{(cross-entropy)}$
   - interpretation: the average branching factor at each decision point, if our dist. were uniform
- the "goodness" of different values of these measures depends on the corpus. e.g., if corpus is "meow, meow, meow,..." then we'd hope cross-entropy would be very low (0).
- smoothing: fixes flaw of MLE, which is that it estimates probs that make training data maximally probable by making everything else (unseen data) minimally probable
   - add-one (Laplace) smoothing: just pretend you're seen everything one more time than you did
   $$P_{+1}(w_i|w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i) + 1}{C(w_{i-2}, w_{i-1}) + V}$$   where $V$ = vocab size

      - large vocabs make it so that way too much prob mass is stolen from seen events
   - add-$\alpha$ (Lidstone) smoothing improves things:

   $$P_{+\alpha}(w_i|w_{i-2}, w_{i-1}) = \frac{C(w_{i-2}, w_{i-1}, w_i) + \alpha}{C(w_{i-2}, w_{i-1}) + \alpha V}$$   where $\alpha < 1$

   2

   - to find $\alpha$: test $\alpha$ on different validation sets to see which gives lowest cross-entropy

## n-gram language models

- good-Turing
  - previous methods change the denominator, which can have big effects on frequent events; Good-Turing changes the numerator
    - usually, we have $P_{ML} = \frac{c}{n}$ → # times you see n-gram
    
    but Good-Turing uses adjusted counts $c^*$: $P_{GT} = \frac{c^*}{n}$ where $c^* = (c+1)\frac{N_{c+1}}{N_c}$,
    
    $N_c$ = # n-grams that appear exactly c times in corpus
    
    $N_0$ = # unseen n-grams
    
  - summary:
    - add-1 and add-$\alpha$ are simple, but not very good
    - Good-Turing more sophisticated and better, but there are even better ones

## More smoothing and the noisy channel model

- Good-Turing says these have same prob: "Scotting beer drinkers" and "Scottish beer eaters"
- solution: use info from lower-order N-grams "beer drinkers" and "beer eaters"
- two ways: interpolation and backoff
- interpolation: combine lower order N-gram models since they have different strengths/weaknesses
  - just a weighted average of probs from different models
  - called a mixture model
  - weights $\lambda_i$ are interpolation parameters or mixture weights
- back-off: trust the highest order language model that contains N-gram
- but diversity of histories matters: think $P(York | New)$ is high but $P(York | !New)$ is low
- Kneser-Ney Smoothing takes diversity of histories into account
  - replace raw counts with counts of histories:
  
  $$P_{ML}(w_i) = \frac{C(w_i)}{\sum_w C(w)} \implies P_{KN}(w_i) = \frac{N_{1+}(\bullet w_i)}{\sum_w N_{1+}(\bullet w)}$$
  
  where $N_{1+}(\bullet w_i) = |\{w_{i-1} : c(w_{i-1}, w_i) > 0\}|$
  
  - best smoothing for n-grams
- quick summary:
  - uniform probs: add-$\alpha$, Good-Turing
  - probs from lower-order n-grams: interpolation, backoff
  - prob of appearing in new contexts: Kneser-Ney
- now dealing with word similarity, e.g., knowing that $P(salmon | caught two)$ tells us something about $P(swordfish | caught two)$
- can use embeddings
- noisy channel model: $P(y) \underset{\text{language model}}{\implies} P(x|y) \underset{\text{noisy model}}{\implies} P(x)$ resulting dist

| Application | Y | X |
|---|---|---|
| Speech recog. | spoken words | acoustic signal |
| machine translation | words in $L_1$ | words in $L_2$ |
| spelling correction | intended words | typed words |

3

<u>More smoothing and the noisy channel model (cont'd)</u>

- we want $\arg\max_y P(y|x)$

$$\arg\max_y P(y|x) = \arg\max_y \frac{P(x|y) P(y)}{P(x)} = \arg\max_y P(x|y) P(y)$$

