

# Homework 1:

## Probability Review and Priors

Huilin Chang  
hc5hq@virginia.edu

### 1 (15)

You are a data Scientist and are choosing between three approaches, A, B, and C to a problem. With approach A you will spend a total of four days coding and running an algorithm and it will not produce useful results. With approach B you will spend a total of three days coding and running an algorithm and it will not produce useful results. With approach C you will spend over day coding and running an algorithm and it will give the results you are looking for. You are equality likely to choose among unselected options. What is the expected time in days for you to obtain the results you are looking for? What is the variance on this time?

Response:

Considering all the possible approach chains

If starting with approach A:

- A-B-C:  $4+3+1 = 8$  days
- A-C:  $4+1 = 5$  days

If starting with approach B:

- B-A-C:  $3+4+1 = 8$  days
- B-C :  $3+1 = 4$  days

If staring with approach C

- C: need one day = 1 days

Let E to represent the days to solve the question = {8, 5, 8, 4, 1} and the mean value is 5.2

The variance of time based on the variance formula:

$$S^2 = \frac{\sum (x_i - \bar{x})^2}{n-1} = \frac{34.8}{4} = 8.7, \text{ the variance} = 2.95$$

## 2 (15)

Suppose if it is sunny or not in Charlottesville depends on the weather of the last three days. Show how this can be modeled as Markov chain.

Response:

We can think about “Markov chain” in this way

- The next term in a sequence could depend on all the previous terms
  - If it only depends on the previous term it is called “first-order” Markov
  - If it depends on the two previous terms it is “second-order” Markov

We can form a Markov chain as follows, take the weather states R(rain), (nice), S(Sunny) to form transition probabilities.

$$P = \begin{bmatrix} P_{11} & P_{12} & P_{13} \\ P_{21} & P_{22} & P_{23} \\ P_{31} & P_{32} & P_{33} \end{bmatrix}$$

$$P_{ij} \geq 0, i = 1, \dots, n, j = 1, \dots, n \text{ and } \sum_{j=1}^n P_{ij} = 1$$

For example:

R      N      S

$$P = \begin{matrix} & \begin{matrix} R & N & S \end{matrix} \\ \begin{matrix} R \\ N \\ S \end{matrix} & \begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix} \end{matrix}$$

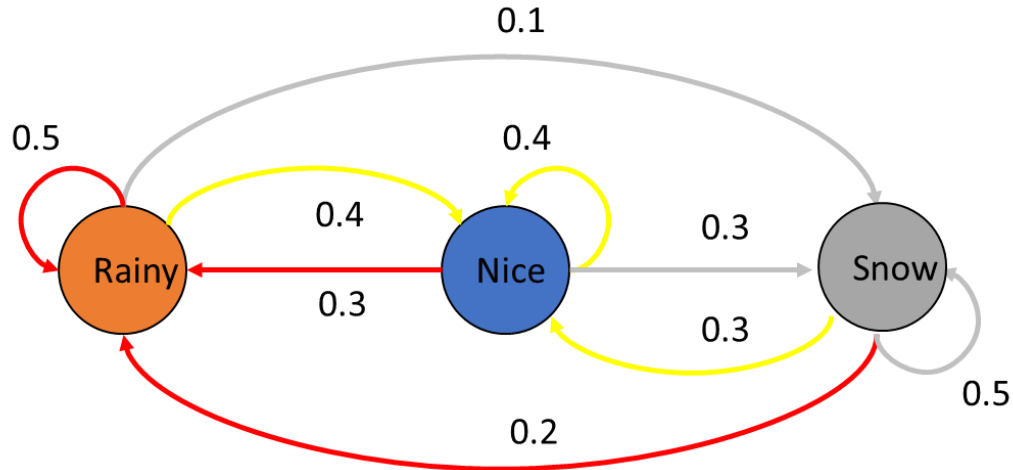
The entries in the first row of the matrix P represent the probabilities for the various kinds of weather following a rainy day. Similarly, the entries in the second and third rows represent the probabilities for the various kinds of weather following nice and snowy day, respectively.

