

CHE384, From Data to Decisions: Measurement, Uncertainty, Analysis, and Modeling

Lecture 75

Bayesian Regression, part 2

Chris A. Mack
Adjunct Associate Professor

<http://www.lithoguru.com/scientist/statistics/>

© Chris Mack, 2016 Data to Decisions 1

Bayesian Regression

$$P(\theta|y) = \frac{P(y|\theta)P(\theta)}{P(y)}$$

Posterior Distribution Likelihood Function Prior Distribution

Normalizing constant

- We already know how to calculate the likelihood function (we make assumptions about the pdf of y)
- The term $P(y)$ is a constant (independent of θ) and used to normalize the posterior pdf
- Prior distribution $P(\theta)$: our understanding of the model parameters and σ_ϵ before we began the experiment
- How do we interpret the posterior distribution?

© Chris Mack, 2016 Data to Decisions 2

Posterior Distribution

- The output of a Bayesian regression is not a set of best fit parameters, but a probability distribution for each parameter and σ_ϵ (called the **posterior distribution**)
- How do we interpret the posterior distribution?
- We need to **summarize** the distribution:
 - Use the mode, mean, median, or range midpoint as an equivalent "best estimate" of the parameter
 - Use the distribution to calculate "credible interval" (quantiles) for the parameter

© Chris Mack, 2016 Data to Decisions 3

Posterior Distribution

- Posterior distribution describes how much the data has changed our prior beliefs
- Bernstein-von Mises Theorem**: for a sufficiently large sample size, the posterior distribution becomes independent of the prior distribution (so long as the prior is not either 0 or 1)
 - The posterior tends towards a normal distribution with a mean equal to the MLE (assuming iid data), a restatement of the central limit theorem
 - The effect of the prior diminishes as the amount of data increases

© Chris Mack, 2016 Data to Decisions 4

Prior Distribution

- The prior distribution $P(\theta)$ is really shorthand notation for $P(\theta|I)$, where I is all the information we have about the problem before we start collecting data
 - If we have NO information about what the parameters could or should be, then $P(\theta|I)$ is a constant (called an **uninformative prior** or **objective prior**), and the posterior distribution equals the likelihood function
 - We almost always have some information

© Chris Mack, 2016 Data to Decisions 5

Prior Choice

- Non-informative** (baseline or objective) prior
 - Ex: a uniform probability over the expected range of possible values
 - Flat priors are not always uninformative! Ex: should we have a uniform distribution of slopes, or uniform distribution of the angle of the line, or its sine?
- Substantive** (informative) prior
 - Use some problem-specific information to provide a prior distribution for each model parameter
 - Based on previous data, experiments, knowledge
 - Sometimes one can assume the prior for each parameter is independent of the others ($P(\theta) = P(\beta)P(\sigma^2)$), but frequently a joint probability distribution is required
 - Setting the prior to a delta function fixes a parameter independent of the data (we never do this in general)

© Chris Mack, 2016 Data to Decisions 6

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Reparameterization

- Often, reparameterization of the model is required to make prior assignment easier
- Consider a linear regression with $\hat{y}_i = \beta_0 + \beta_1 x_i$
 - But, $\beta_0 = \hat{y}_i(x_i = 0)$, and often $x_i = 0$ is not physically meaningful
 - Shift the x-axis by \bar{x} , giving $\hat{y}_i = \beta_0' + \beta_1(x_i - \bar{x})$
 - Now, $\beta_0' = \hat{y}_i(x_i = \bar{x})$, which is meaningful
 - Define the prior for β_0' (which we can assume to be independent of β_1 and its prior)

© Chris Mack, 2016 Data to Decisions 7

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Prior Choice Example

- Consider a linear regression with $\hat{y}_i = \beta_0 + \beta_1 x_i$ Bayes: $P(\theta|y) \propto P(y|\theta)P(\theta)$
- We can set our prior for the slope to favor the “null hypothesis”: $P(\beta_1) \sim N(0, \sigma_b^2)$ where σ_b is large enough to allow for the expected range of possible slopes
- Does the data contain information to push us away from our prior belief that there is no relationship between x and y ?

© Chris Mack, 2016 Data to Decisions 8

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Conjugate Priors

- For a given likelihood distribution, analytical solutions of the posterior distribution are only possible for special cases of priors (called **conjugate priors**)
 - Example: For iid normal errors, the conjugate prior for β is normal, and for σ^2 is inverse gamma
- Usually, we need to solve Bayes' equation numerically
 - Markov Chain Monte Carlo (MCMC) sampling
 - Result is a set of points from the posterior distribution that we then summarize (mean, or maximum a posteriori – MAP – estimate of the mode)

© Chris Mack, 2016 Data to Decisions 9

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Relationship to OLS

- Consider a linear regression with $\hat{y}_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots$
- Assume our likelihood function is a product of Normal probabilities with constant variance
- Choose a baseline prior so that that the prior $P(\beta, \sigma^2) \propto 1/\sigma^2$ (**Jeffreys prior**, it is “improper”)
- The resulting posterior distribution is t-distributed about the MLE parameter estimates
- The results are identical to OLS**, but with a different interpretation (credible intervals rather than confidence intervals)

© Chris Mack, 2016 Data to Decisions 10

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Interpretation

- 95% Confidence Interval** (frequentist)
 - Our parameter is an unknown constant, and with a large number of repeated samples, 95% of such calculated confidence intervals would include the true value of the parameter
- 95% Credible Interval/Region** (Bayesian)
 - Our parameter is a random variable, with a 95% probability of falling within the given interval

© Chris Mack, 2016 Data to Decisions 11

THE UNIVERSITY OF TEXAS AT AUSTIN WHAT STARTS HERE CHANGES THE WORLD

Lecture 75: What have we learned?

- How is the posterior distribution used (summarized) to tell us about model parameters?
- What does the Bernstein–von Mises theorem tell us about the relationship between the prior and posterior distributions?
- What is the difference between uninformative and substantive priors?
- What is the Jeffreys prior and how does it apply to linear regression?
- Explain the difference between confidence intervals and credible intervals

© Chris Mack, 2016 Data to Decisions 12