# Discriminative Estimation (Maxent models and perceptron)

## Generative vs. Discriminative models

# Introduction

- So far we've looked at "generative models"
  - Naive Bayes
- But there is now much use of conditional or discriminative probabilistic models in NLP, Speech, IR (and ML generally)
- Because:
  - They give high accuracy performance
  - They make it easy to incorporate lots of linguistically important features

# Joint Models

- We have some data $\{(d, c)\}$ of paired observations $d$ and hidden classes $c$.

- Joint (generative) models place probabilities over both observed data and the hidden stuff (gene-rate the observed data from hidden stuff):

  $P(c, d)$

  - All the classic StatNLP models:
    - $n$-gram models, Naive Bayes classifiers, hidden Markov models, probabilistic context-free grammars, IBM machine translation alignment models

# Conditional Models

- Discriminative (conditional) models take the data as given, and put a probability over hidden structure given the data:

  $P(c|d)$

  - Logistic regression, conditional loglinear or maximum entropy models, conditional random fields
  - Also, SVMs, (averaged) perceptron, etc. are discriminative classifiers (but not directly probabilistic)

# Joint Likelihood vs. Conditional Likelihood

- A *joint* model gives probabilities $P(d,c)$ and tries to maximize this joint likelihood.
  - It turns out to be trivial to choose weights: just relative frequencies.
- A *conditional* model gives probabilities $P(c|d)$. It takes the data as given and only models the conditional probability of the class.
  - Harder to do.
  - More closely related to classification error.

# Maxent Models and Discriminative Estimation

Generative vs. Discriminative models

# The Maxent Model

# Example features

- $f_1(c, d) \equiv [c = \text{LOCATION} \wedge w_{-1} = \text{"in"} \wedge \text{isCapitalized}(w)]$    weight: 1.8
- $f_2(c, d) \equiv [c = \text{LOCATION} \wedge \text{hasAccentedLatinChar}(w)]$    weight: -0.6
- $f_3(c, d) \equiv [c = \text{DRUG} \wedge \text{ends}(w, \text{"c"})]$    weight: 0.3

LOCATION    LOCATION          DRUG    PERSON
*in Arcadia*    *in Québec*          *taking Zantac*    *saw Sue*

- Models will assign to each feature a *weight:*
  - A positive weight votes that this configuration is likely correct
  - A negative weight votes that this configuration is likely incorrect

# The Maxent Model

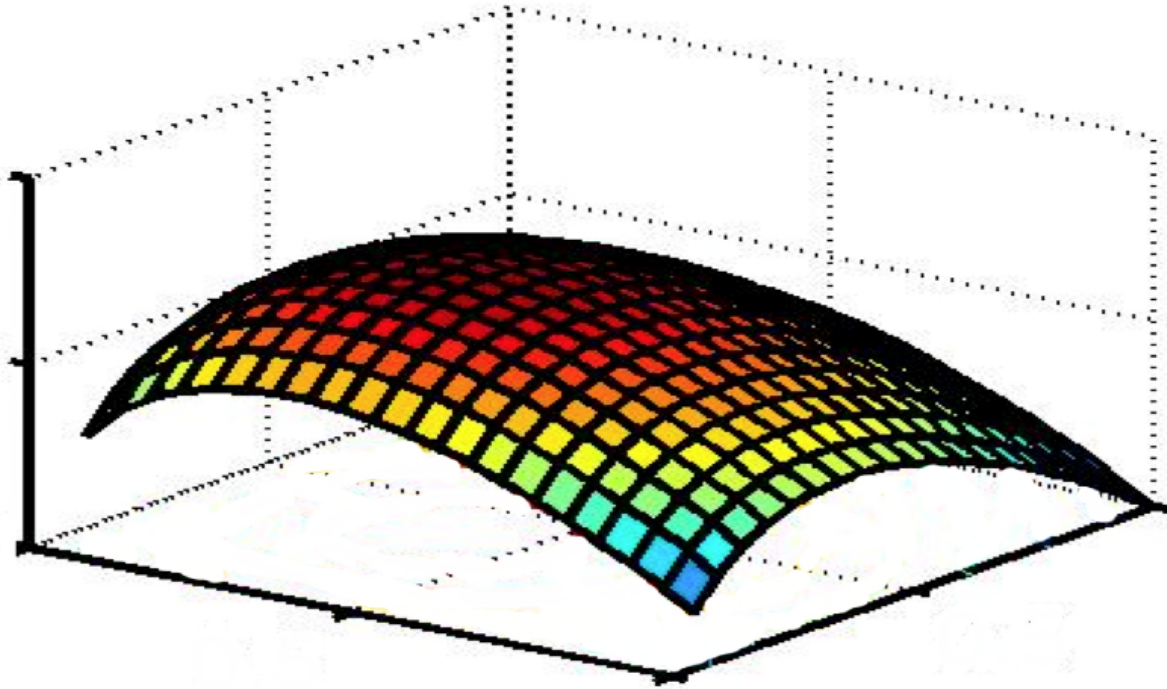- Exponential (log-linear, maxent, logistic, Gibbs) models:

$$P(c \mid d, \lambda) = \frac{\exp \sum_i \lambda_i f_i(c, d)}{\sum_{c'} \exp \sum_i \lambda_i f_i(c', d)}$$

Makes votes positive

Normalizes votes

- P(LOCATION|*in Québec*) = $e^{1.8}e^{-0.6}/(e^{1.8}e^{-0.6} + e^{0.3} + e^0)$ = 0.586
- P(DRUG|*in Québec*) = $e^{0.3}/(e^{1.8}e^{-0.6} + e^{0.3} + e^0)$ = 0.238
- P(PERSON|*in Québec*) = $e^0/(e^{1.8}e^{-0.6} + e^{0.3} + e^0)$ = 0.176

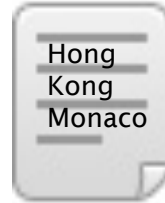# A likelihood surface

# Naive Bayes vs. Maxent Models

- Naive Bayes models multi-count correlated evidence
  - Each feature is multiplied in, even when you have multiple features telling you the same thing

- Maximum Entropy models (pretty much) solve this problem
  - this is done by weighting features, avoid to assign equally high weights to correlated features.

# Text classification: Asia or Europe

# Perceptron

Another Discriminative
Learning algorithm

# Perceptron Algorithm

- Algorithm is Very similar to logistic regression
- Not exactly computing gradients

Initalize weight vector w = 0

Loop for K iterations

    Loop For all training examples x_i

        if sign(w * x_i) != y_i

            w += (y_i - sign(w * x_i)) * x_i

# Regularization in the Perceptron Algorithm

- run different numbers of iterations
- Use parameter averaging, for instance, average of all parameters after seeing each data point