# ECON5529
# Bayesian Theory

Ellis Scharfenaker

Fall 2017, Lecture 4

## Multiparameter Models

■ For most practical models we are interested in more than a single unknown parameter. In this case we are interested in sampling from a higher dimensional joint distribution of all unknowns.

■ In some cases all parameters will be of interest to us, but in others they may merely be a *nuisance*. In this case, the ultimate aim is to obtain the *marginal* posterior distribution of the particular parameters of interest.

■ In general, multiparameter models pose a more challenging optimization problem.

## Marginal Distributions

■ The posterior distribution in a Bayesian analysis describes a joint distribution over all of the parameters of a model.

■ For example, in the linear regression model,

$$y = \beta_0 + \beta_1 x + \epsilon$$
$$\epsilon = y - (\beta_0 + \beta_1 x)$$

(1)

■ We typically use the normal distribution to model the error terms making the likelihood:

$$p[\epsilon \mid \sigma^2, \beta_0, \beta_1] = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{\epsilon^2}{2\sigma^2}} = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y-(\beta_0+\beta_1 x))^2}{2\sigma^2}}$$

(2)

■ The parameters are the intercept $\beta_0$, the slope $\beta_1$, and $\sigma$, the standard deviation by which we scale the errors.

■ But we are often much more interested in what the data analysis tells us about some parameters than about others.

■ In these cases the less interesting parameters are called *nuisance parameters*. In the regression example, the error standard deviation $\sigma$ is often regarded as such a nuisance parameter.

■ What we really want to know is the posterior distribution of $\beta_0$ and $\beta_1$ properly averaging out over all the possibilities for $\sigma$.

■ In principle it is easy enough to see that we would like to know the marginal distribution for the parameter(s) of interest after integrating out over the posterior probability of the nuisance parameters, for example:

$$p[\beta_0, \beta_1 \mid x] = \int p[\beta_1, \beta_0, \sigma \mid x] \, d\sigma$$

(3)

## Nuisance Parameters

■ Suppose that our multiparameter model contains two unknown or unobserved quantities $\theta = \{\theta_1, \theta_2\}$ and that we are only really interested in making inferences with $\theta_1$. In such a case $\theta_2$ is the *nuisance parameter.*

■ In this situation we are interested in the conditional distribution of the parameter of interest given the data, $p[\theta_1 \mid x]$, but face the possible dependence of $\theta_1$ on $\theta_2$ so that $p[\theta_1 \mid x, \theta_2] \neq p[\theta_1 \mid x]$

■ To obtain the conditional distribution of interest ($p[\theta_1 \mid x]$) apply the product rule to the joint posterior for the unknown parameter:

$$p[\theta_1, \theta_2 \mid x] = p[\theta_1 \mid \theta_2, x]\, p[\theta_2 \mid x]\, p[x] \qquad (4)$$

■ But note the $p[x]$ cancels out when normalizing.

$$p[\theta_1, \theta_2 \mid x] = p[\theta_1 \mid \theta_2, x]\, p[\theta_2 \mid x] \qquad (5)$$

■ Averaging over the nuisance parameter $\theta_2$ we properly eliminate it from the posterior distribution:

$$p[\theta_1 \mid x] = \int_{\theta_2} p[\theta_1 \mid \theta_2, x]\, p[\theta_2 \mid x]\, d\theta_2 \qquad (6)$$

■ This result shows that the marginal posterior distribution ($p[\theta_1 \mid x]$) is a mixture of the conditional posterior distributions given $\theta_2$, ($p[\theta_1 \mid \theta_2, x]$) weighted by $p[\theta_2 \mid x]$, the possible values of $\theta_2$.

■ The weights, $p[\theta_2 \mid x]$, are the *posterior density $p[\theta_2 \mid x] \propto p[x \mid \theta_2]\, p[\theta_2]$* of $\theta_2$.

■ Rarely do we evaluate the integral (6) explicitly as posterior distributions can be obtained by marginal and conditional simulation.

■ Simulating an integral like (6) requires a **sampling strategy**. In this case it is quite straightforward.

