# CS 8501 Advanced Topics in Machine Learning

#### **Lecture 02: Generative Modeling**

Yangfeng Ji Information and Language Processing Lab Department of Computer Science University of Virginia https://yangfengji.net/

## **Generative Modeling**

#### $p(y \mid x; \theta) \propto p(y; \theta) \cdot p(x \mid y; \theta)$

- heta represents all the parameters in this model
- $p(x \mid y; heta)$ : likelihood
- $p(y; \theta)$ : prior

### **Concept Learning**

- Pick a concept
- Give some examples of this concept
- Ask someone whether a new example belongs to this concept

## **Concept Learning in Generative Modeling**

Given

- y: concept
- $\mathcal{D} = \{x_i\}$ : observations

Learning

 $p(y \mid \mathcal{D}; heta) \propto p(y; heta) \cdot p(\mathcal{D} \mid y; heta)$ 

#### **Number Game**

- Assume all observed numbers are randomly drawn from  $\{1,\ldots,100\}$  with equal chance
- Hypothesis space  $\mathcal{H}$ : a set of hypotheses
- Version space  $\mathcal{V}$ : the subset of  $\mathcal{H}$  that is consistent with  $\mathcal{D}$ . For example, if  $\mathcal{D}=\{2\}$ 
  - $\circ h_{
    m two}$  = "powers of two"
  - $h_{\rm even}$  = "even numbers"

#### Likelihood

The likelihood function of each hypothesis is

$$p(\mathcal{D} \mid h) = ig(rac{1}{|h|}ig)^N$$

- |h|: the size of numbers can be explained by this hypothesis  $\circ$  For example,  $|h_{
  m two}|=6$
- N: the size of observation  ${\cal D}$

#### Example

- Given  $\mathcal{D} = \{16, 8, 2, 64\}$ ,
  - With hypothesis  $h_{
    m two}$

$$p(\mathcal{D} \mid h_{ ext{two}}) = ig(rac{1}{6}ig)^4$$

- With hypothesis  $h_{
m even}$ 

$$p(\mathcal{D} \mid h_{ ext{even}}) = ig(rac{1}{50}ig)^4$$

#### **Occam's Razor**

- [Mackay, 2006]: "Accept the simplest explanation that fits the data"
- For example, how many boxes in the following image?



### **Model Selection**



• William of Ockham: "Entities are not to be multiplied without necessity"

## **Implication of Occam's Razor**

There are many implications of Occam's Razor in machine learning research. For example,

- You should always select the simpliest models that can solve the problem
- However, it also means you should understand your research problem (or your data)

#### **Example (Cont.)**

Given  $\mathcal{D} = \{16, 8, 2, 64\}$ , now let's compare two similar hypotheses

•  $h_{
m two}$  = "power of two"

$$p(\mathcal{D} \mid h_{ ext{two}}) = ig(rac{1}{6}ig)^6$$

•  $h_{
m another}$  = "power of two except 4 and 32"

$$p(\mathcal{D} \mid h_{ ext{another}}) = ig(rac{1}{4}ig)^4$$

The limitation of will be addressed by the Bayesian version of Occam's razor

# **About Maximizing the Likelihood**

If the goal is solely about maximizing the likelihood function, then we may pick the hypothesis that explains the current data **too well**.

- This is overfitting
- The simple explanation applies to any other learning scenarios

#### **Prior**

Follow the previous discussion, and consider the two hypotheses:

- *h*: "powers of two"
- h': "powers of two except 4 and 32"
- With the previous discussion on likelihood functions,  $h^\prime$  is more likely
- However, in practice, hypotheses like  $h^\prime$  is more complicated to implement

• Or "conceptually unnatural", as discussed in the textbook

# Subjectivity

In Bayesian modeling

- Prior is the mechanism by which background knowledge can be brought to bear on a problem
  - It is also the key of "rapid learning" (e.g., learning with small sample sizes)
- The choice of prior sometimes is subjective
  - The subjectivity is a controversial issue in Bayesian modeling

#### **Posterior**

The posterior distribution of a hypothesis given  ${\cal D}$  is

$$p(h \mid \mathcal{D}) = rac{p(\mathcal{D} \mid h)p(h)}{p(\mathcal{D})}$$

where

- $p(\mathcal{D}) = \sum_{h'} p(\mathcal{D} \mid h') p(h')$  is the major challenge in Bayesian inference
- If we need a single hypothesis from the posterior distribution, we can use the MAP estimate

$$\hat{h}_{ ext{MAP}} = ext{argmax}_h p(h \mid \mathcal{D}) = ext{argmax}_h p(\mathcal{D} \mid h) p(h)$$

### Example



16

### **Example (Cont.)**



17

#### **Posterior Predictive Distribution**

Provide a way to prodict the next number

$$p( ilde{x} \in C \mid \mathcal{D}) = \sum_{h} p(y = 1 \mid ilde{x}, h) p(h \mid \mathcal{D})$$

Instead of considering one single hypothesis h for the prediction, it *averages* the possibility of all hypotheses.

