# Introduction to Machine Learning

Statistics

Mingchen Gao

**Outline**

# Contents

# 1 Generative Models for Discrete Data

- **X**, feature vector, represents the data with multiple discrete attributes
  It is a discrete random variable because it can take $2^d$ values if there are $d$ binary attributes.

- $Y$ represents the class

**Most probable class**
$$P(Y = c | \mathbf{X} = \mathbf{x}, \boldsymbol{\theta}) \propto P(\mathbf{X} = \mathbf{x} | Y = c, \boldsymbol{\theta}) P(Y = c, \boldsymbol{\theta})$$

- $p(\mathbf{x} | y = c, \boldsymbol{\theta})$ - **class conditional density**

- How is the data distributed for each class?

# 2 Steps for Learning a Generative Model

- Example: $D$ is a sequence of $N$ binary values (0s and 1s) (coin tosses)

- What is the best distribution that could describe $D$?

- What is the probability of observing a *head* in future?

**Step 1: Choose the form of the model**

- Hypothesis Space - All possible distributions

  - Too complicated!!

- Revised hypothesis space - All Bernoulli distributions ($X \sim Ber(\theta), 0 \le \theta \le 1$)

  - $\theta$ is the hypothesis
  - Still infinite ($\theta$ can take infinite possible values)

- Likelihood of $D$
$$p(D | \theta) = \theta^{N_1} (1 - \theta)^{N_0}$$

**Maximum Likelihood Estimate**

$$
\begin{aligned}
\hat{\theta}_{MLE} &= \arg\max_{\theta} p(D | \theta) = \arg\max_{\theta} \theta^{N_1} (1 - \theta)^{N_0} \\
&= \frac{N_1}{N_0 + N_1}
\end{aligned}
$$

To compute MLE we set the derivative of the likelihood with respect to $\theta$ to 0.

$$
\begin{aligned}
\frac{d}{d\theta}p(D|\theta) &= \frac{d}{d\theta}\theta^{N_1}(1-\theta)^{N_0} \\
&= N_1\theta^{N_1-1}(1-\theta)^{N_0} - N_0\theta^{N_1}(1-\theta)^{N_0-1} \\
&= \theta^{N_1-1}(1-\theta)^{N_0-1}(N_1(1-\theta) - N_0\theta)
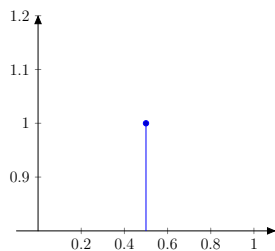\end{aligned}
$$

Setting above to zero:

$$
\begin{aligned}
\theta^{N_1-1}(1-\theta)^{N_0-1}(N_1(1-\theta) - N_0\theta) &= 0 \\
N_1(1-\theta) = N_0\theta \\
\theta = \frac{N_1}{N_0 + N_1}
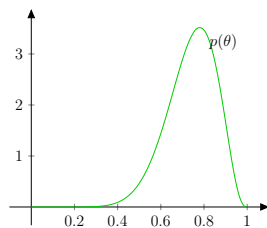\end{aligned}
$$

- **We can stop here (MLE approach)**

- Probability of getting a head next:

$$
p(x^* = 1|D) = \hat{\theta}_{MLE}
$$

## 2.1   Incorporating Prior

- Prior *encodes* our prior belief on $\theta$

- How to set a Bayesian prior?

    1. A point estimate: $\theta_{prior} = 0.5$
    2. A probability distribution over $\theta$ (**a random variable**)
        - Which one?
        - For a bernoulli distribution $0 \leq \theta \leq 1$
        - *Beta* Distribution

## 2.2   Beta Distribution

- Continuous random variables defined between 0 and 1

$$Beta(\theta|a,b) \triangleq p(\theta|a,b) = \frac{1}{B(a,b)}\theta^{a-1}(1-\theta)^{b-1}$$

