

Lecture 11A

MTH6102: Bayesian Statistical Methods

Eftychia Solea

Queen Mary University of London

2023

Today's agenda

Today's lecture

- Learn how to use the law of total probability to compute posterior predictive probabilities.

Review: Predictive probabilities

- Posterior predictive probability describes how likely are different outcomes of a future experiment.
- We have observed data (result of the experiment) $y \sim p(y \mid \theta)$, dependent on parameters θ .
- Then we update our prior distribution for θ , $p(\theta)$, to the posterior distribution $p(\theta \mid y)$.

Posterior predictive probabilities

- Suppose we plan to perform the experiment again to observe new data x
- We want to compute the posterior predictive distribution $p(x | y)$ of x given the observed data y .
- Posterior predictive probabilities are used to predict future data x when the experiment is performed again, and they are computed after observing data y and updating prior to posterior.

Predictive distributions: discrete prior, discrete data

- Discrete observed data: $y \sim p(y \mid \theta)$, with θ unknown
- Discrete likelihood: $p(y \mid \theta)$.
- Discrete hypothesis θ with values $\theta_1, \theta_2, \dots, \theta_K$.
- Prior pmf $p(\theta_i)$ of θ , $p(\theta_i) = p(\theta = \theta_i)$, $i = 1, \dots, K$.
- Posterior pmf $p(\theta_i \mid y) = \frac{p(y|\theta_i)p(\theta_i)}{p(y)}$, $i = 1, \dots, K$.

Hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$p(\theta)$	$p(y \theta)$	$p(y \theta)p(\theta)$	$p(\theta y)$
θ_1	$p(\theta_1)$	$p(y \theta_1)$	$p(y \theta_1) p(\theta_1)$	$p(\theta_1 y)$
θ_2	$p(\theta_2)$	$p(y \theta_2)$	$p(y \theta_2) p(\theta_2)$	$p(\theta_2 y)$
\vdots	\vdots	\vdots	\vdots	\vdots
θ_K	$p(\theta_K)$	$p(y \theta_K)$	$p(y \theta_K) p(\theta_K)$	$p(\theta_K y)$
Total	1	NOT SUM TO 1	$p(y)$	1

Predictive distributions: discrete prior, discrete data

- By the Law of total probability,

$$p(y) = \sum_{i=1}^K p(y|\theta_i)p(\theta_i),$$

is called the **prior predictive probability**.

- **Prior predictive probabilities.** Assign a probability to an outcome of the experiment. They are computed **before we observe any data**.

Predictive distributions: discrete prior, discrete data

- Let x : future data from the same experiment. We assume that x and y are independent given θ . $p(x|y, \theta) = p(x|\theta)$
- By, the **law of total probability**, the **posterior predictive probability** of x given the observed data y is

$$p(x|y) = \sum_{i=1}^K p(x|\theta_i) p(\theta_i|y).$$

$$p(x|y) = \sum_{i=1}^K p(x|\theta_i) p(\theta_i|y)$$

Board example: Three type of coins

There are three type of coins in the drawer with probabilities 0.5, 0.6 and 0.9 of heads, respectively. Each coin is equally likely

Data: Pick one and toss 5 times. You get 1 head out of 5 tosses.

- (a) Compute the posterior probabilities for the type of coin
- (b) Compute the posterior predictive distributions of observing heads in a future toss.
- (c) Compute the posterior predictive distributions of observing 2 heads in 5 future coin tosses.

