# Introduction to Machine Learning

Brown University CSCI 1950-F, Spring 2012 Prof. Erik Sudderth

Lecture 4:

Probability: Continuous Random Variables Naïve Bayes: ML & Bayesian Estimation Preview of Decision Theory

> Many figures courtesy Kevin Murphy's textbook, Machine Learning: A Probabilistic Perspective

## Bayes Rule (Bayes Theorem)

- $\theta \longrightarrow$  unknown parameters (many possible models)
- $\mathcal{D} \longrightarrow$  observed data available for learning
- $p(\theta) \longrightarrow$  prior distribution (domain knowledge)
- $p(\mathcal{D} \mid \theta) \longrightarrow$  likelihood function (measurement model)
- $p(\theta \mid \mathcal{D}) \longrightarrow \text{posterior distribution (learned information)}$

$$p(\theta, \mathcal{D}) = p(\theta)p(\mathcal{D} \mid \theta) = p(\mathcal{D})p(\theta \mid \mathcal{D})$$
$$p(\theta \mid \mathcal{D}) = \frac{p(\theta, \mathcal{D})}{p(\mathcal{D})} = \frac{p(\mathcal{D} \mid \theta)p(\theta)}{\sum_{\theta' \in \Theta} p(\mathcal{D} \mid \theta')p(\theta')}$$
$$\propto p(\mathcal{D} \mid \theta)p(\theta)$$

#### **Continuous Random Variables**



#### **Moments of Random Variables**

#### Mean

$$\mathbb{E}[X] \triangleq \int_{\mathcal{X}} x \ p(x) dx \qquad \qquad \mathbb{E}[X] \triangleq \sum_{x \in \mathcal{X}} x \ p(x)$$

ſ

#### Variance

$$\operatorname{var} [X] \triangleq \mathbb{E} \left[ (X - \mu)^2 \right] = \int (x - \mu)^2 p(x) dx$$
$$= \int x^2 p(x) dx + \mu^2 \int p(x) dx - 2\mu \int x p(x) dx = \mathbb{E} \left[ X^2 \right] - \mu^2$$
$$\mathbb{E} \left[ X^2 \right] = \mu^2 + \sigma^2 \qquad \text{second moment}$$
$$\operatorname{std} [X] \triangleq \sqrt{\operatorname{var} [X]} \qquad \text{standard deviation}$$

Moments & Conditional Moments

$$\mathbb{E}[g(X)] = \int_{\mathcal{X}} g(x)p(x) \, dx \quad \mathbb{E}[g(X) \mid Y = y] = \int_{\mathcal{X}} g(x)p(x \mid y) \, dx$$

#### Gaussian (Normal) Distributions



 $Mode[X] = \arg \max_{x \in \mathbb{R}} \mathcal{N}(x \mid \mu, \sigma^2) = \mu$ 

#### Gaussian (Normal) Distributions



Summaries: Mean, median, mode, variance, standard deviation

#### **Learning Binary Probabilities**

# Bernoulli Distribution: Single toss of a (possibly biased) coin $Ber(x \mid \theta) = \theta^{\mathbb{I}(x=1)} (1-\theta)^{\mathbb{I}(x=0)} \quad 0 \le \theta \le 1$

 Suppose we observe N samples from a Bernoulli distribution with unknown mean:

$$X_i \sim \text{Ber}(\theta), i = 1, \dots, N$$
$$p(x_1, \dots, x_N \mid \theta) = \theta^{N_1} (1 - \theta)^{N_0}$$
$$N_1 = \sum_{i=1}^N \mathbb{I}(x_i = 1) \qquad N_0 = \sum_{i=1}^N \mathbb{I}(x_i = 0)$$

• What is the *maximum likelihood* parameter estimate?

$$\hat{\theta} = \arg\max_{\theta} \log p(x \mid \theta) = \frac{N_1}{N}$$

#### **Beta Distributions**



#### **Beta Distributions**



### **Bayesian Learning of Probabilities**

Bernoulli Likelihood: Single toss of a (possibly biased) coin

$$\operatorname{Ber}(x \mid \theta) = \theta^{\mathbb{I}(x=1)} (1-\theta)^{\mathbb{I}(x=0)} \quad 0 \le \theta \le 1$$
$$p(x_1, \dots, x_N \mid \theta) = \theta^{N_1} (1-\theta)^{N_0}$$

#### **Beta Prior Distribution:**

$$p(\theta) = \text{Beta}(\theta \mid a, b) \propto \theta^{a-1} (1-\theta)^{b-1}$$

**Posterior Distribution:** 

 $p(\theta \mid x) \propto \theta^{N_1 + a - 1} (1 - \theta)^{N_0 + b - 1} \propto \text{Beta}(\theta \mid N_1 + a, N_0 + b)$ 

- This is a conjugate prior, because posterior is in same family
- Estimate by posterior mode (MAP) or mean (preferred)
- Here, posterior predictive equivalent to mean estimate

#### **Sequence of Beta Posteriors**





#### **Constrained Optimization**



- Solution:  $\hat{\theta}_k = \frac{a_k}{a_0}$   $a_0 = \sum_{k=1}^K a_k$
- Proof for K=2: Change of variables to unconstrained problem
- Proof for general K: Lagrange multipliers (see textbook)

## Learning Categorical Probabilities

Multinoulli Distribution: Single roll of a (possibly biased) die  $\operatorname{Cat}(x \mid \theta) = \prod_{k=1}^{K} \theta_k^{x_k}$   $\mathcal{X} = \{0, 1\}^K, \sum_{k=1}^{K} x_k = 1$ 

• If we have  $N_k$  observations of outcome k in N trials:

$$p(x_1,\ldots,x_N \mid \theta) = \prod_{k=1}^{K} \theta_k^{N_k}$$

• The *maximum likelihood* parameter estimates are then:

$$\hat{\theta} = \arg\max_{\theta} \log p(x \mid \theta) \qquad \qquad \hat{\theta}_k = \frac{N_k}{N}$$

• Will this produce sensible predictions when *K* is large?