- $a$ and $b$ are the (hyper-)parameters for the distribution

- $B(a,b)$ is the **beta function**

$$B(a,b) = \frac{\Gamma(a)\Gamma(b)}{\Gamma(a+b)}$$

$$\Gamma(x) = \int_0^\infty u^{x-1}e^{-u}du$$
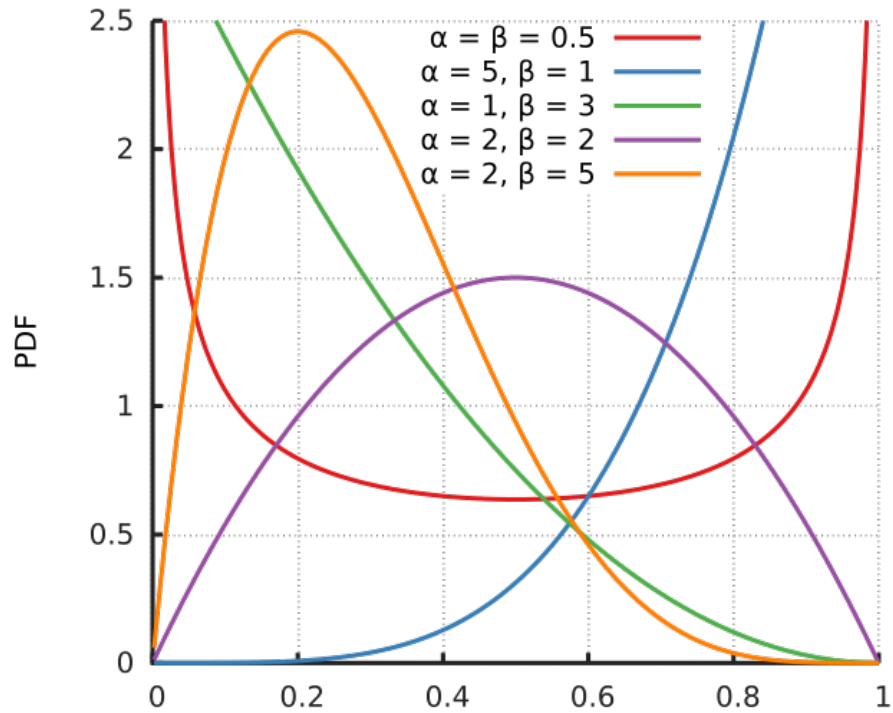
If $x$ is integer

$$\Gamma(x) = (x-1)!$$

- "Control" the shape of the pdf

The *gamma function* is an extension of factorial to real and complex numbers. By varying $a$ and $b$, one can set any prior on $\theta$, including a uniform prior, a close to point estimate, and a Gaussian prior.

- **We can stop here as well (prior approach)**

$$p(x^* = 1) = \theta_{prior}$$

4

## 2.3   Conjugate Priors

- Another reason to choose Beta distribution

$$p(D|\theta) = \theta^{N_1}(1-\theta)^{N_0}$$

$$p(\theta) \propto \theta^{a-1}(1-\theta)^{b-1}$$

- Posterior $\propto$ Likelihood $\times$ Prior

$$
\begin{aligned}
p(\theta|D) \;&\propto\; \theta^{N_1}(1-\theta)^{N_0}\theta^{a-1}(1-\theta)^{b-1} \\
&\propto\; \theta^{N_1+a-1}(1-\theta)^{N_0+b-1}
\end{aligned}
$$

- **Posterior has same form as the prior**

- Beta distribution is a conjugate prior for Bernoulli/Binomial distribution

Conjugate priors are widely used because they simplify the math and are easy to interpret.

## 2.4   Estimating Posterior

- Posterior

$$p(\theta|D) \quad \propto \quad \theta^{N_1+a-1}(1-\theta)^{N_0+b-1}$$
$$= \quad Beta(\theta|N_1 + a, N_0 + b)$$

- After observing $N$ trials in which we observe $N_1$ heads and $N_0$ tails, we update our belief as:

$$\mathbb{E}[\theta|D] = \frac{a + N_1}{a + b + N}$$

- We know that posterior over $\theta$ is a beta distribution

- MAP estimate

$$\hat{\theta}_{MAP} \quad = \quad \arg\max_{\theta} p(\theta|a + N_1, b + N_0)$$
$$= \quad \frac{a + N_1 - 1}{a + b + N - 2}$$

- What happens if $a = b = 1$?

- **We can stop here as well (MAP approach)**

- Probability of getting a head next:

$$p(x^* = 1|D) = \hat{\theta}_{MAP}$$

Using $a = b = 1$ means using an uninformative prior, which essentially reduces the MAP estimate to MLE estimate.

