

1 Probability distributions

1.1 Uniform distribution

1.2 Beta distribution

1.3 Bernoulli distribution

1.4 Binomial distribution

1.5 Beta-binomial distribution

1.6 Categorical distribution

1.7 Dirichlet distribution

1.8 Multinomial distribution

1.9 Pareto distribution

2 Bayesian parameter estimation

2.1 Beta-Bernoulli model

2.1.1 Summary

The model

$$X_i \sim \text{Ber}(\theta), \text{ for } i \in \{1, \dots, N\} \quad (2.1)$$

$$\mathcal{D} = \{x_1, \dots, x_N\} \quad (2.2)$$

$$N_1 = \sum_{i=1}^N \mathbb{I}(x_i = 1) \quad (2.3)$$

$$N_0 = \sum_{i=1}^N \mathbb{I}(x_i = 0) \quad (2.4)$$

Likelihood

$$p(\mathcal{D}|\theta) = \theta^{N_1} (1 - \theta)^{N_0} \quad (2.5)$$

Prior

$$p(\theta) = \text{Beta}(\theta|a, b) \quad (2.6)$$

Posterior

$$p(\theta|\mathcal{D}) = \text{Beta}(\theta|a' = N_1 + a, b' = N_0 + b) \quad (2.7)$$

Posterior predictive

$$p(\tilde{x} = 1|\mathcal{D}) = \frac{a'}{a' + b'} \quad (2.8)$$

Evidence

2.1.2 Derivations

2.2 Beta-binomial model

2.2.1 Summary

The model

$$N_1 \sim \text{Bin}(N, \theta) \quad (2.9)$$

$$\mathcal{D} = \{N_1, N\} \quad (2.10)$$

$$N_1 = \text{number of successes} \quad (2.11)$$

$$N = \text{total number of trials} \quad (2.12)$$

$$\tilde{\mathcal{D}} = \{\tilde{N}_1, \tilde{N}\} \quad (2.13)$$

$$\tilde{N}_1 = \text{number of successes in a new batch of data} \quad (2.14)$$

$$\tilde{N} = \text{total number of trials in a new batch of data} \quad (2.15)$$

Likelihood

$$p(\mathcal{D}|\theta) = \text{Bin}(N_1|N, \theta) \quad (2.16)$$

Prior

$$p(\theta) = \text{Beta}(\theta|a, b) \quad (2.17)$$

Posterior

$$p(\theta|\mathcal{D}) = \text{Beta}(\theta|a' = N_1 + a, b' = N_0 + b) \quad (2.18)$$

Posterior predictive

$$p(\tilde{\mathcal{D}}|\mathcal{D}) = \text{Bb}(\tilde{N}_1; a', b', \tilde{N}) \quad (2.19)$$

Evidence

2.2.2 Derivations

2.3 Dirichlet-categorical model

2.3.1 Summary

The model

$$X_i \sim \text{Cat}(\boldsymbol{\theta} = (\theta_1, \dots, \theta_K)^T), \text{ for } i \in \{1, \dots, N\} \quad (2.20)$$

$$\mathcal{D} = \{x_1, \dots, x_N\} \quad (2.21)$$

$$n_k = \sum_{i=1}^N \mathbb{I}(x_i = k) \quad (2.22)$$

Likelihood

$$p(\mathcal{D}|\theta) = \prod_{k=1}^K \theta_k^{n_k} \quad (2.23)$$

Prior

$$p(\theta) = \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}) \quad (2.24)$$

Posterior

$$p(\theta|\mathcal{D}) = \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}' = \boldsymbol{\alpha} + (n_1, \dots, n_K)^T) \quad (2.25)$$

Posterior predictive

$$p(\tilde{X} = j|\mathcal{D}) = \frac{\alpha'_j}{\sum_{k=1}^K \alpha'_k} \quad (2.26)$$

$$= \frac{\alpha_j + n_j}{\alpha_0 + N} \quad (2.27)$$

$$\text{where } \alpha_0 = \sum_{k=1}^K \alpha_k \quad (2.28)$$

Evidence

2.3.2 Derivations

2.4 Dirichlet-multinomial model

2.4.1 Summary

The model

$$\mathbf{N} \sim \text{Mult}(N, \boldsymbol{\theta}) \in \mathbb{R}^K \quad (2.29)$$

$$\mathcal{D} = \{\mathbf{n} = \text{vector of counts of successes}\} \quad (2.30)$$

$$N = \sum_{i=1}^K n_i \quad (2.31)$$

$$\tilde{\mathcal{D}} = \{\tilde{\mathbf{n}} = \text{vector of counts of successes in a new batch of data}\} \quad (2.32)$$

$$\tilde{N} = \sum_{i=1}^K \tilde{n}_i \quad (2.33)$$

Likelihood

$$p(\mathcal{D}|\theta) = \text{Mult}(\mathbf{n}; N, \boldsymbol{\theta}) \quad (2.34)$$

Prior

$$p(\theta) = \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}) \quad (2.35)$$

Posterior

$$p(\theta|\mathcal{D}) = \text{Dir}(\boldsymbol{\theta}; \boldsymbol{\alpha}' = \boldsymbol{\alpha} + (n_1, \dots, n_K)^T) \quad (2.36)$$

Posterior predictive

$$p(\tilde{\mathcal{D}}|\mathcal{D}) = \frac{\Gamma(\alpha_0 + N)}{\Gamma(\alpha_0 + N + \tilde{N})} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k + \tilde{n}_k)}{\Gamma(\alpha_k + n_k)} \quad (2.37)$$

$$\text{where } \alpha_0 = \sum_{k=1}^K \alpha_k \quad (2.38)$$

Evidence

$$p(\mathcal{D}|\boldsymbol{\alpha}) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_0 + N)} \prod_{k=1}^K \frac{\Gamma(\alpha_k + n_k)}{\Gamma(\alpha_k)} \quad (2.39)$$

