

Categorical and multinomial distributions

David J. H. Shih

Intended learning outcomes

- ▶ Recognize and apply categorical and multinomial distributions in statistical models
- ▶ Derive the multinomial distribution
- ▶ Derive posterior predictive distributions for statistical models

Categorical distribution

Definition

A random variable X with support $\mathcal{X} = \{1, 2, \dots, K\}$ has a categorical distribution if

$$P(X = k \mid \boldsymbol{\theta}) = \theta_k, \quad k \in \mathcal{X}.$$

where $\boldsymbol{\theta} = [\theta_k]$ and θ_k is the probability of success for outcome k
s.t. $\sum_k \theta_k = 1$.

Other representations are

$$P(X = x \mid \boldsymbol{\theta}) = \prod_k^K \theta_k^{I(x=k)}$$

$$p(X = x \mid \boldsymbol{\theta}) = \text{Categorical}(x \mid \boldsymbol{\theta})$$

$$X \sim \text{Categorical}(\boldsymbol{\theta})$$

$$X \sim \boldsymbol{\theta}.$$

Categorical distribution

- ▶ Generalization of Bernoulli distribution
- ▶ Representation for empirical probability mass function

Example

Let $X = \{1, 2, 3, 4\}$ represent the suit (Club, Diamond, Heart, Spade) of a playing card drawn from a cut deck.

$$X \sim \theta$$

Dirichlet distribution

Define a random vector $\mathbf{X} = [X_1, \dots, X_K]$ with $K \geq 2$ such that

$$X_k \in [0, 1] \quad \text{and} \quad \sum_k^K X_k = 1.$$

\mathbf{X} has a Dirichlet distribution with parameter $\boldsymbol{\alpha}$ if

$$f_{\mathbf{X}}(\mathbf{x}) = \frac{1}{B(\boldsymbol{\alpha})} \prod_k^K x_k^{\alpha_k - 1},$$

where $\alpha_k > 0$ and $B(\boldsymbol{\alpha})$ is the multivariate beta function:

$$B(\boldsymbol{\alpha}) = \frac{\prod_k^K \Gamma(\alpha_k)}{\Gamma\left(\sum_k^K \alpha_k\right)}.$$

Dirichlet-categorical model

Given data $X \in \{1, \dots, K\}$, define model

$$\begin{aligned} X &\sim \theta \\ \theta &\sim \text{Dirichlet}(\alpha), \end{aligned}$$

with parameter θ and hyperparameter α .

The posterior distribution of θ is

$$p(\theta | x) = \frac{p(x | \theta) p(\theta)}{p(x)} \quad (\text{Bayes' theorem}),$$

where we need to solve $p(x | \theta)$ $p(\theta)$ and $p(x)$.

We first solve the numerator:

$$\begin{aligned} p(x \mid \boldsymbol{\theta}) p(\boldsymbol{\theta}) &= \theta_x \left(B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha_k - 1} \right) \\ &= \left(\prod_k \theta_k^{I(x=k)} \right) B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha_k - 1} \\ &= B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha_k + I(x=k) - 1} \\ &= B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha'_k - 1}, \end{aligned}$$

where $\alpha'_k = \alpha_k + I(x = k)$.

Now, we solve the denominator:

$$\begin{aligned} p(x) &= \int_{\Theta} p(x | \boldsymbol{\theta}) p(\boldsymbol{\theta}) d\boldsymbol{\theta} \quad (\text{total law of prob.}) \\ &= \int_{\Theta} B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha'_k - 1} d\boldsymbol{\theta} \\ &= B(\boldsymbol{\alpha})^{-1} \int_{\Theta} \prod_k \theta_k^{\alpha'_k - 1} d\boldsymbol{\theta} \\ &= B(\boldsymbol{\alpha})^{-1} B(\boldsymbol{\alpha}') \quad (\text{Dirichlet distribution integrates to 1}) \end{aligned}$$

Finally, we solve the posterior by substituting in the numerator and the denominator.

$$\begin{aligned} p(\boldsymbol{\theta} | x) &= \frac{B(\boldsymbol{\alpha})^{-1} \prod_k \theta_k^{\alpha'_k - 1}}{B(\boldsymbol{\alpha})^{-1} B(\boldsymbol{\alpha}')} = \frac{\prod_k \theta_k^{\alpha'_k - 1}}{B(\boldsymbol{\alpha}')} \\ &= \text{Dirichlet}(\boldsymbol{\theta} | \boldsymbol{\alpha}'), \end{aligned}$$

where $\alpha'_k = \alpha_k + I(x = k)$.