  ■ (1) Draw a sample $\tilde{\theta}_2$ from its marginal posterior distribution $p[\theta_2 \mid x]$

  ■ (2) Sample $\tilde{\theta}_1$ from its conditional posterior distribution, $p[\theta_1 \mid \tilde{\theta}_2, x]$, using the $\tilde{\theta}_2$ draw from step 1.

■ In this way the integral (6) is computed indirectly and much of the remainder of the course will rely on posterior sampling methods such as this.

## Example: Normal Data with Conjugate prior

■ To illustrate multiparameter modeling and posterior marginalization of nuisance parameters let's take the rather typical example of estimating the mean of a population from a sample.

■ The data in this model is a list $x = \{x_1, x_2, ..., x_n\}$ of $n$ independent observations from a univariate normal distribution $x \sim N[\mu, \sigma^2]$.

■ One choice for a noninformative prior for $\mu$ and $\sigma^2$ assuming prior independence is Jeffreys' prior, which in this case is $p[\mu, \sigma^2] \propto \sigma^{-2}$

■ The **joint posterior** is then equal to the likelihood multiplied by the factor $\sigma^{-2}$

$$p[\mu, \sigma^2 \mid x] = \prod_{i=1}^{n} \frac{1}{\sqrt{2\pi}\,\sigma}\, e^{-\frac{(x_i - \mu)^2}{2\sigma^2}}\, \sigma^{-2}$$

$$\propto \sigma^{-n-2}\, e^{-\frac{1}{2\sigma^2} \Sigma_{i=1}^{n} (x_i - \mu)^2} \qquad (7)$$

$$= \sigma^{-n-2}\, e^{-\frac{1}{2\sigma^2} \Sigma_{i=1}^{n} \left((x_i - \bar{x})^2 + n(\bar{x} - \mu)^2\right)}$$

$$= \sigma^{-n-2}\, e^{-\frac{1}{2\sigma^2} \left((n-1)\, s^2 + n(\bar{x} - \mu)^2\right)}$$

- Where $\tilde{x} = \frac{1}{n} \Sigma_i x_i$ is the sample mean and $s^2 = \frac{1}{n-1} \Sigma_i (x_i - \tilde{x})^2$ is the sample variance. But remember we are interested in the the conditional posterior distribution of $\mu$ after eliminating $\sigma^2$.

$$p[\mu \mid x] = \int p[\mu, \sigma^2 \mid x] \, d\sigma^2 = \int p[\mu \mid \sigma^2, x] \, p[\sigma^2 \mid x] \, d\sigma^2 \tag{8}$$

- To calculate the conditional posterior distribution of $\mu$ we first consider the conditional posterior density $p[\mu \mid \sigma^2, x]$ and then the marginal posterior distribution $p[\sigma^2 \mid x]$.

- For the mean of a normal distribution with known variance and a conjugate prior distribution we can use the results derived from the previous lecture $p[\mu \mid \sigma^2, x] \sim \mathcal{N}[\hat{\mu}, \hat{\tau}^2]$ with:

$$\hat{m} = \frac{\frac{1}{\tau^2} m + \frac{1}{\sigma^2} \tilde{x}}{\frac{1}{\tau^2} + \frac{1}{\sigma^2}}, \quad \hat{\tau}^2 = \left( \frac{1}{\tau^2} + \frac{1}{\sigma^2} \right)^{-1}$$

- For the marginal distribution $p[\sigma^2 \mid x]$ we need to average over the joint distribution

$$p[\sigma^2 \mid x] \propto \int \sigma^{-n-2} e^{-\frac{1}{2\sigma^2} \left( (n-1) s^2 + n (\tilde{x} - \mu)^2 \right)} \, d\mu$$

$$= \sigma^{-n-2} e^{-\frac{1}{2\sigma^2} (n-1) s^2} \int e^{-\frac{1}{2\sigma^2} n (\tilde{x} - \mu)^2} \, d\mu \tag{9}$$

