



Introduction

- So far we've looked at "generative models"
 - Language 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
 - They allow automatic building of language independent, retargetable NLP modules



Joint vs. Conditional 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 (generate the observed data from hidden stuff):
 - All the classic StatNLP models:
 - n -gram models, Naive Bayes classifiers, hidden Markov models, probabilistic context-free grammars, IBM machine translation alignment models

$$P(c, d)$$



Joint vs. Conditional Models

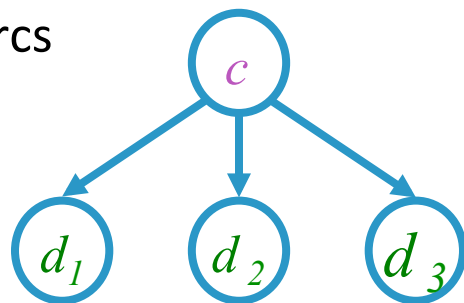
- **Discriminative (conditional) models** take the data as given, and put a probability over hidden structure given the data:
 - Logistic regression, conditional loglinear or maximum entropy models, conditional random fields
 - Also, SVMs, (averaged) perceptron, etc. are discriminative classifiers (but not directly probabilistic)

$$P(c|d)$$



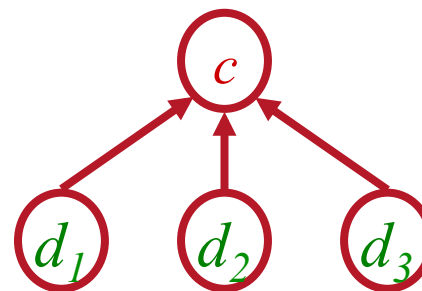
Bayes Net/Graphical Models

- Bayes net diagrams draw circles for random variables, and lines for direct dependencies
- Some variables are observed; some are hidden
- Each node is a little classifier (conditional probability table) based on incoming arcs



Naive Bayes

Generative



Logistic Regression

Discriminative



Conditional vs. Joint 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 models only the conditional probability of the class.
 - We seek to maximize conditional likelihood.
 - Harder to do (as we'll see...)
 - More closely related to classification error.



Conditional models work well: Word Sense Disambiguation

Training Set	
Objective	Accuracy
Joint Like.	86.8
Cond. Like.	98.5

Test Set	
Objective	Accuracy
Joint Like.	73.6
Cond. Like.	76.1

- Even with exactly the same features, changing from joint to conditional estimation increases performance
- That is, we use the same smoothing, and the same word-class features, we just change the numbers (parameters)

(Klein and Manning 2002, using Senseval-1 Data)



Maxent Models and Discriminative Estimation

Generative vs. Discriminative models

Christopher Manning