- new data $x=1$. The posterior predictive probability

$$p(x=1|y=1) = \sum_{\theta \in \{0.5, 0.6, 0.9\}} p(x=1|\theta) p(\theta|y=1)$$

$$= p(x=1|\theta=0.5) p(\theta=0.5|y=1) + p(x=1|\theta=0.6) p(\theta=0.6|y=1) + p(x=1|\theta=0.9) p(\theta=0.9|y=1)$$

$$= 0.5 (0.669) + 0.6 (0.329) + 0.9 (0.002)$$

$$= 0.46634 \approx 0.5$$

• $n=1, x=1 \sim \text{binomial}(1, \theta)$
 $p(x=1|\theta) = \binom{1}{1} \theta^1 (1-\theta)^{1-1} = \theta$

- new data: $x=2$ heads out of 5 tosses.

$$p(x=2|y=1) = \sum_{\theta \in \{0.5, 0.6, 0.9\}} p(x=2|\theta) p(\theta|y=1)$$

$$= p(x=2|\theta=0.5) p(\theta=0.5|y=1) + p(x=2|\theta=0.6) p(\theta=0.6|y=1)$$

$$+ p(x=2|\theta=0.9) p(\theta=0.9|y=1)$$

binomial prob.

$$p(x=2|\theta) = \binom{5}{2} \theta^2 (1-\theta)^3 \quad \theta \in \{0.5, 0.6, 0.9\}$$

$$= \binom{5}{2} 0.5^2 (0.5)^3 (0.669) + \binom{5}{2} 0.6^2 (0.4)^3 (0.329)$$

$$+ \binom{5}{2} 0.9^2 (0.1)^3 (0.002) = 0.28$$

Board example: Three type of coins

- Bayesian updating table

Hypothesis	prior	likelihood	Bayes numerator	posterior
θ	$p(\theta)$	$p(y \theta) \sim \text{binomial}(5, \theta)$	$p(y \theta)p(\theta)$	$p(\theta y)$
$\theta_1 = 0.5$	$p(\theta_1) = 1/3$	$p(y = 1 \theta_1) = 0.15625$	$p(y = 1 \theta_1) p(\theta_1) = 0.0521$	$p(\theta_1 y = 1) = 0.669$
$\theta_2 = 0.6$	$p(\theta_2) = 1/3$	$p(y = 1 \theta_2) = 0.0768$	$p(y = 1 \theta_2) p(\theta_2) = 0.0256$	$p(\theta_2 y = 1) = 0.329$
$\theta_3 = 0.9$	$p(\theta_3) = 1/3$	$p(y = 1 \theta_3) = 0.00045$	$p(y = 1 \theta_3) p(\theta_3) = 0.00015$	$p(\theta_3 y = 1) = 0.00193$
Total	1	NOT SUM TO 1	$p(y = 1) = 0.07785$	1

- Prior predictive probability: $p(y = 1) = p(y = 1|\theta_1)p(\theta_1) + p(y = 1|\theta_2)p(\theta_2) + p(y = 1|\theta_3)p(\theta_3) = 0.07785$

Board example: Three type of coins

- Does the order of the 1 head and 4 tails affect the posterior distribution of the coin type?

(a) Yes

☒ (b) No.

- Does the order of the 1 head and 4 tails affect the posterior predictive distribution of the next flip?

(a) Yes

☒ (b) No.

Board example

- Suppose that y is the number of expensive goods in a shop over 24 days. So $y \sim \text{Poisson}(24\theta)$ where $\theta = 1/2$, $\theta = 1/4$ or $\theta = 1/8$.
- Suppose the prior pmf is

$$p(\theta = 1/2) = p(1/2) = 0.2, \quad p(\theta = 1/4) = p(1/4) = 0.5, \\ p(\theta = 1/8) = p(1/8) = 0.3.$$

- We observe $y = 10$ expensive goods were sold in the last 24 days.
 - 1 Compute the posterior pmf for θ .
 - 2 Compute the posterior predictive distribution that $x = 10$ number of goods will be sold in the next 24 days.

The likelihood in this case is

$$p(y=10|\theta) = \frac{(24\theta)^{10} e^{-24\theta}}{10!}, \quad \theta \in \{0.5, 0.25, \frac{1}{8}\}$$

Predictive distributions: continuous prior, discrete data

- Continuous parameter θ in the range $[a, b]$.
- Prior: $p(\theta)$, $\theta \in [a, b]$.
- Discrete data, y . Likelihood $p(y|\theta)$.
- By, the **law of total probability**, the **prior predictive probability of y** is

$$p(\text{data}) = p(y) = \int_a^b \underbrace{p(y|\theta) p(\theta)} d\theta,$$

where the integral is computed over the entire range of θ .