#### **Dirichlet Probability Densities**



#### **Dirichlet Samples**



## **Bayesian Learning of Probabilities**

Multinoulli Distribution: Single roll of a (possibly biased) die

$$\operatorname{Cat}(x \mid \theta) = \prod_{k=1}^{K} \theta_k^{x_k} \qquad \mathcal{X} = \{0, 1\}^K, \sum_{k=1}^{K} x_k = 1$$
$$p(x_1, \dots, x_N \mid \theta) = \prod_{k=1}^{K} \theta_k^{N_k}$$

K

**Dirichlet Prior Distribution:** 

$$p(\theta) = \operatorname{Dir}(\theta \mid \alpha) \propto \prod_{k=1}^{\infty} \theta_k^{\alpha_k - 1}$$

**Posterior Distribution:** 

$$p(\theta \mid x) \propto \prod_{k=1}^{K} \theta_k^{N_k + \alpha_k - 1} \propto \text{Dir}(\theta \mid N_1 + \alpha_1, \dots, N_K + \alpha_K)$$

• This is a conjugate prior, because posterior is in same family

 Assume we have N training examples independently sampled from an unknown naïve Bayes model:

 $\begin{array}{l} y_i \longrightarrow \text{ observed class label for training example } i \\ x_{ij} \longrightarrow \text{ value of feature } j \text{ for training example } i \\ p(\theta \mid y, x) \propto p(\theta) \prod_{\substack{n=1 \\ N}} p(y_i \mid \theta) p(x_i \mid y_i, \theta) \\ \propto p(\theta) \prod_{\substack{i=1 \\ N}} p(y_i \mid \theta) \prod_{\substack{j=1 \\ i=1}}^{D} p(x_{ij} \mid y_i, \theta) \end{array}$ 

• Learning: ML estimate, MAP estimate, or posterior prediction

**Naïve Bayes:** ML & Bayes  $p(\mathbf{x}_{i}, y_{i}|\boldsymbol{\theta}) = p(y_{i}|\boldsymbol{\pi}) \prod_{j} p(x_{ij}|\boldsymbol{\theta}_{j}) = \prod_{c} \pi_{c}^{\mathbb{I}(y_{i}=c)} \prod_{j} \prod_{c} p(x_{ij}|\boldsymbol{\theta}_{jc})^{\mathbb{I}(y_{i}=c)}$   $\log p(\mathcal{D}|\boldsymbol{\theta}) = \sum_{c=1}^{C} N_{c} \log \pi_{c} + \sum_{j=1}^{D} \sum_{c=1}^{C} \sum_{i:y_{i}=c} \log p(x_{ij}|\boldsymbol{\theta}_{jc})$   $N_{c} \longrightarrow \text{ number of examples of training class } c$ 

 Maximizing the sum of functions of independent parameters can be done by maximizing them independently:

 Similarly, if the parameters for different features are independent under the prior, they remain independent under the posterior, and Bayesian analysis decomposes

#### **Generative Classifiers**

• Compute class *posterior distribution* via Bayes rule:

$$p(y = c \mid x, \theta) = \frac{p(y = c \mid \theta)p(x \mid y = c, \theta)}{\sum_{c'=1}^{C} p(y = c' \mid \theta)p(x \mid y = c', \theta)}$$

- *Inference:* Find label distribution for some input example
- Classification: Make decision based on inferred distribution
- Learning: Estimate parameters heta from labeled training data

## **Decision Theory**

- $y \in \mathcal{Y} \longrightarrow$  unknown hidden state of "nature"
- $x \in \mathcal{X} \longrightarrow$  observed data
- $a \in \mathcal{A} \longrightarrow$  set of possible actions we can take

 $L(y, a) \longrightarrow$  real-valued loss function: the price we pay if we choose action *a*, and *y* is the true hidden state

- Goal: Choose the action which minimizes the expected loss  $\delta(\mathbf{x}) = \operatorname*{argmin}_{a \in \mathcal{A}} \mathbb{E}\left[L(y, a)\right] \qquad \delta \, : \, \mathcal{X} \, \to \, \mathcal{A}$ 
  - Some averaging is necessary because we don't know y
  - Two notions of expectation: Bayesian versus frequentist
- Some communities speak of maximizing expected utility, which is equivalent if utility equals negative loss

#### **Losses for Classification**



Most common default choice is the 0-1 loss:

$$L(y,a) = \mathbb{I}(y \neq a) = \begin{cases} 0 & \text{if } a = y \\ 1 & \text{if } a \neq y \end{cases}$$

• For the special case of binary classification:

predicted	true label $y$	
label $\hat{y}$	0	1
0	0	$\lambda_{01}$
1	$\lambda_{10}$	0