

How to put features into a classifier



- Linear classifiers at classification time:
 - Linear function from feature sets $\{f_i\}$ to classes $\{c\}$.
 - Assign a weight λ_i to each feature f_i .
 - We consider each class for an observed datum d
 - For a pair (c,d), features vote with their weights:

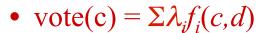
• vote(c) =
$$\sum \lambda_i f_i(c,d)$$

PERSON in Québec LOCATION in Québec DRUG in Québec

• Choose the class c which maximizes $\sum \lambda_i f_i(c,d)$



- Linear classifiers at classification time:
 - Linear function from feature sets $\{f_i\}$ to classes $\{c\}$.
 - Assign a weight λ_i to each feature f_i .
 - We consider each class for an observed datum d
 - For a pair (c,d), features vote with their weights:



PERSON in Québec



0.3 DRUG in Québec

• Choose the class c which maximizes $\sum \lambda_i f_i(c,d) = \text{LOCATION}$



There are many ways to chose weights for features

- Perceptron: find a currently misclassified example, and nudge weights in the direction of its correct classification
- Margin-based methods (Support Vector Machines)



- Exponential (log-linear, maxent, logistic, Gibbs) models:
 - Make a probabilistic model from the linear combination $\sum \lambda_i f_i(c,d)$

$$P(c \mid d, \lambda) = \frac{\exp \sum_{i} \lambda_{i} f_{i}(c, d)}{\sum_{c'} \exp \sum_{i} \lambda_{i} f_{i}(c', d)} \underbrace{\qquad \qquad \text{Makes votes positive}}_{\text{Normalizes votes}}$$

- $P(LOCATION|in\ Qu\'ebec) = e^{1.8}e^{-0.6}/(e^{1.8}e^{-0.6} + e^{0.3} + e^0) = 0.586$
- $P(DRUG|in\ Qu\'ebec) = e^{0.3}/(e^{1.8}e^{-0.6} + e^{0.3} + e^{0}) = 0.238$
- $P(PERSON|in\ Qu\'ebec) = e^0/(e^{1.8}e^{-0.6} + e^{0.3} + e^0) = 0.176$
- The weights are the parameters of the probability model, combined via a "soft max" function

Christopher Manning



Feature-Based Linear Classifiers

- Exponential (log-linear, maxent, logistic, Gibbs) models:
 - Given this model form, we will choose parameters $\{\lambda_i\}$ that maximize the conditional likelihood of the data according to this model.
 - We construct not only classifications, but probability distributions over classifications.
 - There are other (good!) ways of discriminating classes SVMs, boosting, even perceptrons but these methods are not as trivial to interpret as distributions over classes.

Christopher Manning



Aside: logistic regression

- Maxent models in NLP are essentially the same as multiclass logistic regression models in statistics (or machine learning)
 - If you haven't seen these before, don't worry, this presentation is self-contained!
 - If you have seen these before you might think about:
 - The parameterization is slightly different in a way that is advantageous for NLP-style models with tons of sparse features (but statistically inelegant)
 - The key role of feature functions in NLP and in this presentation
 - The features are more general, with f also being a function of the class –
 when might this be useful?

Christopher Manning



Quiz Question

- Assuming exactly the same set up (3 class decision: LOCATION, PERSON, or DRUG; 3 features as before, maxent), what are:
 - P(PERSON | by Goéric) =
 - P(LOCATION | by Goéric) =
 - P(DRUG | by Goéric) =
 - 1.8 $f_1(c, d) = [c = \text{LOCATION } \land w_{-1} = \text{"in"} \land \text{isCapitalized}(w)]$
 - -0.6 $f_2(c, d) = [c = LOCATION \land hasAccentedLatinChar(w)]$
 - 0.3 $f_3(c, d) = [c = DRUG \land ends(w, "c")]$



LOCATION by Goéric



$$P(c \mid d, \lambda) = \frac{\exp \sum_{i} \lambda_{i} f_{i}(c, d)}{\sum_{c'} \exp \sum_{i} \lambda_{i} f_{i}(c', d)}$$



How to put features into a classifier