<u>Spelling correction, edit distance, and EM</u>

- Spelling correction:
  - we don't know $P(y)$, so need to make constraints:
    - words must differ by one character from $x$
  - then can compute $P(x|y) P(y)$
  - assume substitutions, e.g., $o \rightarrow e$, has one prob instead of a prob conditional on context (surrounding letters)

    so:
    $$P(x|y) = \prod_{i=1}^{n} P(x_i|y_i)$$

    e.g., $P(no|not) = P(n|n) P(o|o) P(-|t)$

data needed: subs probs
  - can estimate subs probs with data (confusion matrix)
  - alignments and edit distance
    - want to find <u>optimal</u> character alignment b/t two words ( fewest changes: <u>min edit distance MED</u> )

      STALL
      TALL          MED (stall, table) = 3     also:        STALL -
      TABL                                     (as          d  l  l  s |  i
      TABLE                                    alignment)    - T A B L E

    - exponential number of possibilities: how to choose?
    - use <u>dynamic programming</u> (memoization) algorithms: Viterbi, CKY
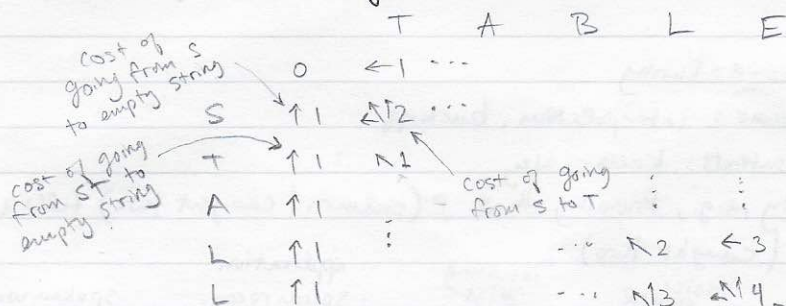    - intuition: $D$ (stall, table) must be min of:
      $D$(stall, tabl) + cost (ins)
      $D$(stal, table) + cost (del)        assume cost (ins) = cost(del) = 1 , cost (sub) = 2
      $D$(stal, tabl) + cost (sub)
    - chart stores two things: MED and <u>backpointers</u>                    start
                                                            fill like this: ↓↓↓↓↓ etc.
                      T    A    B    L    E

cost of s going from s to empty string      0  ←1 ⋯
                      S  ↖↑1 ↖↑2 ⋯
cost of going from ST to empty string   T  ↗↑1 ↖1
                      A  ↑1  (cost of going from S to T)  ⋮    ⋮
                      L  ↑1  ⋮           ⋯↑↖2 ←3
                      L  ↑1              ⋯↖3 ↖↑4        min cost: follow arrows to get possible paths back

    - chicken-and-egg: want to use costs from noise model, but need costs to estimate noise model
    - Solution: <u>EM</u> algorithm
      1. initialize params to arbitrary values (e.g., all costs = 1)
      2. using these params, compute optimal values (run MED to get alignments)
      3. using alignments, recompute params
      4. repeat 2 and 3 until params stop changing

4

spelling correction, edit distance, and EM (cont'd)

- EM just described for spelling correction is "hard" EM, which converges (like soft EM) $\longrightarrow$ to something, i.e., doesn't diverge or oscillate
  but is not nicely defined mathematically (doesn't converge to optimum of likelihood
  function like soft EM does); but it probably works fine in practice (easier to compute)
- likelihood function: $\longrightarrow$ if alignments a are latent

$$P(\text{data} \mid \theta) = \prod_{i=1}^{n} P(x_i \mid y_i) = \sum_a \prod_{i=1}^{n} P(x_i \mid y_i, a)$$

- neither hard nor soft EM guaranteed to converge to global optimum; soft maybe not even local