- The integral to be evaluated is just a simple normal integral

```
Integrate[e^(-1/(2σ²) n (x-μ)²), {μ, -∞, ∞}, Assumptions → {Re[n/σ²] ≥ 0}]
```

$$\frac{\sqrt{2\pi}}{\sqrt{\frac{n}{\sigma^2}}}$$

$$p[\sigma^2 \mid x] \propto \sigma^{-n-2} \sqrt{\frac{2\pi\sigma^2}{n}} \, e^{-\frac{1}{2\sigma^2}(n-1) s^2} \tag{10}$$

$$\propto (\sigma^2)^{\frac{-(n+1)}{2}} e^{-\frac{(n-1)s^2}{2\sigma^2}}$$

- In this case $\sigma^2 \sim \mathcal{I}nv.\chi^2[n-1, s^2]$.

- For $\theta \sim \mathcal{I}nv\text{-}\chi^2$, the density function is:

$$p[\theta] = \frac{2^{-v/2}}{\Gamma[v/2]} \theta^{-\frac{v}{2}+1} e^{-\frac{v s^2}{2\theta}} \tag{11}$$

- The Inv-$\chi^2$ distribution is the same as a $\mathcal{I}G[\alpha = \frac{v}{2}, \beta = \frac{1}{2}]$.

## Sampling Strategy

- (1) Draw a single $\sigma^2$ sample from $p[\sigma^2 \mid x] = \mathcal{I}nv.\chi^2[n-1, s^2]$.

- (2) Given the $\sigma^2$ draw $\mu$ from $p[\mu \mid \sigma^2, x] = \mathcal{N}[\hat{\mu}, \hat{\tau}^2]$
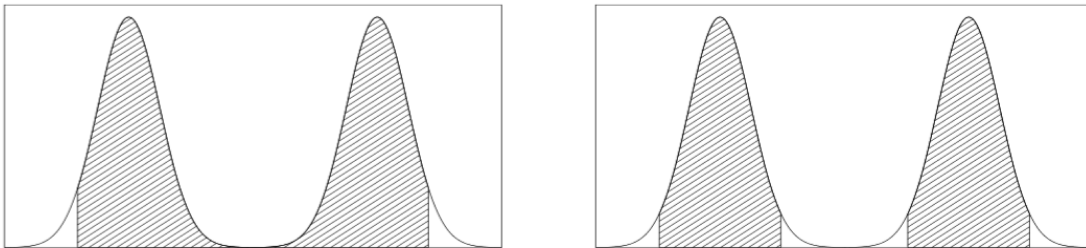
- (3) Repeat

In R:

## Reporting the Posterior

■ Generally the first stage of investigating the posterior probability of the parameters of a model is to find the maximum posterior probability value of the parameters.

■ A local maximum of a probability distribution is often called the *mode*, so the maximum posterior probability parameters are sometimes called the *modal parameters*.

■ Other commonly used summaries of posteriors are the mean, and median of the distribution; variation is commonly summarized by the standard deviation, the interquartile range, and other quantiles.

## High posterior probability sets or intervals

■ If posterior probabilities were always normal, it would suffice to put "error bars" around parameter estimates representing two standard deviations of the posterior probability distribution.

■ It turns out that posterior probabilities for relatively small numbers of data points, such as economists often work with, can deviate significantly from normal distributions.

■ Hence, in addition to point summaries, it is always important to report ***posterior uncertainty***. Our usual approach is to present a central interval of posterior probability.

■ The ***highest posterior density*** (HPD) region is the set of values that contains $100 (1 - \alpha)$ % of the posterior probability.

■ The HPD indicates which points of a distribution are most credible, and which cover most of the distribution.

■ The *highest posterior density* is identical to a central posterior interval if the posterior distribution is unimodal and symmetric. However, if the posterior is skewed or multimodal, the HPD can be quite different from a standard confidence interval.



## Multinomial Model

■ The multinomial model is a generalization of the binomial model for categorical outcomes where the possible qualitative indicators are greater than two.

