# Tagging

- used for disambiguation between more suitable meanings of one word
- between ==morphology== and ==syntax==

## Tagsets

- tag ~ one category

$$T \longleftrightarrow (C_1 \_ C_n) \quad 1:1 \text{ mapping}$$

- Penn treebank
- Brown Corpus

# Lemmatisation   ~ reduced morphological analysis

- lemma ≈ one dictionary entry ref. / lexical unit

→ MA disambiguation

$$L : A^+ \longrightarrow 2^L, \text{ even though } A^+ - L \text{ is what is wanted}$$

# Morphological analysis

- Word form list — books: book-2/VB2
                           book-1/NNS   — it is only there / not there
- direct coding   - from endings etc.      $A^+ \longrightarrow 2^{(L,e)}$
- FSM  ~ finite state machinery ...
- CFG ... - more linguistic than computational phenomena

# FST

- implemented as FSA, one symbol: $(r:s)$ from alphabets $R, S$
- can just run for acceptance or analysis          or      synthesis
                                          input: $S$              input: $r$
                                          output: $r$             output: $s$

Rule-based disambig.
  – rules using / word forms
                 / – tags
                 \ combination...

    – if-then / reg. exp.

  – the biggest problem is that the rules are created by hand

  – trying to eliminate as much possibilities as possible to have just one (correct) tag.

  – today relevant for checking simple rules with NO computational costs

HMM - tagging :

  – tags can be products of „hidden" features of the language

  – these features are represented by hidden states in HMM

  – separation of transition and emission distributions

Transformation - tagging:

    – NOT source channel view

    – NOT probabilistic

    – BUT statistical

        – uses training data to learn rules

    – criterion-based selection of rules selection

    – rules can look to the future (diff from HMM)

    – Maximum Entropy principle...

        – when constraints given, the higher the entropy, the less closed
          does in uncertainty.

Phrase structure tree

- always rooted, can be enclosed by correct bracketing

- projective

Dependency tree

- one word, one node

- not always projective

## PCFG

- Sometimes more parses are correct, but we still need to output just one

$$P(T) = \prod_i P(r_i)$$ , $r_i$ are rules used to generate the sentence.

$\llcorner\!\rightarrow$ very strong independence assumptions...

$$\sum_{r \in R_A} P(r) = 1 \quad \longrightarrow \text{aha. given left hand, all right hands must sum up to 1.}$$

estimating the probs: $\quad P(r) = c(r) / C(A) \quad$ for $\quad r: A \rightarrow \alpha$

$$P(w) = \sum_j P(T_j) \quad \sim T_j \text{ generates word } w$$

$$\beta_N(p,q) = \sum_{A,B} \sum_{d=p-q-1} \beta_A(p,d) P(N \rightarrow AB) \beta_B(d+1, q)$$

inside probability



| | dog | saw | out | with | telescope | |
|---|---|---|---|---|---|---|
| | 1 | 2 | 3 | 4 | 5 | |
| 1 | N 0,3<br>NP 0,21 | | S 0,0441 | | S 0,00966 | |
| 2 | | V 1,0 | VP 0,21 | | VP 0,018+0,028= 0,046 | this is the probability |
| 3 | | | N 0,5<br>NP 0,35 | | NP 0,03 | of the final parse. |
| 4 | | | | PREP 1,0 | PP 0,2 | |
| 5 | | | | | N 0,2<br>NP 0,14 | |

## Statistical machine translation

Entropy:
$$H(P) = - \sum_x P(x) \cdot \log P(x) \qquad NLL$$

$$Perplexity = 2^{H(P)}$$

$$H(P_{bigram}) = - \sum_x P(x,y) \cdot \log P(x|y)$$

Cross Entropy:

$$H(P,Q) = - \sum_x P(x) \cdot \log Q(x)$$

$$H(P,x) \geq H(P) \quad \forall x$$

$$H(P,P) = H(P)$$

# Maximum entropy methods

The form of the constraints for the Maximum Entropy model is defined as

Vyberte jednu nebo více možností:

✓ a. $1/|T| \sum_{t=1..T} \sum_{y \in Y} p(y|x_t) f_i(y,x_t) - d_i = 0$, where T is the training data (and $|T|$ its size), $p(y|x)$ the conditional model probability distribution, y the predicted variable and x the context ($x_t$ is the concrete context at t-th data item), and $f_i$ the i-th feature. $d_i$ is the true feature count as extracted from the training data.

✓    Yes, this is the approximation formula using the data-oriented expected value computation due to the complexity of summing over all possible xs (which is often impossible to enumerate).

✓ b. $\sum_{y,x} p(y|x) f_i(y,x) - d_i = 0$, where p is the conditional model distribution, $f_i$ are the features, and $d_i$ is the true count as extracted from the training data.

✗    No. The weight in the expected value formula computation must always be the joint distribution, not the conditional one.

✓ c. $E_p(f_i(y,x)) - d_i = 0$, where E is the expected value of feature count expressed in terms of the probability distribution p as $E_p(f_i) = \sum_{y,x} p(x,y) f_i(y,x)$, with p being the model joint distribution, and $d_i$ is the true count as extracted from the training data.

✓    Yes. It simply says that the (joint) distribution p must be such that it models the feature count in such a way that it equals to the true count as found in the data, by using the standart optimization technique known as Lagrange multipliers (where the "multipliers" then serve as the feature weights).

?

---

The model (distribution) that fulfills the constraints while maximizing the entropy has the following form:

Vyberte jednu nebo více možností:

☐ a. $p(y|x) = 1/Z(x) \, e^{\sum_{i=1..N} \lambda_i f_i(y,x)}$, where y is the predicted variable, x is the context of y used (together with y) for computing the feature values, N is the number of features used in the model, $f_i(y,x)$ are the feature functions, and $\lambda_i$ are their weights. $Z(x)$ is the normalization factor for a given context x computed either from the special $f_0$ feature or simply by normalizing to $<0..1>$.

✓ b. $E_p(f_i) = \sum_{y,x} p(x,y) f_i(y,x)$, where E is the expected value for each $f_i$, and $p(x,y)$ is the joint distribution.

✗    No, this is the expected count of the features used in the constraints where this value (expressed by means of the distribution) must equal to the empirical count from the training data (the $d_i$ constant).

✓ c. $p(y,x) = (1/Z) \, e^{\sum_{i=1..N} \lambda_i f_i(y,x)}$, where y is the predicted variable, x is the context of y used (together with y) for computing the feature values, N is the number of features used in the model, $f_i(y,x)$ are the feature functions, and $\lambda_i$ are their weights. Z is the normalization factor, computed either from the special $f_0$ feature or simply by normalizing to $<0..1>$.

✓    Yes. This is the joint probability "version". The normalization is needed for the model to be a probability distribution (and in fact that's what makes the computation, especially at training time, slow in terms of complexity based on N).

?