## text classification

- have y values like spam/not spam, sentiment, topic
- bag of words models: word order doesn't matter
- Naive Bayes

$$\hat{c} = \underset{c \in C}{\text{argmax}}\ P(c \mid d) = \underset{c \in C}{\text{argmax}}\ P(d \mid c)\, P(c)$$

category $\nearrow$ document (over $P(c|d)$); likelihood

model $P(d \mid c)$ as set of features (words) it contains: $P(d \mid c) = P(f_1, f_2, \ldots, f_n \mid c)$

$\approx P(f_1 \mid c) P(f_2 \mid c) \ldots P(f_n \mid c)$

So: $\hat{c} = \underset{c \in C}{\text{argmax}}\ P(c) \prod_{i=1}^{n} P(f_i \mid c)$ 

prior $\nearrow$ (over $P(c)$)

(naive Bayes assumption)

- can estimate $P(f_i \mid c)$ with simple smoothing
- alternative feature values and feature sets:
  - use only binary values for $f_i$
  - only use subset of vocab for F
  - use more complex features (e.g., bigrams, syntactic features)
  - for sentiment analysis: whether word is +/−
- quick and easy model
- but independence assumption may be strong/unrealistic

data: need features (can be extracted?)

- MaxEnt better if you have enough data
  - a.k.a. multinomial logistic regression

$$\hat{c} = \underset{c \in C}{\text{argmax}}\ P(c \mid d) \quad (\text{same as NB})$$

   ...but we model $P(c \mid d)$ directly

- features are functions of both observations $\vec{x}$ and class c

$$P(c \mid \vec{x}) = \frac{1}{Z} \exp\left( \sum_i w_i f_i(\vec{x}, c) \right)$$

(logistic) $\nearrow$ and using softmax

$\Rightarrow$ as if just running c separate regressions...

- training:

$$\hat{w} = \underset{w}{\text{argmax}} \sum_j \log P(c^{(j)} \mid x^{(j)})$$

can take some time (gradient ascent): conditional MLE (CMLE)

Part-of-speech tagging and HMMs
- first step toward syntactic analysis
- deciding on correct label depends on word to be labeled and labels of surrounding words
- parts of speech: two main classes: open class ("content" words) and closed class (function words)
- commonly 40-100 different tags
- hard b/c of ambiguity, sparsity

|        | sequences | hidden vars |
|--------|-----------|-------------|
| n-gram | ✓         | ✗           |
| NB     | ✗         | ✓           |
| HMM    | ✓         | ✓           |

- relationship to previous models:
  - n-gram model: a model for sequences that also makes a Markov assumption, but has no hidden vars
  - Naive Bayes: a model w/ hidden vars (classes) but no sequential dependencies
  - HMM: a model for sequences with hidden variables
- HMM:

$$\operatorname*{argmax}_{T} p(T|S) = \operatorname*{argmax}_{T} p(S|T) p(T)$$

$$P(T) = \prod_i P(t_i | t_{i-1})$$

$$P(S|T) = \prod_i P(w_i | t_i)$$

where $P(S|T) P(T) = \prod_i P(w_i | t_i) P(t_i | t_{i-1})$
$= P(S, T)$

  - too many tag sequences $T$ to enumerate
    $\Rightarrow$ use Viterbi again (see next lecture)

algorithms for HMMs
- change of notation:
  - A: transition matrix w/ possible states $q^1, q^2, \ldots, q^N$ for N tags
  - B: output matrix w/ possible outputs $o^1, o^2, \ldots, o^V$ for V words
  - $\lambda = (A, B)$: params of HMM
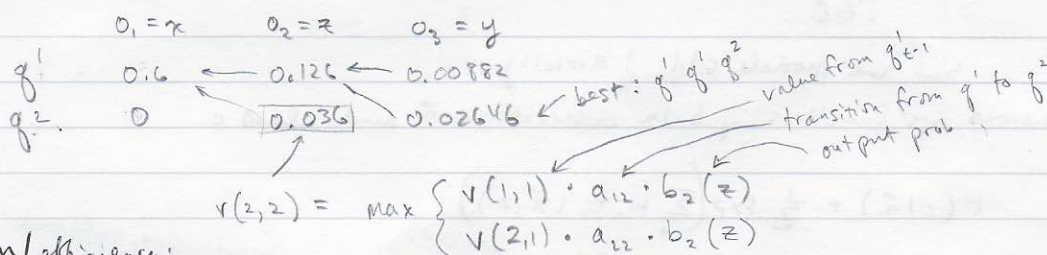  - Q: state sequence $q_1, q_2, \ldots, q_T$ } T: time
  - O: output sequence $o_1, o_2, \ldots, o_T$ } T: time
- using new notation:   $w_i$  $t_i$     $t_i$  $t_{i-1}$  $\leftarrow$ old notation

