

CS778 Fall 2006

## Primer on Generative vs. Discriminative Learning

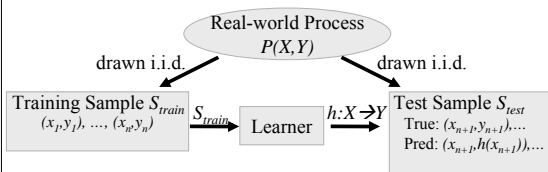
Thorsten Joachims

Cornell University  
Department of Computer Science

## Project

- **Project Topic**
  - Self-defined: comparison of methods, application of exiting method to new domain, explore modifications of existing method, etc.
  - Can be centered around a paper, but should envision some originality
  - I will help refine the topic
- **Project Proposal**
  - Outline general idea for project
  - Discuss resources and timeline
  - Due: September 19 (in class)

## Supervised Learning



- **Learning Task:**  $P(X,Y) = P(X) P(Y|X)$ 
  - Input Space:  $X$  (e.g. feature vectors, word sequence, etc.)
  - Output Space:  $Y$  (e.g. class label  $l..k$ )
  - Training Data:  $S_{train} = ((x_1, y_1), \dots, (x_n, y_n)) \sim_{iid} P(X, Y)$
- **Goal:** Find  $h: X \rightarrow Y$  with low prediction error  $Err_P(h)$

## Generalization Error and Sample Error

**Definition:** The prediction error/generalization error/true error/expected loss/risk  $Err_P(h)$  of a hypothesis  $h$  for a learning task  $P(X, Y)$  is

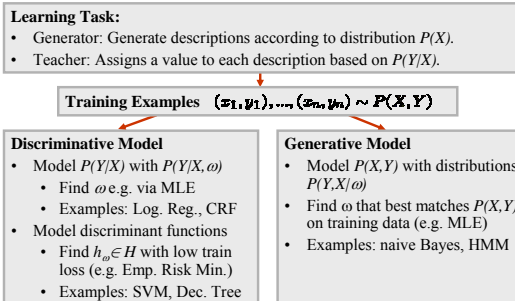
$$Err_P(h) = \sum_{\bar{x} \in X, y \in Y} \Delta(h(\bar{x}), y) P(X = \bar{x}, Y = y).$$

**Definition:**  $\Delta(a, b)$  is a loss function that measures the cost of making a wrong prediction. A commonly used loss function is the 0/1-loss

$$\Delta(a, b) = \begin{cases} 0 & \text{if } (a == b) \\ 1 & \text{else} \end{cases}$$

**Definition:** The error on sample  $S$   $Err_S(h)$  of a hypothesis  $h$  is  $Err_S(h) = \frac{1}{n} \sum_{i=1}^n \Delta(h(\bar{x}_i), y_i)$ .

## Generative vs. Discriminative Models



## Generative Model: Model P(X,Y)

- **Bayes' Decision Rule:** Optimal Decision is
 
$$h(x) = \underset{y \in Y}{\operatorname{argmin}} \left[ \sum_{y' \in Y} \Delta(y, y') P(Y = y' | X = x) \right]$$
- **Equivalent Reformulations:** For 0/1-Loss  $\Delta(y, y') = 1$ , if  $y \neq y'$ , 0 else
 
$$h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(Y = y | X = x)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} \left[ \frac{P(X = x | Y = y) P(Y = y)}{P(X = x)} \right]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(X = x | Y = y) P(Y = y)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(X = x, Y = y)]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [P(Y = y | X = x) P(X = x)]$$
- **Learning: maximum likelihood (or MAP, or Bayesian)**
  - Assume model class  $P(X, Y | \omega)$  with parameters  $\omega \in \Omega$
  - Find  $\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y = y_i, X = x_i | \omega)]$

## Naïve Bayes' Classifier (Multivariate)

- **Input Space X: Feature Vector**

- **Output Space Y: {1,-1}**

- **Model:**

- Prior class probabilities

$$P(Y = +1) \quad P(Y = -1)$$

- Class conditional model (one for each class)

$$P(X=x|Y=+1) = \prod_{j=1}^n P(X^{(j)}=x^{(j)}|Y=+1)$$

$$P(X=x|Y=-1) = \prod_{j=1}^n P(X^{(j)}=x^{(j)}|Y=-1)$$

- **Classification rule:**

$$h_{naive}(x) = \underset{y \in \{+1, -1\}}{\operatorname{argmax}} \left\{ P(Y=y) \prod_{j=1}^n P(X^{(j)}=x^{(j)}|Y=y) \right\}$$

fever	cough	pukes	flu?
(3)	(2)	(2)	
high	yes	no	1
high	no	yes	1
low	yes	no	-1
low	yes	yes	1
high	no	yes	???

