Feel free to work with other students, but make sure you write up the homework and code on your own (no copying homework *or* code; no pair programming). Feel free to ask students or instructors for help debugging code or whatever else, though.

---

**1** (**Murphy 2.16**) Suppose $\theta \sim \text{Beta}(a, b)$ such that

$$\mathbb{P}(\theta; a, b) = \frac{1}{B(a, b)} \theta^{a-1}(1 - \theta)^{b-1} = \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \theta^{a-1}(1 - \theta)^{b-1}$$

where $B(a, b) = \Gamma(a)\Gamma(b)/\Gamma(a + b)$ is the Beta function and $\Gamma(x)$ is the Gamma function. Derive the mean, mode, and variance of $\theta$.

---

**Mean** *The mean value is the expected value;*

$$\mathbb{E}(\theta) = \int \theta \mathbb{P}(\theta; a, b) d\theta = \int \theta \frac{1}{B(a, b)} \theta^{a-1}(1 - \theta)^{b-1} d\theta$$

$$= \frac{1}{B(a, b)} \int \theta^a (1 - \theta)^{b-1} d\theta$$

$$= \frac{B(a + 1, b)}{B(a, b)}$$

$$= \left[ \frac{\Gamma(a + 1)\Gamma(b)}{\Gamma(a + b + 1)} \right] \left[ \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \right]$$

$$= \left[ \frac{a\Gamma(a)\Gamma(b)}{a + b\Gamma(a + b)} \right] \left[ \frac{\Gamma(a + b)}{\Gamma(a)\Gamma(b)} \right]$$

$$= \boxed{\frac{a}{a + b}}$$

**Variance** *We know that this is* $\mathbb{E}[\theta^2] - \mathbb{E}[\theta]^2$;

$$\mathbb{E}[\theta^2] - \mathbb{E}[\theta]^2 = \int \theta^2 \mathbb{P}(\theta; a, b) d\theta - \frac{a^2}{(a+b)^2}$$

$$= \frac{1}{B(a,b)} \int \theta^{a+1}(1-\theta)^{b-1} d\theta - \frac{a^2}{(a+b)^2}$$

$$= \frac{B(a+2,b)}{B(a,b)} - \frac{a^2}{(a+b)^2}$$

$$= \frac{\Gamma(a+2)\Gamma(b)}{\Gamma(a+b+2)} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} - \frac{a^2}{(a+b)^2}$$

$$= \frac{a(a+1)\Gamma(a)\Gamma(b)}{(a+b)(a+b+1)\Gamma(a+b)} \frac{\Gamma(a+b)}{\Gamma(a)\Gamma(b)} - \frac{a^2}{(a+b)^2}$$

$$= \frac{a(a+1)}{(a+b)(a+b+1)} - \frac{a^2}{(a+b)^2}$$

$$= \frac{(a^2+a)(a+b) - a^2(a+b+1)}{(a+b+1)(a+b)^2}$$

$$= \boxed{\frac{ab}{(a+b+1)(a+b)^2}}$$

**Mode** *If we visualize a distribution, the point at which there is the greatest number of values is the mode, it is also the only global extrema, as such, it is the only point where* $\nabla_\theta = 0$;

$$0 = \nabla_\theta \left( \frac{1}{B(a,b)} \theta^{a-1}(1-\theta)^{b-1} \right)$$

$$= \nabla_\theta(\theta^{a-1}(1-\theta)^{b-1}$$

$$= (a-1)\theta^{a-2}(1-\theta)^{b-1} - (b-1)\theta^{a-1}(1-\theta)^{b-2}$$

*Therefore we can now solve for the specific value of* $\theta$ *that makes this true;*

$$(a-1)\theta^{a-2}(1-b)^{b-1} = (b-1)\theta^{a-1}(1-\theta)^{b-2}$$

$$(a-1)(1-\theta) = (b-1)\theta$$

$$\theta = \frac{a-1}{a+b-2}$$

$$\blacksquare$$

**2 (Murphy 9)** Show that the multinoulli distribution

$$\text{Cat}(\mathbf{x}|\boldsymbol{\mu}) = \prod_{i=1}^{K} \mu_i^{x_i}$$

is in the exponential family and show that the generalized linear model corresponding to this distribution is the same as multinoulli logistic regression (softmax regression).

We are searching for something of the form $\mathbf{p}(\mathbf{y}; \boldsymbol{\eta} = b(\boldsymbol{\eta}) \exp\left(\boldsymbol{\eta}^T T(\mathbf{y}) + a(\boldsymbol{\eta})\right)$, as we did in class I am going to force the multinoulli distribution $\prod_{i=1}^{K} \mu_i^{x_i}$ into its exponential form:

$$\prod_{i=1}^{K} \mu_i^{x_i} = \exp\left[\log\left(\prod_{i=1}^{K} \mu_i^{x_i}\right)\right] \tag{1}$$

$$= \exp\left[\sum_{i=1}^{K} \log(\mu_i^{x_i})\right] \tag{2}$$

$$= \exp\left[\sum_{i=1}^{K} x_i \log(\mu_i)\right] \tag{3}$$

Now lets note that since this is a multinoulli distribution;

$$\mu_k = 1 - \sum_{i=1}^{K-1} \mu_i \text{ and } x_k = 1 - \sum_{i=1}^{K-1} x_i$$

$$= \exp\left[\sum_{i=1}^{K-1} x_i \log(\mu_i) + \left(1 - \sum_{i=1}^{K-1} x_i\right) \log(\mu_k)\right] \tag{4}$$

$$= \exp\left[\sum_{i=1}^{K-1} x_i \log\left(\frac{\mu_i}{\mu_k}\right) + \log(\mu_k)\right] \tag{5}$$

Lets continue by defining $\boldsymbol{\eta}$ as follows;

$$\boldsymbol{\eta} = \begin{pmatrix} \log\left(\frac{\mu_1}{\mu_k}\right) \\ \vdots \\ \log\left(\frac{\mu_{k-1}}{\mu_k}\right) \end{pmatrix}$$

It follows that we can define $\mu_k$ as $1 - \sum_{i=1}^{K-1} \mu_k \exp(\eta_i)$ or equivalently as $\frac{1}{1+\sum_{i=1}^{K-1} exp(\eta_i)}$
Thus picking up where we left off;

$$= \exp\left[\boldsymbol{\eta}^T \mathbf{x} + \log\left(\frac{1}{1 + \sum_{i=1}^{K-1} \exp(\eta_i)}\right)\right] \tag{6}$$

This is in exponential form with $b(y) = 1$, $T(y) = \mathbf{x}$ and $a(\boldsymbol{\eta}) = -\log(1 + \sum_{i=1}^{K-1} \exp(\eta_i))$

∎