- **Note:** $p(y)$ is a probability mass function, i.e., $p(y) = P(Y = y)$

Similar

$$p(y) = \sum_{i=1}^K p(y|\theta_i) p(\theta_i)$$

Predictive distributions: continuous prior, discrete data

- Posterior: $p(\theta|y) = \frac{p(\theta) \times p(y|\theta)}{p(y)}$
- x : future data of the same experiment. We assume that x and y are independent given θ
- By, the **law of total probability**, the **posterior predictive probability of x** (given y) is

$$p(x|y) = \int_a^b \underbrace{p(x|\theta) p(\theta|y)} d\theta.$$

$$p(x|y) = \int p(x|\theta) p(\theta|y) d\theta$$

Example

We have a coin with unknown probability θ of heads, $\theta \in [0, 1]$.
Prior: $p(\theta) = 2\theta$, $\theta \in [0, 1]$.

- Find the prior predictive probability of throwing heads on the first toss.
- Suppose the first flip was heads. Find the posterior predictive probabilities of both heads and tails on the second flip.

Solution

• Let y be the result of the first toss.

$$\begin{aligned} p(y=1) &= \int_0^1 p(y=1|\theta) p(\theta) d\theta \\ &= \int_0^1 \theta \cdot (2\theta) d\theta = \frac{2}{3} \end{aligned}$$

• Data $y=1$ (first flip was heads). First, we need to compute the posterior pdf, $p(\theta|y=1)$.

By Bayes' Theorem,

$$p(\theta|y=1) = \frac{p(\theta) \times p(y=1|\theta)}{p(y=1)} = \frac{(2\theta) \cdot \theta}{2/3} = \underline{3\theta^2}$$

Let x be the result of the second flip. Then,

$$\begin{aligned} p(x=1|y=1) &= \int_0^1 p(x=1|\theta) p(\theta|y=1) d\theta \\ &= \int_0^1 \theta (3\theta^2) d\theta = \frac{3}{4} \end{aligned}$$

Example: beta prior/ binomial data

- Data, $k \sim \text{binomial}(n, q)$
 - Prior, $q \sim \text{beta}(\alpha, \beta)$.
-
- Find the posterior predictive probability to observe success on the next Bernoulli trial.
 - Find the posterior predictive probability to observe new x successes on the next m Bernoulli trials.

Solution

First, the posterior pdf of z given the data x is

$$p(z|x) \sim \text{beta}(\alpha+x, \beta+n-x)$$

The posterior, predictive distribution of x given x is

$$p(x|x) = \int_0^1 p(x|z) p(z|x) dz$$

$$\text{Now, } p(x|z) = \binom{m}{x} z^x (1-z)^{m-x}$$

$$p(z|x) = \frac{z^{\alpha+x-1} (1-z)^{\beta+n-x-1}}{\text{Beta}(\alpha+x, \beta+n-x)}$$

$$\text{Thus, } p(x|x) = \int_0^1 \binom{m}{x} z^x (1-z)^{m-x} \frac{z^{\alpha+x-1} (1-z)^{\beta+n-x-1}}{\text{Beta}(\alpha+x, \beta+n-x)} dz$$

$$= \binom{m}{x} \frac{1}{\text{Beta}(\alpha+x, \beta+n-x)} \int_0^1 z^{x+\alpha+x-1} (1-z)^{m-x+\beta+n-x-1} dz$$

$X \sim \text{beta}(a, b)$ the pdf is

$$f_X(x) = \frac{x^{a-1} (1-x)^{b-1}}{\text{Beta}(a, b)} \quad x \in [0, 1] \\ a > 0, b > 0$$

$$\int_0^1 f_X(x) dx = 1$$

$$\text{So } \int_0^1 \frac{x^{a-1} (1-x)^{b-1}}{\text{Beta}(a, b)} dx = 1$$

$$\Rightarrow \frac{1}{\text{Beta}(a, b)} \int_0^1 x^{a-1} (1-x)^{b-1} dx = 1$$

$$\Rightarrow \int_0^1 x^{a-1} (1-x)^{b-1} dx = \text{Beta}(a, b)$$

Use ④ with $a+x+k$ instead of a and $\theta+m-x+n-k$ instead of θ , to find

$$p(x/k) = \binom{m}{x} \frac{\text{Beta}(x+a+x, \theta+m-x+n-k)}{\text{Beta}(a+k, \theta+n-x)}.$$