## Estimating the Parameters of Naïve Bayes

- **Count frequencies in training data**

- $n$ : number of training examples
- $n_+$  /  $n_-$ : number of pos/neg examples
- $\#(X^{(j)}=x^{(j)}, y)$ : number of times feature  $X^{(j)}$  takes value  $x^{(j)}$  for examples in class  $y$
- $|X^{(j)}|$ : number of values attribute of  $X^{(j)}$

fever	cough	pukes	flu?
(3)	(2)	(2)	
high	yes	no	1
high	no	yes	1
low	yes	no	-1
low	yes	yes	1
high	no	yes	???

- **Estimating:  $\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n \left[ P(Y=y) \prod_{j=1}^n \left[ P(X_i^{(j)}=x_i^{(j)}|Y=y) \right] \right]$**

- P(Y): Maximum Likelihood Estimate

$$P(Y=1) = \frac{n_+}{n} \quad P(Y=-1) = \frac{n_-}{n}$$

- P(X|Y): Maximum Likelihood Estimate

$$P(X^{(j)}=x^{(j)}|Y=y) = \frac{\#(X^{(j)}=x^{(j)}, y)}{n_y}$$

- P(X|Y): Smoothing with Laplace estimate

$$P(X^{(j)}=x^{(j)}|Y=y) = \frac{\#(X^{(j)}=x^{(j)}, y) + 1}{n_y + |X^{(j)}|}$$

## Generative vs. Discriminative Models

### Learning Task:

- Generator: Generate descriptions according to distribution  $P(X)$ .
- Teacher: Assigns a value to each description based on  $P(Y|X)$ .

Training Examples  $(x_1, y_1), \dots, (x_n, y_n) \sim P(X, Y)$

### Discriminative Model

- Model  $P(Y|X)$  with  $P(Y|X, \omega)$ 
  - Find  $\omega$  e.g. via MLE
  - Examples: Log. Reg., CRF
- Model discriminant functions
  - Find  $h_\omega \in H$  with low train loss (e.g. Emp. Risk Min.)
  - Examples: SVM, Dec. Tree

### Generative Model

- Model  $P(X, Y)$  with distributions  $P(Y, X|\omega)$ 
  - Find  $\omega$  that best matches  $P(X, Y)$  on training data (e.g. MLE)
  - Examples: naive Bayes, HMM

## Discriminative Model: Model $P(Y|X)$

- **Bayes' Decision Rule:**

$$\text{General: } h(x) = \underset{y \in Y}{\operatorname{argmin}} \left[ \sum_{y' \in Y} \Delta(y, y') P(Y=y'|X=x) \right]$$

- Assume 0/1 Loss  $\Delta(y, y') = 1$ , if  $y \neq y'$ , 0 else

$$h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(Y=y|X=x)]$$

- **Learning: maximum likelihood (or MAP, or Bayesian)**

- Assume model class  $P(Y|X, \omega)$  with parameters  $\omega \in \Omega$

- Find

$$\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y=y_i|X=x_i, \omega)]$$

## Logistic Regression/"Maximum Entropy"

- **Assume:**

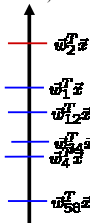
$$P(Y=y|X=x, \omega = (w_1, \dots, w_k)) = \frac{e^{w_y^T x}}{\sum_{y' \in Y} e^{w_{y'}^T x}}$$

→ Learn one weight vector  $w_y$  for each class  $y \in Y$  (linear discriminant)

$$h(x) = \underset{y \in Y}{\operatorname{argmax}} [w_y^T x]$$

- **Maximum Likelihood training:**

$$\begin{aligned} \omega' &= \underset{\omega \in \Omega}{\operatorname{argmax}} \left[ \prod_{i=1}^n P(Y=y_i|X=x_i, \omega) \right] \\ &= \underset{\omega \in \Omega}{\operatorname{argmax}} \left[ \sum_{i=1}^n \log(P(Y=y_i|X=x_i, \omega)) \right] \\ &= \underset{(w_1, \dots, w_k)}{\operatorname{argmax}} \left[ \sum_{i=1}^n w_{y_i}^T x_i - \log \left( \sum_{y' \in Y} e^{w_{y'}^T x_i} \right) \right] \end{aligned}$$



## Generative vs. Discriminative Models

### Learning Task:

- Generator: Generate descriptions according to distribution  $P(X)$ .
- Teacher: Assigns a value to each description based on  $P(Y|X)$ .