Considering given the chain is in state i today, it will be in state j two days from now given this probability by  $p_{ij}$

Markov property: The state of the system at time  $t+1$  only depends on the state of the system at time  $t$

$$P[X_{t+1} = x_{t+1} | X_t = x_t, X_{t-1} = x_{t-1}, \dots, X_1 = x_1, X_0 = x_0] =$$

$$P[X_{t+1} = x_{t+1} | X_t = x_t]$$



Weather forecasting

- Two days:  $\begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix} \begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix} = \begin{pmatrix} 0.39 & 0.39 & 0.22 \\ 0.33 & 0.37 & 0.30 \\ 0.29 & 0.35 & 0.36 \end{pmatrix}$

- Four days:

$$\begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix}^2 \begin{pmatrix} 0.5 & 0.4 & 0.1 \\ 0.3 & 0.4 & 0.3 \\ 0.2 & 0.3 & 0.5 \end{pmatrix}^2 = \begin{pmatrix} 0.3446 & 0.3734 & 0.2820 \\ 0.3378 & 0.3706 & 0.2916 \\ 0.3330 & 0.3686 & 0.2984 \end{pmatrix}$$

The graph 3 (15)

Assume a Gaussian distribution for observations,  $X_i, i=1, \dots, N$  with unknown mean,  $M$  and known variance 5. Suppose the prior for  $M$  is Gaussian with variance 10. How large a random sample must be taken (i.e., what is the minimum value for  $N$ ) to specify an interval having unit length of 1 such that the probability that  $M$  lies in this interval is 0.95?

Response:

Gaussian with Unknown mean and known variance

From lecture

- Likelihood with N trials,  $x = (x_1, \dots, x_N)$  with unknown mean  $M$  and known variance  $\sigma^2$ 
  - $f(x|m, \sigma^2) \propto \frac{1}{\sigma} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - m)^2\right)$
- Prior for  $M$  is  $N(\mu_0, \sigma_0^2)$
- Posterior is  $N(\mu_N \text{ and } \sigma_{post}^2)$  where
  - $\mu_N = \frac{\mu_0 \sigma^2 + N \bar{x} \sigma_0^2}{\sigma^2 + N \sigma_0^2}$
  - $\sigma_{post}^2 = \frac{\sigma_0^2 \sigma^2}{\sigma^2 + N \sigma_0^2}$

The variance of the Gaussian distribution is 5 ( $\sigma$ ) and the variance for the prior is 10 ( $\sigma_0^2$ )

Since the posterior is a function of N, we need to find an N that makes this variance (here is 5) that the probability M is centered at the posterior mean  $\mu_N = 0.95$

$$\sigma_{post}^2 = \frac{\sigma_0^2 \sigma^2}{\sigma^2 + N \sigma_0^2} \rightarrow \sigma_{post}^2 = \frac{10 \cdot 5}{5 + N \cdot 10} = \frac{50}{5 + 10N}$$

$$CDF = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{x - \mu}{\sigma \sqrt{2}} \right) \right]$$

$$0.975 = \frac{1}{2} \left[ 1 + \operatorname{erf} \left( \frac{0.5}{\sqrt{\frac{50}{5 + 10N}} \sqrt{2}} \right) \right]$$

$$1.3859 * \left( \sqrt{\frac{50}{5 + 10N}} \sqrt{2} \right) = 0.5$$

$$1.9207 * \left( \frac{100}{5 + 10N} \right) = 0.25, N = 77$$

**4 (15)**

You have started an online business selling books that are of interest to your customers. A publisher has just given you a large book with photos from famous 20<sup>th</sup> century photographers. You think this book will appeal to people who have bought art books, history books and coffee table books. In an initial offering of the new book you collect data on purchases of the new book and combine these data with data from the past purchases (see ArtHistBooks.csv).