Board example

Data: 10 patients have 6 successes. $\theta \sim \text{beta}(5, 5)$

- Find the posterior distribution of θ .
- Find the posterior predictive probability of success with the next patient.

Posterior predictive distribution: continuous prior, continuous data

- Continuous parameter θ in the range $[a, b]$.
- Prior pdf: $p(\theta)$, $\theta \in [a, b]$.
- Continuous data, y . Likelihood $p(y|\theta)$.
- The prior predictive pdf of y is

$$p(y) = \int_a^b p(y|\theta) p(\theta) d\theta,$$

where the integral is computed over the entire range of θ .

- **Note:** $p(y)$ is a pdf.

Posterior predictive distribution: continuous prior, continuous data

- Posterior pdf: $p(\theta|y)$
- x : future data of the same experiment.
- The posterior predictive distribution of x is

$$p(x|y) = \int_a^b \underbrace{p(x|y, \theta) p(\theta|y)} d\theta.$$

- As usual, we usually assume x and y are conditionally independent given θ . That is, $p(x|y, \theta) = p(x|\theta)$.
- In this case,

$$\underbrace{p(x|y) = \int_a^b p(x|\theta) p(\theta|y) d\theta}_{\text{red circled equation}}$$

Posterior predictive distribution

The posterior predictive distribution for x given the observed data y is

$$p(x \mid y) = \int p(x \mid \theta) p(\theta \mid y) d\theta$$

- This is the probability distribution for unobserved or future data x .
- This distribution includes two types of uncertainty:
 - the uncertainty remaining about θ after we have seen y ;
 - the random variation in x .

Board example: Exponential data/Gamma prior

- The time until failure for a type of light bulb is exponentially distributed with parameter $\theta > 0$, where θ is unknown.
 - We observe n bulbs, with failure times t_1, \dots, t_n . *$t_i \sim \exp(\theta)$*
 - We assume a $\text{Gamma}(\alpha, \beta)$ prior distribution for θ , where $\alpha > 0$ and $\beta > 0$ are known.
- 1 Determine the predictive posterior distribution for future data x

Solution

Observed data $t = (t_1, \dots, t_n)$, where each $t_i \sim \text{Exp}(\theta)$, $\theta > 0$.

Since Gamma(a, θ) is conjugate to the exponential likelihood, the posterior of θ given the data t , is

$$p(\theta|t) \sim \text{Gamma}(\underbrace{a+n}_{\tilde{a}}, \underbrace{\theta+S}_{\tilde{\theta}}), \text{ where } S = \sum_{i=1}^n t_i$$

Future data $x \sim \text{Exp}(\theta)$

$$\begin{aligned}\tilde{a} &= a+n \\ \tilde{\theta} &= \theta+S\end{aligned}$$

The posterior predictive distribution of x given t

is

$$p(x|t) = \int_0^{\infty} p(x|\theta) p(\theta|t) d\theta$$

$$= \int_0^{\infty} \theta e^{-\theta x} \frac{(\tilde{\theta})^{\tilde{a}} \theta^{\tilde{a}-1} e^{-\tilde{\theta}\theta}}{\Gamma(\tilde{a})} d\theta$$

$$= \frac{(\tilde{\theta})^{\tilde{a}}}{\Gamma(\tilde{a})} \int_0^{\infty} \underbrace{\theta^{(\tilde{a}+1)-1} \exp(-(x+\tilde{\theta})\theta)}_{\text{Gamma PDF}} d\theta$$