Training Examples  $(x_1, y_1), \dots, (x_n, y_n) \sim P(X, Y)$

### Discriminative Model

- Model  $P(Y|X)$  with  $P(Y|X, \omega)$ 
  - Find  $\omega$  e.g. via MLE
  - Examples: Log. Reg., CRF
- Model discriminant functions
  - Find  $h_\omega \in H$  with low train loss (e.g. Emp. Risk Min.)
  - Examples: SVM, Dec. Tree

### Generative Model

- Model  $P(X, Y)$  with distributions  $P(Y, X|\omega)$ 
  - Find  $\omega$  that best matches  $P(X, Y)$  on training data (e.g. MLE)
  - Examples: naive Bayes, HMM

### Discriminative Model: Model Discriminant Function $h$ Directly

- **Discriminant Function:**  $h_\omega: X \times Y \rightarrow \mathfrak{R}$ 

$$h(x) = \underset{y \in Y}{\operatorname{argmin}} \left[ \sum_{y \in Y} \Delta(y, y) P(Y=y | X=x) \right]$$

$$= \underset{y \in Y}{\operatorname{argmax}} [h(x, y)] \quad (\text{e.g. } h(x, y) = [w_y^T x])$$
- **Consistency of Empirical Risk:**
  - Training Error (i.e. Empirical Risk):  $Err_S(h) = \sum_{i=1}^n \Delta(y_i, h(x_i))$
  - For sufficiently "small"  $H_\Omega$  and "large"  $S$ : Rule  $h' \in H_\Omega$  with best  $Err_S(h')$  has  $Err_P(h')$  close to  $\min_{h \in H_\Omega} [Err_P(h)]$
- **Learning: Empirical Risk Minimization (ERM)**
  - Assume class  $H_\Omega$  of discriminant functions  $h_\omega: X \rightarrow Y$
  - Find  $h'_\omega = \underset{h_\omega \in H_\Omega}{\operatorname{argmin}} \left[ \sum_{i=1}^n \Delta(y_i, h_\omega(x_i)) \right]$

### Support Vector Machine

- **Training Examples:**  $(x_1, y_1), \dots, (x_n, y_n)$   $x \in \mathfrak{R}^N$   $y \in \{+1, -1\}$
- **Hypothesis Space:**  $H_\Omega = \{h(x) = \operatorname{sgn}[w^T x + b] : \|w\| < \infty\}$
- **Training Loss:**  $Err_{S_{\min}}(h) = \sum_{i=1}^n \Delta(y_i, h(x_i)) \leq \sum_{i=1}^n \xi_i$

**Optimization Problem:**

$$\min_{w, \xi, b} \frac{1}{2} w^T w + C \sum_{i=1}^n \xi_i$$

s.t.  $y_1 (w^T x_1 + b) \geq 1 - \xi_1$   
 $\dots$   
 $y_n (w^T x_n + b) \geq 1 - \xi_n$

[Vapnik et al.]

[Crammer & Singer 02]

Training: Find  $\langle w_1, \dots, w_k \rangle$  that solve

$$\min_{w_1, \dots, w_k, \xi} \sum_{i=1}^k w_i^2 + C \sum_{i=1}^n \xi_i$$

s.t.  $\forall j \neq y_1: w_{y_1}^T x_1 \geq w_j^T x_1 + 1 - \xi_1$   
 $\dots$   
 $\forall j \neq y_n: w_{y_n}^T x_n \geq w_j^T x_n + 1 - \xi_n$

### So what about Complex Outputs?

- **Approach: view as multi-class classification task**
  - Every complex output  $y \in Y$  is one class

$X$  The bear chased the cat  $\rightarrow$   $Y$  Det  $\rightarrow$  N  $\rightarrow$  V  $\rightarrow$  Det  $\rightarrow$  N