Use Bayesian analysis to give the posterior probabilities for purchases of art books, history books and coffee table books, as well as, the separate probabilities for purchases of new books given each possible combination of prior purchases of art books, history books and coffee table books. Do this by first using beta priors with values of the hyperparameters that represent lack of prior information. Then compute these probabilities again with beta priors that show strong weighting for low likelihood of a book purchase. Compare your results.

**Response: see the notebook**

**5 (15)**

The data set CHDdata.csv contains cases of coronary heart disease(CHD) and variables associated with the patient's condition: systolic blood pressure, yearly tobacco use (in kg), how density lipoprotein (ldl), adiposity, family history (0 or 1), type A personality score (typea), obesity (body mass index), alcohol use, and the diagnosis of CHD (0 or 1). Perform a Bayesian analysis of these data that finds the posterior marginal probability distributions for the means for the data of patients with and without CHD. You should first standard scale (subtract the mean and divide by the standard deviation) all the numeric variables (remove family history and do not scale CHD). Then separate the data into two sets, one for patients with CHD and one for patients without CHD.

Your priors for both groups should assume means of 0 for all variables and a correlation of 0 between all pairs of variables. You should assume all variances for the variables are 1. Use a prior alpha equal to one plus the number of predictor

variables. Compute and compare the Bayesian estimates for the posterior means for each group.

For 5 extra credit points, compute the probability of observing a point at least as extreme as the posterior mean of patients without coronary heart disease under the posterior distribution for the patients with coronary heart disease. Then compute the probability of observing a point at least as extreme as the posterior mean of patients with coronary heart disease under the posterior distribution for the patients without coronary heart disease

**Response: see the notebook**

## 6 (10)

For each of the following types of distributions, state the support type (single or multivariable and discrete or continuous), the formula for the PMP or PDF, the parameters, the support, the mean, and some typical uses of the distribution. You may use whatever source(s) you want, including for example Wikipedia.

- (a) Bernoulli Distribution
- (b) Binomial Distribution
- (c) Poisson Distribution
- (d) Uniform Distribution
- (e) Beta Distribution
- (f) Gamma Distribution
- (g) Gaussian Distribution
- (h) t Distribution
- (i) Cauchy Distribution
- (j) Multinomial Distribution
- (k) Dirichlet Distribution
- (l) Multivariate Gaussian Distribution
- (m) Multivariate t Distribution
- (n) Wishart Distribution

## Distribution, PMP, PDF, support type

a	Bernoulli Distribution	<p>Discrete probability distribution</p> <p>Support <math>k \in \{0, 1\}</math></p> <p>PMF <math>\begin{cases} q = 1 - p &amp; \text{if } k = 0 \\ p &amp; \text{if } k = 1 \end{cases}</math></p> <p>Mean = <math>p</math></p> <p>Bernoulli is the discrete probability distribution of a random variable which takes the value 1 with probability <math>p</math> and the value 0 with probability <math>q=1-p</math>.</p> <p>A Bernoulli distribution can be thought of as a model for the set of possible outcomes of any single experiment that asks a yes/true/one with probability <math>p</math>, and failure/no/false/zero with probability <math>q</math> such as coin toss problem.</p>
b	Binomial Distribution	<p>Binomial distribution is the discrete probability distribution.</p> <p>Support <math>k \in \{0, 1, \dots, n\}</math> – number of successes</p> <p>PMF <math>\binom{n}{k} p^k q^{n-k}</math></p> <p>Mean <math>np</math></p> <p>The binomial distribution is frequently used to model the number of successes in a sample of size <math>n</math> drawn with replacement from a population of size <math>N</math>.</p>
c	Poisson Distribution	<p>The Poisson distribution is a discrete probability distribution</p> <p>Support <math>k \in N_0</math>, Natural numbers starting from 0</p> <p>PMF <math>= \frac{\lambda^k e^{-\lambda}}{k!}</math></p> <p>Mean = <math>\lambda</math></p> <p>For instance: expresses the probability of a given number of events occurring in a fixed interval of time or space if these events occur with a known constant mean rate and independently of the time since the last event.</p>
d	Uniform Distribution	<p>Uniform distribution is the continuous uniform distribution or rectangular distribution</p> <p>Support <math>x \in [a, b]</math></p>