$X \sim \text{Gamma}(a, \theta)$ then

$$\int_0^{\infty} f_X(x) dx = 1 \Leftrightarrow \int_0^{\infty} \frac{\theta^a}{\Gamma(a)} x^{a-1} e^{-\theta x} dx = 1$$

$$\Leftrightarrow \frac{\theta^a}{\Gamma(a)} \int_0^{\infty} x^{a-1} e^{-\theta x} dx = 1 \quad \Leftrightarrow$$

$$\Rightarrow \int_0^{\infty} x^{a-1} e^{-\theta x} dx = \frac{\Gamma(a)}{\theta^a} \quad (*)$$

Use $(*)$ but substitute in $\tilde{a}+1$ instead of a and $x+\tilde{e}$ instead of θ to get

$$p(x|t) = \frac{(\tilde{e})^{\tilde{a}}}{\Gamma(\tilde{a})} \frac{\Gamma(\tilde{a}+1)}{(x+\tilde{e})^{\tilde{a}+1}}$$

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

$$\begin{aligned} p(x|t) &= \frac{(\tilde{e})^{\tilde{a}} \tilde{a} \cancel{\Gamma(\tilde{a})}}{\cancel{\Gamma(\tilde{a})} (x+\tilde{e})^{\tilde{a}+1}} \\ &= \frac{(\tilde{e})^{\tilde{a}} \tilde{a}}{(x+\tilde{e})^{\tilde{a}+1}} \end{aligned}$$

Finding the posterior predictive distribution

$$p(x \mid y) = \int p(x \mid \theta) p(\theta \mid y) d\theta$$

- In conjugate examples, one can usually derive $p(x \mid y)$.
- It is generally easier to find the mean and variance of $p(x \mid y)$ than deriving the full distribution.

Conditional mean and variance in general

- Suppose that X and W are general random variables.
- Then

$$E(X) = E(E(X | W)) \quad \text{law of iterated expectation}$$

and

$$\text{Var}(X) = \text{Var}(E(X | W)) + E(\text{Var}(X | W)) \quad \text{law of total variance}$$

- In Bayesian inference, we replace W with parameters and X with the new data we would like to predict.

$$\underline{\mathbb{E}(X|W) = g(W)} \Rightarrow \text{random variable}$$

$$\text{if } W=w, w \in \mathbb{R}$$

$$\underbrace{\mathbb{E}(X|W=w) = g(w)}_{\text{non-random}}$$

$$\bullet \mathbb{E}(g(W)) = \mathbb{E}(\mathbb{E}(X|W)) = \mathbb{E}(X)$$

Mean and variance of posterior predictive distribution

- For new data x and parameter(s) θ

$$E(x) = E(E(x \mid \theta))$$

$$Var(x) = Var(E(x \mid \theta)) + E(Var(x \mid \theta))$$

$$E(x \mid y) = E(E(x \mid \theta, y))$$

Mean and variance of posterior predictive distribution

- Add conditioning on observed data y , since we want *posterior* predictions

$$E(x | y) = E(E(x | \theta, y)) \quad \text{law of iterated expectation}$$

$$\Rightarrow Var(x | y) = Var(E(x | \theta, y)) + E(Var(x | \theta, y)) \quad \text{law of total variance}$$

- These are the posterior predictive mean and posterior predictive variance of x , respectively.

Because x and y are independent given θ

$$\underline{f_{x|y,\theta}}(x|y,\theta) = f_{x|\theta}(x|\theta)$$

$$E(x|y,\theta) = E(x|\theta)$$

$$\text{var}(x|y,\theta) = \text{var}(x|\theta)$$

Example: beta prior, binomial data

- Data, $k \sim \text{binomial}(n, q)$
- Prior, $q \sim \text{beta}(\alpha, \beta)$.
- New data, $x \sim \text{binomial}(m, q)$, m is known.