2.4.2 Derivations

2.5 Poisson-gamma model

2.5.1 Summary

The model

$$x \sim \text{Poi}(\lambda) \quad (2.40)$$

$$\mathcal{D} = \{x_1, \dots, x_N\} \quad (2.41)$$

Likelihood

$$p(\mathcal{D}|\lambda) = \prod_{i=1}^N \frac{\lambda^{x_i}}{x_i!} \exp(-\lambda) \quad (2.42)$$

Prior

$$p(\lambda) = \text{Gamma}(\lambda; a, b) \quad (2.43)$$

Posterior

$$p(\lambda|\mathcal{D}) = \text{Gamma}\left(\lambda; a' = a + \sum_{i=1}^N x_i, b' = b + N\right) \quad (2.44)$$

Posterior predictive

$$p(\tilde{x}|\mathcal{D}) = \text{NB}(\tilde{x}|a', \frac{1}{1+b'}) \quad (2.45)$$

Evidence

$$p(\mathcal{D}) = \quad (2.46)$$

2.5.2 Derivations

3 Sampling algorithms

3.1 Introduction

Let p be a probability distribution with pdf $p(\mathbf{x})$, which we assume can be evaluated only up to a constant of proportionality (i.e. we can only evaluate $p^*(\mathbf{x}) = Z_p p(\mathbf{x})$, where $Z_p = \int p^*(\mathbf{x}) d\mathbf{x}$). We want to achieve the following:

Problem 1 Generate samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(R)}\}$ (shorthand notation $\{\mathbf{x}^{(r)}\}$) from the probability distribution p .

Problem 2 Estimate the expectation of an arbitrary function $f(\mathbf{X})$ given that $\mathbf{X} \sim p$, $\mathbb{E}[f(\mathbf{X})]$.

3.2 Importance sampling

Assume that we can sample from a proposal distribution q with a pdf $q(\mathbf{x})$, which can be evaluated only up to a constant of proportionality (i.e. we can only evaluate $q^*(\mathbf{x}) = Z_q q(\mathbf{x})$). To solve problem 2, we follow

1. Generate samples from q , $\{\mathbf{x}^{(r)}\}$.
2. Calculate importance weights $w_r = \frac{p^*(\mathbf{x}^{(r)})}{q^*(\mathbf{x}^{(r)})}$.
3. $\hat{\mathbf{y}} = \frac{\sum_r w_r f(\mathbf{x}^{(r)})}{\sum_r w_r}$ is the estimator of $\mathbb{E}[f(\mathbf{X})]$.

3.2.1 Convergence of estimator as R increases

We want to prove that if $q(\mathbf{x})$ is non-zero for all \mathbf{x} where $p(\mathbf{x})$ is non-zero, the estimator $\hat{\mathbf{y}}$ converges to $\mathbb{E}[f(\mathbf{X})]$, as R increases. We consider the the expectations of the numerator and denominator separately:

$$\mathbb{E}_q[\text{numerator}] = \mathbb{E}_q \left[\sum_r w_r f(\mathbf{x}^{(r)}) \right] \quad (3.1)$$

$$= \sum_r \mathbb{E}_q \left[w_r f(\mathbf{x}^{(r)}) \right] \quad (3.2)$$

$$= \sum_r \mathbb{E}_q \left[\frac{p^*(\mathbf{x}^{(r)})}{q^*(\mathbf{x}^{(r)})} f(\mathbf{x}^{(r)}) \right] \quad (3.3)$$

$$= \sum_r \mathbb{E}_q \left[\frac{Z_p p(\mathbf{x}^{(r)})}{Z_q q(\mathbf{x}^{(r)})} f(\mathbf{x}^{(r)}) \right] \quad (3.4)$$

$$= \frac{Z_p}{Z_q} \sum_r \int_{\mathbf{x}^{(r)}} p(\mathbf{x}^{(r)}) f(\mathbf{x}^{(r)}) d\mathbf{x}^{(r)} \quad (3.5)$$

$$= \frac{Z_p}{Z_q} \sum_r \mathbb{E}_p \left[f(\mathbf{x}^{(r)}) \right] \quad (3.6)$$

$$= \frac{Z_p}{Z_q} R \mathbb{E}_p [f(\mathbf{x})] \quad (3.7)$$

$$\mathbb{E}_q[\text{denom}] = \mathbb{E}_q \left[\sum_r w_r \right] \quad (3.8)$$

$$= \sum_r \mathbb{E}_q \left[\frac{p^*(\mathbf{x}^{(r)})}{q^*(\mathbf{x}^{(r)})} \right] \quad (3.9)$$

$$= \sum_r \mathbb{E}_q \left[\frac{Z_p p(\mathbf{x}^{(r)})}{Z_q q(\mathbf{x}^{(r)})} \right] \quad (3.10)$$

$$= \frac{Z_p}{Z_q} \sum_r \int_{\mathbf{x}^{(r)}} p(\mathbf{x}^{(r)}) d\mathbf{x}^{(r)} \quad (3.11)$$

$$= \frac{Z_p}{Z_q} R \quad (3.12)$$

Hence $\hat{\mathbf{y}}$ converges to $\mathbb{E}_p[f(\mathbf{x})]$ as R increases (but is not necessarily an unbiased estimator because $\mathbb{E}_q[\hat{\mathbf{y}}]$ is not necessarily $= \mathbb{E}_p[f(\mathbf{x})]$).