$$P(O, Q | \lambda) = \prod_{t=1}^{T} P(o_t | q_t) P(q_t | q_{t-1})$$

$$= \prod_{t=1}^{T} b_{q_t}(o_t)\, a_{q_{t-1} q_t}$$

- example grid:

|        | $o_1 = x$ | $o_2 = z$ | $o_3 = y$ |
|--------|-----------|-----------|-----------|
| $q^1$  | 0.6 ← | 0.126 ← | 0.00882 |
| $q^2$  | 0        | 0.036     | 0.02646 ← |

best: $q^1 q^1 q^2$
value from $q^{t-1}$
transition from $q^1$ to $q^2$
output prob

$$v(2,2) = \max \begin{cases} v(1,1) \cdot a_{12} \cdot b_2(z) \\ v(2,1) \cdot a_{22} \cdot b_2(z) \end{cases}$$

- implementation / efficiency:
  - enumeration is $O(N^T)$, $O(T)$ space
  - Viterbi is $O(N^2 T)$, $O(NT)$ space (much better on run time)
  - use neg log probs
- may also want to compute likelihood: $P(O|\lambda) = \sum_Q P(O, Q | \lambda)$
  - use Viterbi (w/ a slightly different update), sum last column for $P(O|\lambda)$
- learning $\lambda = (A, B)$: use EM if unsupervised (no tag labels), using expected counts in E step
- forward-backward algo: avoid enumerating all possible sequences $\Rightarrow$ an instance of EM

$\to$ i.e., a language model

$P(q^j | q^i) = \dfrac{c(q^i \to q^j)}{c(q^i)}$
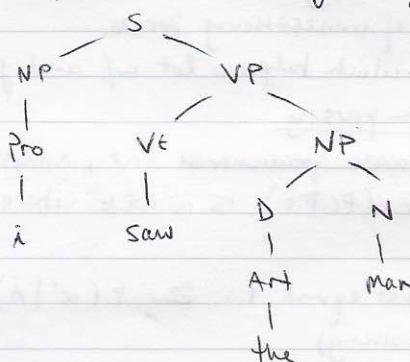
6

FNLP
final exam notes

## Methods in annotation and evaluation

- Annotation not easy, even for humans because of ambiguity and gray areas (both syntax and semantics)
- Penn Treebank has 36 POS tags (excluding punctuation) w/ extensive guidelines for annotators
  - even so, there will still be errors
- inter-annotator agreement (IAA)
  - have multiple people annotate, compute raw agreement rate
  - very low rates ⇒ examine sources of disagreement and consider rewriting guidelines
  - agreement rate can be considered upper bound (human ceiling) on accuracy of a system evaluated on that dataset
  - can also weight errors by severity
- can make hypotheses about linguistic objects (e.g., author), output of a language system, and human beings
- can use cross-validation to choose hyperparams (e.g., for smoothing)
- measures of performance:
  - precision: $P = \frac{TP}{TP+FP}$ ⇒ fixes for cases where you over-predict 1 when there are a lot of 1's
  - recall: $R = \frac{TP}{TP+FN}$ ⇒ fixes for cases where you over-predict 0 when there are few 0's
  - $F_1$-score: $F_1 = \frac{PR}{P+R}$
- use significance tests to determine whether measurements are different