(1) Find the posterior predictive mean and variance of x

Solution

By the law of iterated expectation,

$$\begin{aligned} \mathbb{E}(x|x) &= \mathbb{E}(\mathbb{E}(x|z, x)) \\ &= \mathbb{E}(\mathbb{E}(x|z)) \end{aligned}$$

Now given z , $x \sim \text{binomial}(m, z)$

$$\text{So } \mathbb{E}(x|z) = m \cdot z$$

$$\underbrace{\mathbb{E}(x|x)}_{\text{predictive mean}} = \mathbb{E}(m \cdot \underbrace{z}_{\text{prior expectation}}) = m \cdot \underbrace{\frac{a}{a+b}}_p = m \cdot p$$

$$\text{Var}(x|x) = \mathbb{E}(\underbrace{\text{Var}(x|x, z)}_{\text{conditional variance}}) + \underbrace{\text{Var}(\mathbb{E}(x|x, z))}_{\text{variance of conditional expectation}}$$

$$\mathbb{E}(\text{Var}(x|x, z)) = \mathbb{E}(\text{Var}(x|z))$$

$$\text{Var}(x|z) = m z (1-z)$$

$$\begin{aligned} \text{So } \mathbb{E}(\text{Var}(x|z)) &= \mathbb{E}(m z (1-z)) = m \mathbb{E}(z - z^2) \\ &= m [\mathbb{E}(z) - \mathbb{E}(z^2)] \end{aligned}$$

$$\text{Var}(z) = \mathbb{E}(z^2) - [\mathbb{E}(z)]^2 \Rightarrow \boxed{\mathbb{E}(z^2) = \text{Var}(z) + [\mathbb{E}(z)]^2}$$

We have

$$\begin{aligned} E(\text{Var}(x|x, z)) &= m E(z) - m [\text{Var}(z) + E(z)^2] \\ &= mp - m(v + p^2) \quad (1) \end{aligned}$$

where v is the prior variance, p is the prior mean.

$$\begin{aligned} \bullet \text{Var}(E(x|x, z)) &= \text{Var}(E(x|z)) = \text{Var}(mz) \\ &= m^2 \text{Var}(z) \\ &= m^2 v \quad (2) \end{aligned}$$

where $v = \text{Var}(z)$

Using simulation (Monte Carlo)

- Suppose we know the posterior distribution $p(\theta \mid y)$, or we have a sample from it.
- Then it is easy to use simulation to generate a sample from the posterior predictive distribution of a new data-point x .
- Because we know the distribution of x for any given value of θ : it's the same as the distribution of the original data y .

Simulating the posterior predictive distribution

- Suppose that we have a sample from the posterior distribution

$$\theta_1, \theta_2, \dots, \theta_M$$

- We can simulate the posterior predictive distribution $p(x \mid y)$.
- We just generate

$$x_j \text{ from } p(x \mid \theta_j, y) = p(x \mid \theta_j), \quad j = 1, 2, \dots, M$$

- Then

$$x_1, x_2, \dots, x_M$$

is a sample from the posterior predictive distribution $p(x \mid y)$.

- (Since

$$(x_1, \theta_1), (x_2, \theta_2), \dots, (x_M, \theta_M)$$

is a sample from $p(x, \theta \mid y) = p(\theta \mid y) p(x \mid \theta, y)$).

Simulating the posterior predictive distribution

- When do we have a sample from $p(\theta \mid y)$?
- Almost always, because we use MCMC to make inferences about θ .
- Or in simpler conjugate cases, we can directly generate an independent sample from $p(\theta \mid y)$.
- The latter is an example of simple Monte Carlo.

Using the the posterior predictive sample

- Suppose we have generated a sample from the posterior predictive distribution x_1, x_2, \dots, x_M .
- We can summarize the sample for whatever interests us:
 - Posterior predictive mean, median, variance - just summarize sample x_1, x_2, \dots, x_M
 - Prediction intervals, e.g. with 95% probability, x will be in some interval- just take the 0.025 and 0.975 sample quantiles of the sample x_1, x_2, \dots, x_M .
 - Posterior predictive probability that $x = 0$ - just count what proportion of sample are 0.
 - Posterior predictive probability that $x > c$, for some c - count what proportion of sample are $> c$.