■ In the multinomial model the data is a list of counts of the number of observations of each outcome $x = \{x_1, ..., x_k\}$ where $x_i$ is the number of observations for the $i^{th}$ outcome category. The probability of a particular outcome is $\theta_i$ where $\sum \theta_i = 1$.

■ The likelihood of the data is

$$p[x \mid \theta] \propto \prod_{i=1}^{k} \theta_i^{x_i} = \theta_1^{\sum x_1} ... \theta_k^{\sum x_k} \tag{12}$$

- The conjugate prior for the multinomial model is called the *Dirichlet prior.* Just as the multinomial distribution is a multivariate generalization of the binomial distribution, the Dirichlet distribution is a multivariate generalization of the Beta distribution and is defined as:

$$\theta \sim \mathcal{D}[\alpha], \quad p[\theta] \propto \prod_{i=1}^{k} \theta_i^{\alpha_i - 1}$$

- Similar to the beta distribution we get a uniform prior on all $\theta_j$ when $\alpha_j = 1 \; \forall \; j$.

- The posterior distribution for the unknown probabilities $\theta_j$ using the conjugate Dirichlet prior is:

$$p[\theta \mid x] \propto \prod_{i=1}^{k} \theta_i^{x_i} \prod_{i=1}^{k} \theta_i^{\alpha_i - 1} = \prod_{i=1}^{k} \theta_i^{(x_i + \alpha_i - 1)} \propto \mathcal{D}[x + \alpha]$$

## Example

- During the 2016 presidential election an NYT/CBS News Poll surveyed $n = 1433$ individuals about. The results were:

|  | Registered Voters |
|---|---|
| Clinton | 659 |
| Trump | 588 |
| Undecided / Other | 129 |

- The difference $\theta_1 - \theta_2$ is the difference between Clinton vs. Trump support and represents the margin of Clinton over Trump.

- If the prior over $(\theta_1, \theta_2, \theta_3)$ is $\mathcal{D}[\theta_1, \theta_2, \theta_3]$ then the posterior is

$$\mathcal{D}[\alpha_1 + x_1, \; \alpha_2 + x_2, \; \alpha_3 + x_3] \propto \prod_{i=1}^{3} \theta_i^{(x_i + \alpha_i - 1)} = \theta_1^{(x_1 + \alpha_1 - 1)} \theta_2^{(x_2 + \alpha_2 - 1)} \theta_3^{(x_3 + \alpha_3 - 1)}$$

- Note that most of the mass of this distribution is on positive values, indicating that there is strong evidence that the proportion of voters for Clinton exceeds the proportion for Trump.

- Now suppose we wish to predict the total number of electoral votes ($EV_0$) obtained by Clinton. Let $\theta_{Cj}$ and $\theta_{Tj}$ denote the unknown proportion of voters respectively for Clinton and Trump in the $j^{th}$ state.

- The number of electoral votes for Clinton (using the indicator function ) is:

$$EV_0 = \sum_{j=1}^{51} EV_j \, I\big[\theta_{Cj} > \theta_{Tj}\big]$$

- Let $q_{Cj}$ and $q_{Tj}$ denote the observed data of sample proportions of voters for Clinton and Trump in the $j^{th}$ state.

- Further, assume that each poll is based on a sample of 500 voters.

- We can assign a uniform prior on the vector of proportions, $(\theta_{C1}, \theta_{T1}), \dots, (\theta_{C51}, \theta_{T51})$, i.e. for sate $j$ $(0_{Cj}, 100_{Tj})$ has the same probability as $(100_{Cj}, 0_{Tj})$.

- Thus, the vectors $(\theta_{C1}, \theta_{T1}), \dots, (\theta_{C51}, \theta_{T51})$ have independent posterior distributions, where the proportions favoring the candidates in the $j^{th}$ state, $(\theta_{Cj}, \theta_{Tj}, 1 - \theta_{Cj} - \theta_{Tj})$, have a Dirichlet distribution.

- It would be possible to improve our prediction by using more than a single poll in each state as is commonly done by Nate Silver on FiveThirtyEight.com