## Syntax and parsing

- theory of syntax should explain which sentences are well-formed (grammatical); well-formed distinct from meaningful
  - two theories: context-free grammar and dependency grammar
- Context-free grammar (CFG)
  - two types of grammar symbols: terminals (t): words; and non-terminals (NT): phrasal categories like S, NP, VP, PP
  - has rules of the form NT→β, where β is any string of NT's and t's
  - e.g.,

```
            S
          /   \
        NP     VP
        |     /  \
       Pro   Vt   NP
        |    |    / \
        i   saw  D   N
                 |   |
                Art man
                 |
                the
```

  - structural ambiguity: caused by POS ambiguity; e.g., I saw her duck
  - attachment ambiguity: e.g., I saw the man with the telescope.
- all parsers have two fundamental properties
  - directionality: sequence in which structures are constructed (top-down, bottom-up, mixed)
  - search strategy: the order in which the search space of possible analyses is explored: depth-first, breadth-first, best-first
  - Example: the dog bit:
    (recursive descent strategy) ⇒ top-down, depth-first
    - problems: can be many backtracks or infinite loops (NP → NP PP)

data: need POS tags and rules

| Step | op | subgoals | input | operations: |
|---|---|---|---|---|
| 0 | | S | the dog bit | |
| 1 | E | NP VP | " | E = expand |
| 2 | E | DT NN VP | " | M = match |
| 3 | E | the NN VP | " | Bn = backtrack to step n |
| 4 | M | NN VP | dog bit | |
| 5 | E | bit VP | " | |
| 6 | B4 | NN VP | " | |

7

## Syntax and parsing (cont'd)

- shift-reduce: depth-first, bottom-up
  - e.g.,

- depth-first parsers very efficient for unambiguous structures but massively inefficient when faced with local ambiguity
  - but can use probabilistic model to learn which choices to make (next lesson...)
- breadth-first search using dynamic programming
  - avoid re-analyzing any substring because its analysis is independent of the rest of the parse
  - memoized in chart, a.k.a., well-formed substring table (WFST).
  - CKY algorithm: bottom-up, breadth-first
    - takes $O(Gn^3)$, where $G$ is the number of grammar rules

|   | 1 | 2 | 3 | 4 |
|---|---|---|---|---|
| 0 | Pro, NP ① | | | S ⑤ |
| 1 | | Vt, Vp, N | ② | VP ④ |
| 2 | | | Pro, PosPro, D | NP ③ |
| 3 | | | | N, Vi |

$_0 he_1$  $_1 saw_2$  $_2 her_3$  $_3 duck_4$

① NP → Pro
② D → PosPro
③ NP → D N
④ VP → V+ NP
⑤ S → NP VP

(tree: S → NP  VP)

- start with all POS allowed for word
- ignore lower-diagonal: would be backward strings
- to fill cell $(i,j)$, we use cell from row $i$ and a cell from column $j$
  ⇒ so need to fill cells down and left of $(i,j)$ before filling $(i,j)$
- even though avoids re-computing substructures, much more efficient than depth-first (in worst case)
- but still may be a lot of unnecessary parses.
- next: statistical parsing, which helps a lot w/ ambiguity and efficiency

## CKY parsing, treebanks and statistical parsing

- for probabilistic parsing, use treebank grammars (i.e., annotated sentences)
- a probabilistic context-free grammar (PCFG) is a CFG where each rule $A \to \alpha$ is assigned a probability $P(\alpha|A)$
  - sum over expansions of A must equal 1: $\sum_{\alpha'} P(\alpha'|A) = 1$
  - use MLE for probs (w/ smoothing)
  - a generative model, like with HMMs
  - prob of parse t is product of all rules of parse: $P(t) = \prod_{A \to \alpha \in t} A \to \alpha$