		<p>PDF <math>\begin{cases} \frac{1}{b-a} &amp; \text{for } x \in [a, b] \\ 0 &amp; \text{otherwise} \end{cases}</math></p> <p>Mean <math>\frac{1}{2}(a+b)</math></p> <p>For instance, density function, uniform probability density function</p>
e	Beta Distribution	<p>The beta distribution is a family of continuous probability distributions defined on the interval <math>[0, 1]</math> parameterized by two positive shape parameters, denoted by <math>\alpha</math> and <math>\beta</math>, that appear as exponents of the random variable and control the shape of the distribution.</p> <p>Support <math>x \in [0, 1]</math> or <math>x \in (0, 1)</math></p> <p>PDF <math>\frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}</math> Where <math>B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha+\beta)}</math> and <math>\Gamma</math> is the Gamma function</p> <p>Mean:</p> $E[X] = \frac{\alpha}{\alpha + \beta}$ $E[\ln X] = \psi(\alpha) - \psi(\alpha + \beta)$ $E[X \ln X] = \frac{\alpha}{\alpha + \beta} [\psi(\alpha + 1) - \psi(\alpha + \beta + 1)]$
f	Gamma Distribution	<p>Gamma distribution is a two-parameter family of continuous probability distributions.</p> <p>Support <math>x \in (0, \infty)</math></p> <p>PDF:</p> $f(x) = \frac{1}{\Gamma(k)\theta^k} x^{k-1} e^{-\frac{x}{\theta}} \quad (k > 0, \theta > 0)$ $f(x) = \frac{\beta^\alpha}{\Gamma(\alpha)} x^{\alpha-1} e^{-\beta x} \quad (\alpha > 0, \theta > 0)$ <p>Mean: <math>k\theta</math> (<math>k &gt; 0, \theta &gt; 0</math>)</p> $\frac{\alpha}{\beta} \quad (\alpha > 0, \theta > 0)$ <p>For instance: in life testing, the waiting time until death is a random variable that is frequently modeled with a gamma distribution.</p>
g	Gaussian Distribution	<p>Gaussian distribution is a type of continuous probability distribution</p> <p>Support <math>x \in R</math></p> $\text{PDF} = \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$ <p>Mean <math>\mu</math></p>
h	t Distribution	<p>t-distribution is a member of family of continuous probability distribution</p> <p>support <math>x \in (-\infty, \infty)</math></p>



