

STAT 431 — Applied Bayesian Analysis — Course Notes

Posterior Predictive Checks

Fall 2022

So far, we considered only ways to compare different data model/prior combinations to each other.

How can we check whether a data model/prior combination is a good fit to the data?

Frequentist approach: *lack-of-fit test* (such as a *chi-square test* — later)

Usually produces a p -value — smaller indicates more evidence against the data model.

A Bayesian wants to assess the prior and data model together.

Discrepancies

For notation:

\mathbf{y} = the data (vector)

$\boldsymbol{\theta}$ = the parameter (vector) in data model \mathcal{M}

We choose a numerical function called a **discrepancy**:

$$T(\mathbf{y}; \boldsymbol{\theta})$$

We intended it to measure how far observed data \mathbf{y} depart from what would be expected under data model \mathcal{M} with parameter value $\boldsymbol{\theta}$. Larger values should indicate greater departures from the data model.

An example of a discrepancy:

$$T(\mathbf{y}; \boldsymbol{\theta}) = \sum_{i=1}^n \frac{(y_i - \mathbb{E}(Y_i \mid \boldsymbol{\theta}))^2}{\text{Var}(Y_i \mid \boldsymbol{\theta})}$$

where the observed data $\mathbf{y} = (y_1, \dots, y_n)$ form a numerical vector, and (Y_1, \dots, Y_n) has its distribution under the data model.

This is larger when the y_i s are generally farther from their means than their variances would suggest, under the data model with parameter value $\boldsymbol{\theta}$.

(Note: Similar in form to the classical *chi-square statistic*.)

Frequentist Approach

A frequentist can't use $T(\mathbf{y}; \boldsymbol{\theta})$ directly, since it depends on the unknown $\boldsymbol{\theta}$.

Replacing $\boldsymbol{\theta}$ with an estimate $\hat{\boldsymbol{\theta}}$ might yield a reasonable test statistic

$$D(\mathbf{y}) = T(\mathbf{y}; \hat{\boldsymbol{\theta}})$$

but its distribution under data model \mathcal{M} might still depend on the unknown $\boldsymbol{\theta}$.

If the distribution is (approximately) known, a frequentist can compute the (approximate) p -value

$$p = \text{Prob}(D(\tilde{\mathbf{Y}}) \geq D(\mathbf{y}))$$

where $\tilde{\mathbf{Y}}$ is a *replication* of the data under its model.

In the earlier example:

$$D(\mathbf{y}) = \sum_{i=1}^n \frac{(y_i - \mathbb{E}(Y_i \mid \hat{\boldsymbol{\theta}}))^2}{\text{Var}(Y_i \mid \hat{\boldsymbol{\theta}})}$$

is the classical (Pearson) chi-square statistic (where $\hat{\boldsymbol{\theta}}$ is usually an MLE).

When the data model is correct, asymptotic theory often suggests

$$D(\mathbf{Y}) \mid \boldsymbol{\theta} \quad \dot{\sim} \quad \chi_{n-k}^2$$

if $\boldsymbol{\theta}$ effectively has k elements.

When this approximation is valid, an approximate p -value is the χ_{n-k}^2 PDF tail area to the right of $D(\mathbf{y})$.

In contrast, a Bayesian wants to

- ▶ assess the prior, not just the data model
- ▶ average over a distribution of θ (to avoid substituting an estimate $\hat{\theta}$)
- ▶ avoid any asymptotic approximations

Bayesian Approach

Suppose

$$\tilde{\mathbf{Y}} \mid \boldsymbol{\theta} \sim \mathcal{M}(\boldsymbol{\theta})$$

is conditionally independent of the data.

Note: Averaging its distribution over the posterior gives the *posterior predictive distribution* of the data.

Note: We can generate $\tilde{\mathbf{Y}}$ for a given $\boldsymbol{\theta}$ by simulating from the data model, which is usually easy.

A simulated value $\tilde{\mathbf{y}}$ would be called a **replicate** data set.

Then, instead of a frequentist p -value, a Bayesian could use a **posterior predictive p -value**

$$p_b = \text{Prob}(T(\tilde{\mathbf{Y}}; \boldsymbol{\theta}) \geq T(\mathbf{y}; \boldsymbol{\theta}) \mid \mathbf{y})$$

where the probability is over the joint posterior distribution of $\boldsymbol{\theta}$ and $\tilde{\mathbf{Y}}$.

Sufficiently small p_b indicates a problem with the data model and/or prior.

Usually p_b can't be directly computed. Instead, it can be approximated by posterior simulation (e.g., MCMC):

1. For each posterior-generated value θ , generate a replicate data set $\tilde{\mathbf{y}}$ (conditionally) and compute

$$T(\tilde{\mathbf{y}}; \theta) \quad \text{and} \quad T(\mathbf{y}; \theta)$$

2. Approximate p_b as the fraction of generated pairs $(\theta, \tilde{\mathbf{y}})$ for which

$$T(\tilde{\mathbf{y}}; \theta) \geq T(\mathbf{y}; \theta)$$

Example: Shark Attacks

Recall:

Y_i = number of shark attacks (worldwide)

X_i = year (2005–2017)

Our chosen data model and prior:

$$Y_i \mid \lambda_i \sim \text{indep Poisson}(\lambda_i)$$

$$\ln(\lambda_i) = \beta_1 + (X_i - \bar{X})\beta_2$$

$$\beta_1, \beta_2 \sim \text{iid Normal}(0, 100^2)$$

We will investigate fit using the chi-square discrepancy

$$T(\mathbf{y}; \beta_1, \beta_2) = \sum_{i=1}^n \frac{(y_i - \lambda_i)^2}{\lambda_i}$$

(Why is this the chi-square discrepancy?)

If there are problems with the Poisson regression or the (vague) normal priors, we expect $T(\mathbf{y}; \beta_1, \beta_2)$ to be large relative to its replicate version, leading to small p_b .

```

data {
  xmean <- mean(x)
  for(i in 1:length(x)) {
    xcent[i] <- x[i] - xmean
  }
}

model {
  for(i in 1:length(y)) {
    y[i] ~ dpois(lambda[i])
    log(lambda[i]) <- beta1 + xcent[i] * beta2

    yrep[i] ~ dpois(lambda[i])
  }

  beta1 ~ dnorm(0, 0.0001)
  beta2 ~ dnorm(0, 0.0001)

  chisq <- sum((y - lambda)^2 / lambda)
  chisqrep <- sum((yrep - lambda)^2 / lambda)
  pb.ind <- chisqrep >= chisq
}

```

Notes:

- ▶ The data set used here has only years 2005 to 2017 (and not the missing 2018 observation).
- ▶ JAGS automatically vectorizes arithmetic operations.

We will also try an alternative version of the Bayesian model that has a badly mis-specified prior ...

R/JAGS Example 5.2:

Checking a Poisson Regression

Remark: For some other choices of T , obtaining an especially *large* value of p_b (near 1) would also indicate a problem with the data model and/or prior.