## 2.5   Using Predictive Distribution

- All values of $\theta$ are possible

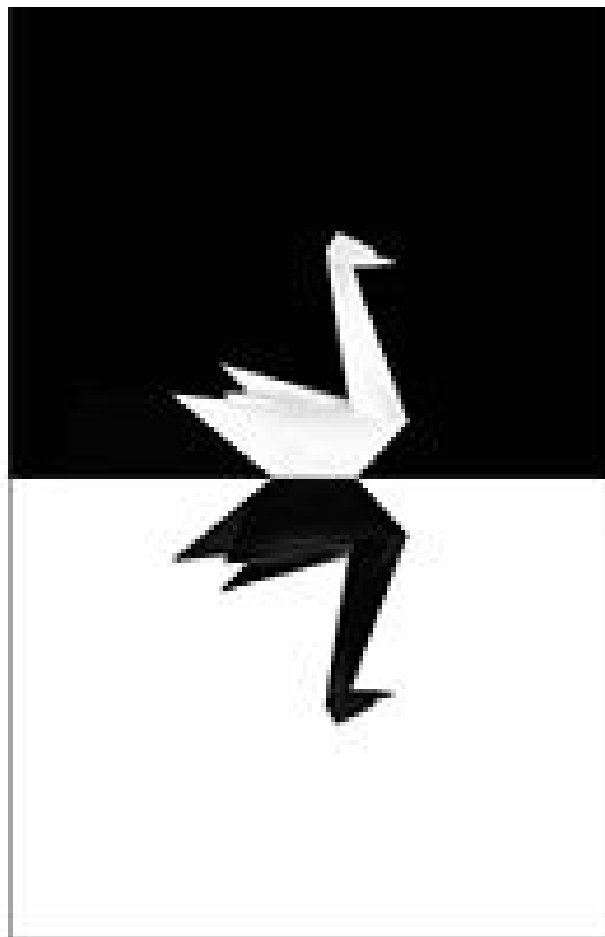- Prediction on an unknown input $(x^*)$ is given by *Bayesian Averaging*

$$
\begin{aligned}
p(x^* = 1|D) &= \int_0^1 p(x^* = 1|\theta)p(\theta|D)d\theta \\
&= \int_0^1 \theta Beta(\theta|a + N_1, b + N_0) \\
&= \mathbb{E}[\theta|D] \\
&= \frac{a + N_1}{a + b + N}
\end{aligned}
$$

- This is same as using $\mathbb{E}[\theta|D]$ as a point estimate for $\theta$

## 2.6   Need for Prior

- Why use a *prior?*

- Consider $D = $ `tails, tails, tails`

- $N_1 = 0, N = 3$

- $\hat{\theta}_{MLE} = 0$

- $p(x^* = 1|D) = 0!!$

  - Never observe a heads
  - The *black swan* paradox

- How does the Bayesian approach help?

$$
p(x^* = 1|D) = \frac{a}{a + b + 3}
$$

The black swan paradox (made famous by an eponymous book by Taleb) essentially states that since one does not observe a phenomenon in the past, he/she incorrectly induces that it can never occur.

## 2.7   Need for Bayesian Averaging

- MAP is only one part of the posterior

  - $\theta$ at which the posterior probability is maximum
  - But is that enough?
  - What about the posterior variance of $\theta$?

$$var[\theta|D] = \frac{(a + N_1)(b + N_0)}{(a + b + N)^2(a + b + N + 1)}$$

- If variance is high then $\theta_{MAP}$ is not trustworthy

- Bayesian averaging helps in this case

## 2.8 Likelihood

- Why choose one hypothesis over other?

- Avoid **suspicious coincidences**

- Choose concept with higher *likelihood*

$$p(D|h) = \prod_{x \in D} p(x|h)$$

- *Log Likelihood*
$$\log p(D|h) = \sum_{x \in D} \log p(x|h)$$

- Always choose the simpler explanation

- A general problem-solving philosophy

There are many ways to describe the Occam's Razor principle. In simple words, if there are two possible explanations for a certain phenomenon, Occam's Razor advocates choosing the "simpler" explanation.