		<p>PDF <math>\frac{\Gamma(\frac{\nu+1}{2})}{\sqrt{\nu\pi}\Gamma(\frac{\nu}{2})} (1 + \frac{x^2}{\nu})^{-\frac{\nu+1}{2}}</math></p> <p>Mean 0 for <math>\nu &gt; 1</math></p>
i	Cauchy Distribution	<p>The Cauchy distribution is a continuous probability distribution.</p> <p>Support <math>x \in (-\infty, \infty)</math></p> <p>PDF = <math>\frac{1}{\pi\gamma \left[1 + \left(\frac{x-x_0}{\gamma}\right)^2\right]}</math></p> <p>Mean undefined</p>
j	Multinomial Distribution	<p>In a probability theory, the multinomial distribution is a generalization of the binomial distribution. For example, it models the probability of counts for each side of a k-sided die rolled n times. For n independent trials each of which leads to a success for exactly one of k categories, with each category having a given fixed success probability, the multinomial distribution gives the probability of any particular combination of numbers of successes for the various categories.</p> <p>Support <math>x_i \in \{0, \dots, n\}, i \in \{1, \dots, k\}</math></p> $\sum x_i = n$ <p>PMF <math>\frac{n!}{x_1! \dots x_k!} p_1^{x_1} \dots p_k^{x_k}</math></p> <p>Mean <math>E(X_i) = np_i</math></p>
k	Dirichlet Distribution	<p>The Dirichlet distribution is a family of continuous multivariate probability distributions parametrized by a vector <math>\alpha</math> of positive reals.</p> <p>Support <math>x_1, \dots, x_K</math> where <math>x_i \in (0, 1)</math> and <math>\sum_{i=1}^K x_i = 1</math></p> <p>PDF:</p> $\frac{1}{B(\alpha)} \prod_{i=1}^K x_i^{\alpha_i - 1}$ <p>Where <math>B(\alpha) = \frac{\prod_{i=1}^K \Gamma(\alpha_i)}{\Gamma(\sum_{i=1}^K \alpha_i)}</math></p> <p>Where <math>\alpha = (\alpha_1, \dots, \alpha_K)</math></p> <p>Mean</p> $E[X_i] = \frac{\alpha_i}{\sum_{k=1}^K \alpha_k}$ $E[\ln X_i] = \psi(\alpha_i) - \psi(\sum_{k=1}^K \alpha_k)$

l	Multivariate Gaussian Distribution	<p>In probability theory, the multivariate normal distribution is a generalization of the one-dimensional normal distribution to higher dimensions.</p> <p>Support <math>x \in \mu + \text{span}(\Sigma) \subseteq R^k</math></p> <p>PDF <math>(2\pi)^{-\frac{k}{2}} \det(\Sigma)^{-\frac{1}{2}} e^{-\frac{1}{2}(x-\mu)^T \Sigma^{-1}(x-\mu)}</math> exists only when <math>\Sigma</math> is positive-definite</p> <p>Mean <math>\mu</math></p>
m	Multivariate t Distribution	<p>In statistics, the multivariate t-distribution is a multivariate probability distribution.</p> <p>Support <math>x \in R^p</math></p> <p>PDF <math>\frac{\Gamma[(\nu+p)/2]}{\Gamma(\nu/2) \nu^{p/2} \pi^{p/2}  \Sigma ^{1/2}} \left[ 1 + \frac{1}{\nu} (x - \mu)^T \Sigma^{-1} (x - \mu) \right]^{-(\nu+p)/2}</math></p> <p>Mean <math>\mu</math> if <math>\nu &gt; 2</math>, else undefined.</p>
n	Wishart Distribution	<p>In statistics, the Wishart distribution is a generalization to multiple dimensions of the gamma distribution.</p> <p>Support <math>X(p \times p)</math></p> <p>PDF <math>f_x(x) = \frac{ x ^{(n-p-1)/2} e^{-\text{tr}(V^{-1}x)/2}}{2^{\frac{np}{2}}  V ^{n/2} \Gamma_p(\frac{n}{2})}</math></p> <p>Mean <math>E[X] = nV</math></p>

## 7 (10)

Using the Python Notebook <http://www.kaggle.com/billbasener/pt2-probabilities-likelihoods-and-bayes-theorem>, complete the challenge questions from Section 6: Modify the code from Section 5 to add the ability to use the `posterior_from_conjugate_prior` function to output the posterior probability parameters given parameters and for a Gaussian Likelihood with known variance  $\sigma^2$ , and use your modified function to create the Prior, Likelihood, Posterior plots as in Section 5 of the notebook.

See the notebook in details

Regarding normal prior and normal likelihood

A normal prior is conjugate to a normal likelihood with known  $\sigma$

- Data:  $x_1, x_2, \dots, x_n$
- Normal likelihood  $x_1, x_2, \dots, x_n \sim N(\theta, \sigma^2)$
- Normal prior  $\theta \sim N(\mu_{prior}, \sigma_{prior}^2)$
- Normal Posterior  $\theta \sim N(\mu_{post}, \sigma_{post}^2)$