    *i.e., each non-terminal node has a prob; just multiply these together*
  - we also have a language model since $t$ implicitly includes words $\vec{w}$: $P(t) = P(t, \vec{w})$
    - sentence prob is then obtained by summing over $T(\vec{w})$, the set of valid parses of $\vec{w}$:
    $$P(\vec{w}) = \sum_{t \in T(\vec{w})} P(t, \vec{w}) = \sum_{t \in T(\vec{w})} P(t)$$
  - straightforward to extend CKY parsing to probabilistic case
    - goal: return highest prob parse of sentence (analogous to Viterbi)
  - best-first parsing help by not having to do exhaustive parsing
    - constituents have scores and are added to agendas, which are ordered by scores
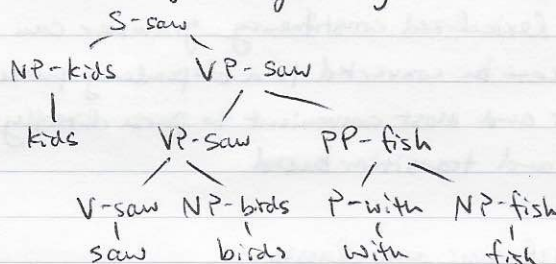    - scores computed as an average over words to normalize tree sizes
- not-so-great consequence: words w/ same POS will give same scores/probs even though we know this isn't true

heads, dependency parsing
- ways to fix PCFGs
  1. automatically create new categories that include the old category as the parent
     parent annotation →
     - e.g., an NP in subject position becomes NP^S
  2. lexicalization: create new categories by adding the lexical head of the phrase
     - e.g.,

         S-saw
        /      \
   NP-kids    VP-saw
      |       /    \
    kids   VP-saw  PP-fish
          /    \    /    \
     V-saw NP-birds P-with NP-fish
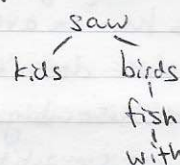      saw   birds   with   fish

- but this leads to huge grammar blowup and very sparse data
- evaluating parse accuracy
  - output considered correct if there's a gold constituent that spans the same sentence positions
  - use precision / recall / $F_1$ score
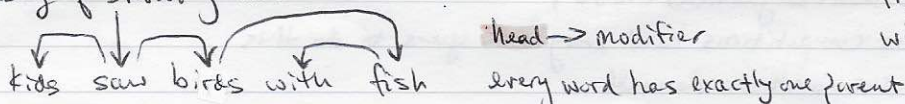- dependencies
  - transforming constituency → dependency parse
    - start w/ lexicalized constituency parse, remove phrasal categories, remove duplicated terminals, and collapse chains of duplicates
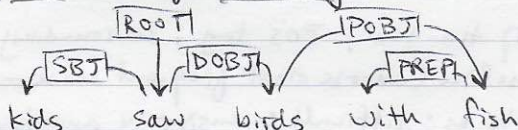    - so the lexicalized constituency parse above becomes:
                            saw
                           /    \
                        kids    birds
                                /    \
                              fish
                               |
                              with
    - another way of showing this:

      kids  saw  birds  with  fish     head → modifier
                                       every word has exactly one parent

  - some treebanks prefer content heads, while others functional heads
    functional                content
    - e.g., ...were always watching...
  - edge labels:
            [ROOT]              [POBJ]
          [SBJ]  [DOBJ]        [PREP]
       kids   saw   birds   with   fish

  - projectivity: parse said to be projective if every subtree (node and all its descendants) occupies a contiguous span of the sentence
    - i.e., no crossing edges
    - nonprojectivity is rare in English, but quite common in other languages
  - constituency → dependency parse (again): but how do we find each phrase's head in the first place?
    - use head rules: designate one RHS nonterminal as containing the head
- direct dependency parsing:
  - transition-based parsing
    - adapts shift-reduce methods: stack and buffer (note: shift-reduce more efficient than CKY)
    - idea: train classifier to predict next action (SHIFT, REDUCE, ATTACH-LEFT, ATTACH-RIGHT), and proceed left-to-right through sentence. $O(n)$ time complexity!
    - only finds projective trees
  - graph-based parsing
    - global algorithm: from the fully connected directed graph of all possible edges, choose the best ones that form a tree
    - e.g., maximum spanning tree algorithm ($O(n^2)$)