#### **Example of Number Prediction**



# **The Beta-Binomial model**

#### **Binomial Distribution**

Let  $X_i \sim \text{Bernoulli}(\theta)$ , where  $p(X_i = 1) = \theta$  and  $p(X_i = 0) = 1 - \theta$ , then the probability of  $\sum_{i=1}^n X_i = k$  is

$$p(k \mid n, heta) = inom{n}{k} heta^k (1 - heta)^{n-k}$$

- Example: tossing a coin n times, the probability of getting the head k times

#### **Prior**

A popular prior distribution for  $\boldsymbol{\theta}$  is the Beta distribution

$$p( heta;\gamma_1,\gamma_2) \propto heta^{\gamma_1-1}(1- heta)^{\gamma_2-1}$$

where  $\gamma_1$  and  $\gamma_2$  are the parameters for the prior.

• Regarding  $\theta$ , these two are called **hyper-parameters** 

#### **A Formal Definition**

$$p( heta;\gamma_1,\gamma_2)=B(\gamma_1,\gamma_2) heta^{\gamma_1-1}(1- heta)^{\gamma_2-1}$$

where  $B(\gamma_1, \gamma_2)$  is the Beta function, which is also the normalization consistent.

#### Mean

$$E( heta) = \int_{ heta} heta p( heta; \gamma_1, \gamma_2) = rac{\gamma_1}{\gamma_1 + \gamma_2}$$

#### Mode

$$\hat{ heta} = rg\max_{ heta} p( heta; \gamma_1, \gamma_2) = rac{\gamma_1 - 1}{\gamma_1 + \gamma_2 - 2}$$

#### **Beta Distribution**

With different  $\gamma_1$  and  $\gamma_2$  ( $\alpha$  and  $\beta$  in the following plot)



#### **Posterior**

The posterior distribution of  $\theta$ 

$$p( heta \mid n,k;\gamma_1,\gamma_2) = rac{p( heta;\gamma_1,\gamma_2)p(k \mid n, heta)}{p(k \mid n;\gamma_1,\gamma_2)}$$

where

$$p(k \mid n; \gamma_1, \gamma_2) = \int_{ heta} p(k, heta \mid n; \gamma_1, \gamma_2) d heta$$

# **Posterior (II)**

Without the denominator, we have

$$p( heta \mid n,k;\gamma_1,\gamma_2) \propto heta^k (1- heta)^{n-k} \cdot heta^{\gamma_1-1} (1- heta)^{\gamma_2-1}$$

or

$$p( heta \mid n,k;\gamma_1,\gamma_2) \propto heta^{k+\gamma_1-1}(1- heta)^{n-k+\gamma_2-1}$$

• Beta distribution is the conjugate prior, because the posterior has the same form as the prior

# **The Dirichlet-Multinomial Model**

### **Distributions**

Multinomial distribution

$$p(x; heta) = rac{n!}{x_1!\cdots x_K!} heta_1^{x_1}\cdots heta_k^{x_K}$$

where x and  $\theta$  are both K-dimensional vectors, and  $\sum_{k=1}^{K} \theta_k = 1$ Dirichlet distribution

$$p( heta;lpha) = rac{1}{B(lpha)} \prod_{k=1}^K heta_k^{lpha_k-1}$$

where  $\alpha$  is a K-dimensional vector too, with  $\alpha_k > 0$ 

#### **Dirichlet Distribution**



# **Dirichlet Distribution (II)**



# **Naive Bayes Classifiers**

### **Problem Setup**

Consider a classification problem, with  $x \in \mathbb{R}^N$  as input and  $y \in \{0,1\}$  as output

• Focus on the posterior distribution of y, instead of the model parameter heta

 $p(y \mid x; heta)$ 

• We can also give  $\theta$  a prior, which is called Bayesian naive Bayes classifier (section 3.5.1.2)

# Likelihood

Given  $\{(x_i,y_i)\}_{i=1}^N$  , where  $x_i \in \mathbb{R}^K$  is a K-dimensional vector.

Navie Bayes assume that different dimensions in input are independent from each other

$$p(x_i \mid y_i; heta_{x \mid y}) = \prod_{j=1}^K p(x_{i,j} \mid y_i; heta_{x \mid y,j})$$

- This is a **naive** assumption
- Choose  $p(x_{i,j} \mid y_i, heta_{x \mid y,j})$  based on your inputs. For example,
  - Continuous variables: Gaussian
  - Discrete variables: Binomial or Multinomial

# Likelihood (II)

The overall likelihood given the training set  $\{(x_i,y_i)\}_{i=1}^N$  is defined as

$$ext{lik}( heta_{x|y}) = \prod_{i=1}^N \prod_{j=1}^K p(x_{i,j} \mid y_i; heta_{x|y,j})$$

#### **Prior**

 $p(y; heta_y)$ 

Choices of the distribution:

- Uniform distribution
- Bernoulli distribution with the parameter estimated from data

#### **MAP Estimate**

$$\log p(y; heta_y) + \sum_{i=1}^N \sum_{j=1}^K \log p(x_{i,j} \mid y_i; heta_{x \mid y, j})$$

Re-arrange it:

$$\log p(y; heta_y) + \sum_{j=1}^K \{\sum_{i=1}^N \log p(x_{i,j} \mid y_i; heta_{x|y,j})\}$$

We need to solve K + 1 one-dimensional problems

# **Thank You!**