- **Problem: Exponentially many classes!**
  - Generative Model:  $P(X, Y | \omega)$
  - Discriminative Model:  $P(Y | X, \omega)$
  - Discriminant Functions:  $h_\omega: X \times Y \rightarrow \mathfrak{R}$
- **Challenges**
  - How to compactly represent model?
  - How to do efficient inference with model (i.e.  $\underset{y \in Y}{\operatorname{argmax}} [h(x, y)]$ )?
  - How to effectively estimate model from data? (e.g. compute  $h'_\omega = \underset{h_\omega \in H_\Omega}{\operatorname{argmin}} \sum_{i=1}^n \Delta(y_i, h_\omega(x_i))$ )

### Predicting Sequences: Hidden Markov Model

- Bayes rule:  $h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(X=x | Y=y) P(Y=y)]$
- Assumptions for compact representation
 
$$P(Y = (y^{(1)}, \dots, y^{(t)})) = \prod_{t=0}^{t-1} P(y_t = y^{(t)} | y_p = y^{(t-1)})$$

$$P(X = (x^{(1)}, \dots, x^{(t)})) = \prod_{t=1}^t P(x_t = x^{(t)} | y_c = y^{(t)})$$

	Y	Det	→	N	→	V	→	Det	→	N
X	↓	↓	↓	↓	↓	↓	↓	↓	↓	↓
	The bear chased the cat									

- **Prediction rule:**

$$h(x) = \underset{y \in Y}{\operatorname{argmax}} [P(X=x | Y=y) P(Y=y)]$$

$$= \underset{(y^{(1)}, \dots, y^{(t)}) \in Y}{\operatorname{argmax}} \left[ \prod_{t=1}^t P(y_t = y^{(t)} | y_p = y^{(t-1)}) P(x_t = x^{(t)} | y_c = y^{(t)}) \right]$$

→ Viterbi (dynamic prog) algorithm computes argmax efficiently

### Predicting Sequences: HMM Training

- **Maximum Likelihood (or alternative estimator)**
  - Find  $\omega' = \underset{\omega \in \Omega}{\operatorname{argmax}} \prod_{i=1}^n [P(Y = y_i, X = x_i | \omega)]$ 

$$= \underset{\omega \in \Omega}{\operatorname{argmax}} \left[ \prod_{t=1}^n \prod_{j=1}^t P(y_t = y^{(t)} | y_p = y^{(t-1)}) P(x_t = x^{(t)} | y_c = y^{(t)}) \right]$$
- **For Hidden Markov Model**
  - Closed-form solutions
 
$$P(Y_c = y_c | Y_p = y_p) = \frac{\# \text{ of Times State A Follows State B}}{\# \text{ of Times State B Occurs}}$$

$$P(X_c = x_c | Y_c = y_c) = \frac{\# \text{ of Times Output A Is Observed In State B}}{\# \text{ of Times State B Occurs}}$$
  - Need for smoothing the estimates (e.g. max a posteriori)

## Questions in this Class

- **Important applications for which conventional methods don't fit!**
  - Noun-phrase co-reference: two step approaches of pair-wise classification and clustering as postprocessing, e.g [Ng & Cardie, 2002]
  - Directly optimize complex loss functions (e.g. F1, AvgPrec)
- **Improve upon existing methods!**
  - Part-of-Speech Tagging: generative models like HMM
  - Is it possible to train discriminatively? Log Regression, SVMs, Boosting, etc.
  - SVM outperforms naïve Bayes for text classification [Joachims, 1998] [Dumais et al., 1998]
- **More**
  - 
  -
- **Training**
  - 
  -
- Support Vector Machines

Precision/Recall Break-Even Point	Naïve Bayes	Linear SVM
Reuters	72.1	87.5
WebKB	82.0	90.3
Ohsumed	62.4	71.6

## Reading

- **Generative vs. Discriminative Training**
  - Tony Jebara, Machine Learning: Discriminative and Generative, Kluwer, 2004. (Chapters 1 and 2)
- **Hidden-Markov Models**
  - C. Manning and H. Schuetze, Foundations of Statistical NLP, MIT Press, 1999.
- **Logistic Regression**
  - T. Hastie, R. Tibshirani, and J. Friedman, The Elements of Statistical Learning, Springer, 2001. (Section 4.4)
- **Support Vector Machines**
  - B. Schoelkopf and A. Smola, Learning with Kernels, MIT Press, 2001. (Chapters 1 and 7, available at <http://www.learning-with-kernels.org/>)