**Maximum A Priori Estimate**
$$\hat{h}_{prior} = \arg\max_{h} p(h)$$

**Maximum Likelihood Estimate (MLE)**

$$\begin{aligned}
\hat{h}_{MLE} &= \arg\max_{h} p(D|h) = \arg\max_{h} \log p(D|h) \\
&= \arg\max_{h} \sum_{x \in D} \log p(x|h)
\end{aligned}$$

**Maximum a Posteriori (MAP) Estimate**
$$\hat{h}_{MAP} = \arg\max_{h} p(D|h)p(h) = \arg\max_{h} (\log p(D|h) + \log p(h))$$

9

- $\hat{h}_{prior}$ - Most likely hypothesis based on prior

- $\hat{h}_{MLE}$ - Most likely hypothesis based on evidence

- $\hat{h}_{MAP}$ - Most likely hypothesis based on posterior

$$\hat{h}_{prior} = \arg\max_{h} \log p(h)$$

$$\hat{h}_{MLE} = \arg\max_{h} \log p(D|h)$$

$$\hat{h}_{MAP} = \arg\max_{h} (\log p(D|h) + \log p(h))$$

MLE and MAP give the most likely hypothesis before and after considering the prior.

- As data increases, MAP estimate converges towards MLE

  - Why?

- MAP/MLE are **consistent estimators**

  - If concept is in $\mathcal{H}$, MAP/ML estimates will converge

- If $c \notin \mathcal{H}$, MAP/ML estimates converge to $h$ which is closest possible to the truth

As we have seen in our numbers example, MAP estimate can be written as the sum of log likelihood and log prior for each hypothesis. As data increases, the log likelihood will increase while the log prior will stay constant. Eventually, enough data will overwhelm the prior.

## 2.9 Posterior Predictive Distribution

- New input, $x^*$

- What is the probability that $x^*$ is also generated by the same concept as $D$?

  - $P(x^* \in C | x^*, D)$?

- **Option 0:** Treat $h^{prior}$ as the true concept
$$P(x^* \in C | x^*, D) = P(x^* \in h^{prior} | x^*, h^{prior})$$

- **Option 1:** Treat $h^{MLE}$ as the true concept
$$P(x^* \in C | x^*, D) = P(x^* \in h^{MLE} | x^*, h^{MLE})$$

- **Option 2:** Treat $h^{MAP}$ as the true concept
$$P(x^* \in C | x^*, D) = P(x^* \in h^{MAP} | x^*, h^{MAP})$$

- **Option 3:** *Bayesian Averaging*
$$P(x^* \in C | x^*, D) = \sum_h P(x^* \in h | x^*, h) p(h|D)$$

Posterior provides a notion of *belief* about the world. How does one use it? One possible use is to estimate if a new input example belongs to the same concept as the training data, $D$.

Bayesian averaging assumes that every hypothesis in $\mathcal{H}$ is possible, but with different probabilities. So the output is also a probability distribution.

# 3 Learning Gaussian Models

- pdf for MVN with $d$ dimensions:
$$\mathcal{N}(\mathbf{x}|\boldsymbol{\mu}, \boldsymbol{\Sigma}) \triangleq \frac{1}{(2\pi)^{d/2}|\boldsymbol{\Sigma}|^{1/2}} exp\left[-\frac{1}{2}(\mathbf{x}-\boldsymbol{\mu})^\top \boldsymbol{\Sigma}^{-1}(\mathbf{x}-\boldsymbol{\mu})\right]$$

## 3.1 Estimating Parameters

**Problem Statement**
Given a set of $N$ **independent and identically distributed** (iid) samples, $D$, learn the parameters $(\boldsymbol{\mu}, \boldsymbol{\Sigma})$ of a Gaussian distribution that generated $D$.

- MLE approach - maximize log-likelihood

- Result
$$\widehat{\boldsymbol{\mu}}_{MLE} = \frac{1}{N}\sum_{i=1}^{N} \mathbf{x_i} \triangleq \bar{\mathbf{x}}$$

$$\widehat{\boldsymbol{\Sigma}}_{MLE} = \frac{1}{N}\sum_{i=1}^{N} (\mathbf{x_i} - \bar{\mathbf{x}})(\mathbf{x_i} - \bar{\mathbf{x}})^\top$$

# References

Chapter 4.1 - 4.2.5, 4.6 Murphy Book