9

## heads, dependency parsing (cont'd)

- conversion-based parser: first constituency-parse, then convert to dependencies
- Summary:
  - while constituency parses give hierarchically nested phrases, dependency parses represent syntax with trees whose edges connect words in a sentence
  - head rules govern how a lexicalized constituency grammar can be extracted from a treebank, and how a constituency parse can be converted to a dependency parse
  - for English, often fastest and most convenient to parse directly to dependencies; two main paradigms: graph-based and transition-based

## lexical semantics: word senses, relations, and classes

- why it's hard: different senses, synonyms, hyponym/hypernym (subset/superset), similarity and gradation, need for inference
  - animal → Polar bear (above "hyponym/hypernym")
  - → vacation/holiday
  - good and great (top right)
- WordNet a hand-built resource w/ 117,000 synsets: sets of synonymous words
  - bank: river/finance → multiple senses
  - although this doesn't solve issues w/ homonyms, polysemes
  - Synsets organized into a network by several kinds of relations, including:
    - hypernymy (is-a): hyponym {ambulance} is a kind of hypernym car
    - meronymy (part-whole): meronym {air bag} is a part of holonym car
  - nouns have on average 1.24 senses, verbs 2.17
  - incomplete: doesn't have all multiword phrases (stress out), neologisms (facepalm), names (Microsoft)
- word sense disambiguation (WSD)
  - there are competitions held every 1-3 years to do this
  - data-driven methods do well
    - e.g., Naive Bayes, decision lists, decision trees
    - as features, can use content words in some window; can also use syntactically related words, syntactic role in sense, topic of the text, POS tag, surrounding POS tags
  - evaluation: extrinsic, intrinsic, baseline (choose most frequent sense — sometimes hard to beat)
  - issues: not clear how fine-grained to be; difficult expensive to annotate w/ fine-grained; classifiers must be trained separately for each word
- named entity recognition and supersense tagging: categories like person, location, artefact
  - → for extracting propn/names
  - supersenses (e.g., attribute, food, object, quantity, motion, social) are language-neutral

## distributional lexical semantics

- distributional hypothesis: a word is the company it keeps
- use word-context matrices (like in NLU): each word is a vector of context word frequencies/probs
- two kinds of co-occurence between two words: first order (near each other), second order (similar neighbors)
- for syntactic similarity, use small context windows; for semantic, use larger (sometimes whole document)
- collocation: occurring unusually often in the context of word w
  - can use $PMI(x,y) = \log_2 \frac{P(x,y)}{P(x)P(y)}$
    - b/c uses MLE, it's over-sensitive to chance co-occurence of infrequent words
  - alternatives: t-test, $\chi^2$, others...
  - or can improve PMI: positive PMI (PPMI) = $\max(PMI, 0)$
- measuring similarity: cosine: $\cos(\vec{v}, \vec{w}) = \frac{\vec{v} \cdot \vec{w}}{|\vec{v}||\vec{w}|}$
- sparse → dense: SVD/LSA, skip-grams, etc.
  - e.g., brown clusters: binary tree

10

## distributional lexical semantics (cont'd)

- class-based language model: $P(w_i | w_{i-1}) = P(c_i | c_{i-1}) P(w_i | c_i)$

## Semantic role labeling and argument structure

- can't just use syntax to trivially infer semantic roles
- two computational datasets/approaches that describe sentences in terms of semantic roles:
  - PropBank: simpler, more data
  - FrameNet: richer, less data
- PropBank has structure and its roles (Arg0, Arg1,...)
  - abstracts away from syntax to predicate-argument structures
  - lexicon and annotations on full WSJ PTB corpus and other data
  - limitation: roles (Arg0,...) are predicate-specific
  - originally verbs only, but now has other POS
- FrameNet                    ⟶ like PropBank
  - no longer (primarily) verb-based, but frame-based
  - FrameNet parsing: SEMAFOR
    - pipeline similar (but not identical) to traditional SRL

## Machine translation [non-examinable]