Multiple categorical random variables

Given observed iid data $Z_i \in \{1, \dots, K\}$ for $i \in \{1, \dots, N\}$ and known parameter θ , suppose

$$Z_i \sim \theta.$$

In other words,

$$P(Z_i = k \mid \theta) = \theta_k$$

$$P(Z_i = z_i \mid \theta) = \prod_k^K \theta_k^{I(z_i=k)}$$

$$\begin{aligned} \prod_i P(Z_i = z_i \mid \theta) &= \prod_i^N \prod_k^K \theta_k^{I(z_i=k)} \\ &= \prod_k^K \theta_k^{\sum_i^N I(z_i=k)} \end{aligned}$$

For $k = 1 \dots K$, define $x_k = \sum_i^N I(z_i = k)$. Then,

$$\prod_i P(Z_i = z_i \mid \boldsymbol{\theta}) = \prod_k^K \theta_k^{x_k}$$

Define random vector $\mathbf{X} = [X_k]$, where $X_k = \sum_i^N I(Z_i = k)$.

$$P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) \propto \prod_k^K \theta_k^{x_k},$$

where each \mathbf{x} can arise from different number of occurrences of each realization $\mathbf{z} = [z_i]$.

Example where $K = 4$ and $N = 6$:

z	x
[1 1 1 1 1 1]	[6 0 0 0]
[2 2 2 2 2 2]	[0 6 0 0]
...	
[1 1 1 2 2 2]	[3 3 0 0]
[1 2 1 1 2 2]	[3 3 0 0]
...	[3 3 0 0]
...	
[1 2 2 1 3 4]	[2 2 1 1]
[2 1 2 1 3 4]	[2 2 1 1]
[3 1 2 1 4 2]	[2 2 1 1]
[4 1 2 1 3 2]	[2 2 1 1]
...	[2 2 1 1]
...	

There are $N!$ ways of drawing random vector $[Z_i]$, which consists of N random variables.

We cannot distinguish among all ways of arranging all the Z_i such that $Z_i = 1$, so we divide by $X_1!$.

Similarly for X_2, X_3, \dots

Therefore, each \mathbf{x} has $\frac{N!}{x_1! x_2! \dots x_K!}$ occurrences, which means that

$$P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) = \frac{N!}{\prod_k^K x_k!} \prod_k^K \theta_k^{x_k}.$$

Multinomial distribution

Definition

A random vector $\mathbf{X} = [X_1, \dots, X_K]$ with support

$$\mathcal{X} = \left\{ [x_1, \dots, x_K] : x_k \geq 0, \sum_k x_k = N \right\}$$

has a multinomial distribution if

$$\begin{aligned} P(\mathbf{X} = \mathbf{x} \mid \boldsymbol{\theta}) &= \frac{N!}{\prod_k x_k!} \prod_k \theta_k^{x_k} \\ &= \frac{\Gamma(N+1)}{\prod_k \Gamma(x_k + 1)} \prod_k \theta_k^{x_k}, \end{aligned}$$

where θ_k is probability of success for outcome k s.t. $\sum_k \theta_k = 1$.
Usually, each x_k is a non-negative integer.

Multinomial distribution

Given $\boldsymbol{\theta} = [\theta_k]$, we can also write

$$\begin{aligned} p(X = \mathbf{x} \mid \boldsymbol{\theta}) &= \text{Multinomial}(\mathbf{x} \mid N, \boldsymbol{\theta}) \\ X &\sim \text{Multinomial}(N, \boldsymbol{\theta}). \end{aligned}$$

Multinomial theorem

Theorem

Given positive integers N and K . Let \mathcal{A} be a set of vectors $\mathbf{x} = [x_1, \dots, x_K]$ s.t. each x_k is a nonnegative integer and $\sum_k x_k = N$. Then, for any real numbers p_1, \dots, p_K ,

$$\left(\sum_k p_k \right)^N = \sum_{\mathbf{x} \in \mathcal{A}} \frac{N!}{\prod_k x_k!} \prod_k p_k^{x_k}.$$

Applying the multinomial theorem with on $\boldsymbol{\theta} = [\theta_1, \dots, \theta_K]$, we get

$$\frac{N!}{\prod_k x_k!} \prod_k \theta_k^{x_k} = \left(\sum_k \theta_k \right)^N = 1^N,$$

since $\sum_k \theta_k = 1$. This proves that the multinomial distribution satisfies Kolmogorov's axiom 2.

Posterior predictive distribution

Given data $\mathbf{x} = [x_i]$, if a model with parameter $\theta \in \Theta$ has a posterior distribution $p(\theta | \mathbf{x})$, then the posterior prediction distribution for a new observation \tilde{x} is

$$p(\tilde{x} | \mathbf{x}) = \int_{\Theta} p(\tilde{x} | \theta) p(\theta | \mathbf{x}) d\theta,$$

which follows from the application of the total law of probability.

