

Introduction

- 1) Chapter Scope
 - A. Examples of probability distributions
 - B. Their properties
- 2) Purpose of Introducing Distributions
 - A. a building blocks for more complex models
 - B. a recipe to discuss some essential statistical concept, e.g., Bayesian inference
 - C. to model the probability distribution $p(\mathbf{x})$, i.e., density estimation
 - * Model Selection becomes an issue since density estimation is fundamentally ill-posed problem in that infinitely many distributions can fit the observed data set.
- 3) Parametric distribution vs. Non-Parametric distribution
 - A. Parametric distribution
 - i. binomial distribution, multinomial distribution, Gaussian distribution (continuous R.V.)
 - ii. For density estimation, the parameters shall be determined with an observed data set.
 1. Frequentist: specific values for parameters (earned by optimizing some criterion, e.g., likelihood function)
 2. Bayesian: estimate posterior distribution with introduced prior distributions over the parameters as well as the observed data
 - iii. Conjugate Priors: To simplify the Bayesian analysis, use conjugate prior which let posterior distribution be in the same form of prior distribution.
 1. Exponential family of distributions is presented as it possesses a number of important properties.
 - B. Non-Parametric distribution
 - i. Distribution form is not forced by a user but typically depends on the size of the data set
 - ii. Still has the parameters but they do not determine the distribution form but the complexity
 - iii. Histogram, nearest-neighbors, kernels

Table. 1 Conjugate prior with posterior distribution in exponential family

| Conjugate Prior | Posterior Distribution |
|------------------------|--------------------------|
| Dirichlet distribution | Multinomial distribution |
| Gaussian distribution | Gaussian distribution |

1. Binary Variables

1) Bernoulli distribution

A. Definition

$$\text{Bern}(x|\mu) = \mu^x(1-\mu)^{1-x}, \text{ where } 0 \leq \mu < 1 \text{ and } x \in \{0, 1\}$$

B. Properties

$$E[x] = \mu$$

$$\text{var}[x] = \mu(1-\mu)$$

C. Density estimation

$$\mathcal{D} = \{x_1, \dots, x_N\}$$

i. Frequentist

1. Estimate μ by maximizing the likelihood function, i.e., maximize the log of likelihood

$$\ln p(\mathcal{D}|\mu) = \sum_{n=1}^N \ln p(x_n|\mu) = \sum_{n=1}^N \{x_n \ln \mu + (1-x_n) \ln(1-\mu)\}$$

2. The above log likelihood function depends on the N observations only through their sum, i.e., *sufficient statistics*: $\sum_n x_n$.

3. $\mu_{ML} = \frac{m}{N} = \text{sample mean}$

ii. Bayesian

1. Flip a coin 3 times resulting all heads \rightarrow what is the reasonable prediction? (overfitting)

2) Binomial distribution

A. Definition

$$\text{Bin}(m|N, \mu) = \binom{N}{m} \mu^m (1-\mu)^{N-m}$$

B. Properties

For independent events, 1) the means of the sum is the sum of the mean and 2) the variance of the sum is the sum of the variance

$$E[m] = N\mu$$

$$\text{var}[m] = N\mu(1-\mu)$$

1.1 The beta distribution (conjugate prior for the binomial distribution)

1) Motivation for the conjugate prior distribution

- A. Prior distribution is required in order to develop a Bayesian treatment.
- B. Make posterior distribution have the same functional form as the prior (conjugacy).

2) Definition

$$\text{Beta}(\mu|a, b) = \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} \mu^{a-1} (1-\mu)^{b-1}, \text{gamma coefficient for the normalization purpose}$$

$$\text{gamma function: } \Gamma(x) \equiv \int_0^{\infty} u^{x-1} e^{-u} du$$

$$\Gamma(x+1) = x\Gamma(x), \Gamma(1) = 1, \Gamma(x+1) = x!$$

a and b controls the distribution of the parameter μ , and thus called the *hyperparameters*

3) Posterior distribution

$$p(\mu|m, l, a, b) = \text{Beta}(\mu|a, b) \times \text{Bin}(m|N, \mu) = \frac{\Gamma(m+a+l+b)}{\Gamma(m+a)\Gamma(l+b)} \mu^{m+a-1} (1-\mu)^{l+b-1}$$

$$l = N - m = \# \text{ of tails}$$

$$m = \# \text{ of heads}$$