$p(\tilde{x} | \mathbf{x})$ may be used as likelihood function to evaluate the probability of held-out observations \tilde{x}_i , or as a sampling distribution to sample new observations \tilde{x}_i .

Posterior predictive distribution

In Bayesian modelling, we often sample from $p(\tilde{x} | \boldsymbol{x})$ and evaluate how the empirical distribution of posterior predictive sample $[\tilde{x}_i]$ agrees with the empirical distribution of \boldsymbol{x} .

Iterative update under categorical-Dirichlet model

Let $X = \{1, 2, 3, 4\}$ represent the suit (Club, Diamond, Heart, Spade) of a playing card drawn from a cut deck (i.e. $K = 4$).

$$X \sim \theta$$

$$\theta \sim \text{Dirichlet}(\alpha)$$

Draw a card, which has suit x_1 , and update $p(\theta | x_1)$. Put card back into the deck. Repeat for x_2, \dots, x_N . What is the probability distribution for the suit of the $N + 1$ card, X_{N+1} ?

$$p(\boldsymbol{\theta}) = \text{Dirichlet}(\boldsymbol{\theta} | [\alpha_k])$$

$$\begin{aligned} p(\boldsymbol{\theta} | x_1) &= p(x_1 | \theta) p(\theta) p(x_1)^{-1} \\ &= \text{Dirichlet}(\boldsymbol{\theta} | [\alpha_k + I(x_1 = k)]) \end{aligned}$$

$$\begin{aligned} p(\boldsymbol{\theta} | x_1, x_2) &= p(x_2 | \theta, x_1) p(\theta | x_1) p(x_2 | x_1)^{-1} \\ &= \text{Dirichlet}(\boldsymbol{\theta} | [\alpha_k + I(x_1 = k) + I(x_2 = k)]) \end{aligned}$$

...

$$\begin{aligned} p(\boldsymbol{\theta} | x_{1:N}) &= \frac{p(x_N | \theta, x_{1:(N-1)}) p(\theta | x_{1:(N-1)})}{p(x_N | x_{1:(N-1)})} \\ &= \text{Dirichlet}(\boldsymbol{\theta} | [\alpha_k + y_k]), \end{aligned}$$

where $y_k = \sum_i^N I(x_i = k)$.

$$\begin{aligned}
p(X_{N+1} = l \mid x_{1:N}) &= \int_{\Theta} p(X_{N+1} = l \mid \boldsymbol{\theta}) p(\boldsymbol{\theta} \mid x_{1:N}) d\boldsymbol{\theta} \\
&= \int_{\Theta} \theta_l p(\boldsymbol{\theta} \mid x_{1:N}) d\boldsymbol{\theta} \\
&= \int_0^1 \int_{\Theta_{-l}} \theta_l p(\boldsymbol{\theta} \mid x_{1:N}) d\boldsymbol{\theta}_{-l} d\theta_l \\
&= \int_0^1 \theta_l \int_{\Theta_{-l}} p(\boldsymbol{\theta} \mid x_{1:N}) d\boldsymbol{\theta}_{-l} d\theta_l \\
&= \int_0^1 \theta_l p(\theta_l \mid x_{1:N}) d\theta_l \quad (\text{total law of prob.}) \\
&= \mathsf{E}_{p(\theta_l \mid x_{1:N})}(\theta_l) \quad (\text{def'n of expectation})
\end{aligned}$$

Under Dirichlet (θ | α),

$$\mathsf{E}(\theta_l) = \frac{\alpha_l}{\sum_k \alpha_k}.$$

Therefore,

$$\begin{aligned} p(X_{N+1} = l \mid x_{1:N}) &= \mathsf{E}_{p(\theta_l \mid x_{1:N})}(\theta_l) \\ &= \mathsf{E}_{\text{Dirichlet}(\theta_l \mid [\alpha'_k])}(\theta_l) \\ &= \frac{\alpha_l^{(N)}}{\sum_k \alpha_k^{(N)}} \end{aligned}$$

where $\alpha_k^{(N)} = \alpha_k + \sum_i^N I(x_i = k)$.

Summary

"In God we trust; all others bring data." - W. Edwards Deming

Casella & Berger 2002, section 4.6, pages 177-182.

Intended learning outcomes

- ▶ Recognize and apply categorical and multinomial distributions in statistical models
- ▶ Derive the multinomial distribution
- ▶ Derive posterior predictive distributions for statistical models