

Xinyu Tan

– 2. Model Checking

- (a) Set up posterior predictive test quantities to check the following assumptions: 1) independent Poisson distributions 2) no trend over time.

To test Poisson distributions, we use Fano factor (σ^2/μ). For Poisson distribution, it is close to 1. To test if the data are from independent distributions, I think it can be tested along with "no trend over time" using autocorrelation.

- (b) Use simulations from the posterior predictive distributions to measure the discrepancies.

- (i) Assume that the numbers of fatal accidents in each year are independent with a Poisson distribution.

Hence fatal accident follows:

$$p(y|\theta) = \frac{1}{y!} \theta^y \exp(-\theta)$$

The likelihood:

$$p(y|\theta) = \prod_{i=1}^{10} \frac{1}{y_i!} \theta^{y_i} \exp(-\theta) \sim \theta^{10\bar{y}} \exp(-10\theta)$$

where $\bar{y} = \frac{1}{10} \sum_{i=1}^{10} y_i$.

Since I don't have any extra information, let's set prior distribution for $\theta \sim \text{beta}(0, 0)$.

Hence the posterior distribution

$$p(\theta|y) \sim \text{Gamma}(10\bar{y}, 10)$$

To compare the posterior predictive test quantities, we will perform the following sampling 1000 times:

- (1)) $\theta^s \sim p(\theta|y_1, \dots, y_{10})$
 (2)) $\tilde{y}_i^s \sim p(y|\theta^s), \forall i \in \{1, 2, \dots, 10\}$

For Fano factor, the p-value and graphical result are shown in 1. p-value is near 0.5, indicating that posterior predictive's fano number is similar to data.

Use lag $k = 1$ autocorrelation, defined as

$$r_1 = \frac{\sum_{i=1}^{N-1} (y_i - \bar{y})(y_{i+1} - \bar{y})}{\sum_{i=1}^{N-1} (y_i - \bar{y})^2}$$

The p-value and the graphical result are shown in 2.

We notice that $p(\text{autocorrelation}(\text{sample}) > \text{autocorrelation}(\text{data})) \approx 0.025$, which means that the assumption that year-to-year fatal accidents are independent is inadequate.

- (ii) Assumes that the numbers of fatal accidents in each year follows independent Poisson distribution with a constant rate and an exposure in each year proportional to the number of passenger miles flown. Denote y to be the number of fatal accidents, θ the rate, x the number of passenger miles flown. We have

$$y|\theta, x \sim \text{Poisson}(x\theta) = \frac{1}{y!} (x\theta)^y \exp(-x\theta)$$

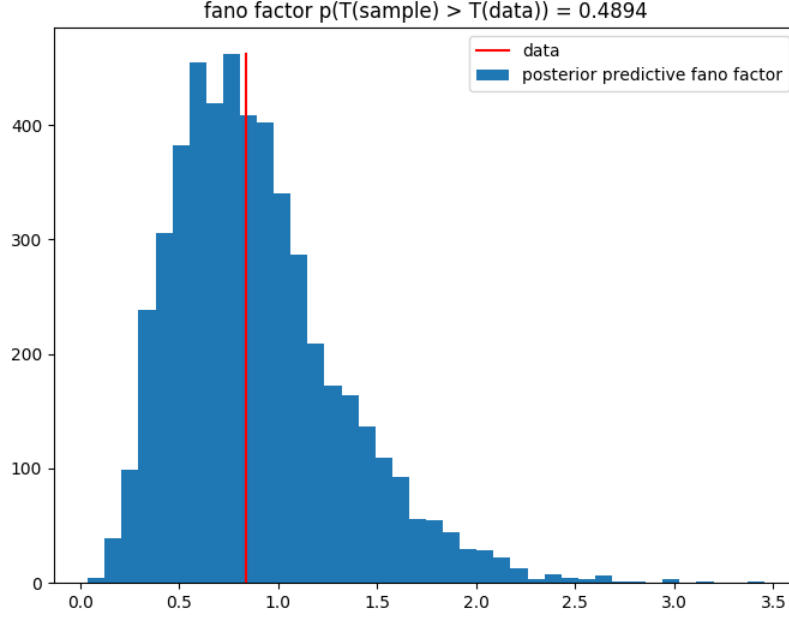


Figure 1: Posterior predictive Fano factor distribution for model 1

The likelihood of the data

$$p(y|\theta, x) = \prod_{i=1}^{10} \frac{1}{y_i!} (x_i \theta)^{y_i} \exp(-x_i \theta) \sim \theta^{\sum_{i=1}^{10} y_i} \exp\left(-\theta \sum_{i=1}^{10} x_i\right)$$

Similarly, if we choose prior to be $\text{Gamma}(0, 0)$, then the posterior distribution is

$$p(\theta|y) = \text{Gamma}(10\bar{y}, 10\bar{x})$$

where $\bar{y} = \frac{1}{10} \sum_{i=1}^{10} y_i$ and $\bar{x} = \frac{1}{10} \sum_{i=1}^{10} x_i$. Not until I did the sampling in python, did I realize that β is too large to have a reliable simulation. According to https://en.wikipedia.org/wiki/Gamma_distribution#Scaling, if

$$X \sim \text{Gamma}(\alpha, \beta),$$

then

$$cX \sim \text{Gamma}\left(\alpha, \frac{\beta}{c}\right).$$

The simulation results are shown in 3. We notice that this model better captured the autocorrelation within the data, but not the Fano factor. Our model significantly increases the ratio between variance and the mean. Notice in the simulation, I first sample a θ , and then use the same θ^s times the passenger miles flown to get the rate for Poisson Distribution. Hence, if there were more passenger miles flown, the model would have a greater fatal accident rate. However, this model is limited, not considering many more factors. For one thing, technology improves over the years too along with that the flights get more popular.

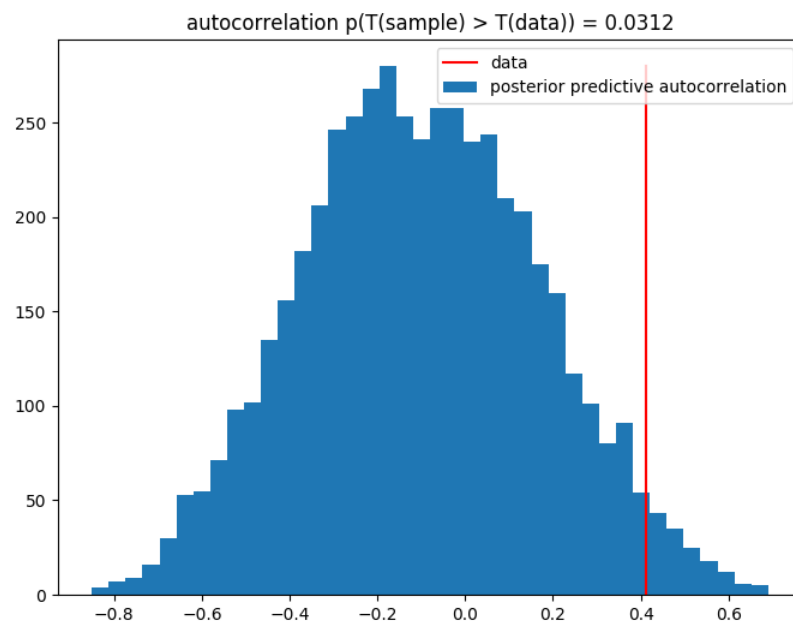


Figure 2: Posterior predictive autocorrelation distribution for model 1

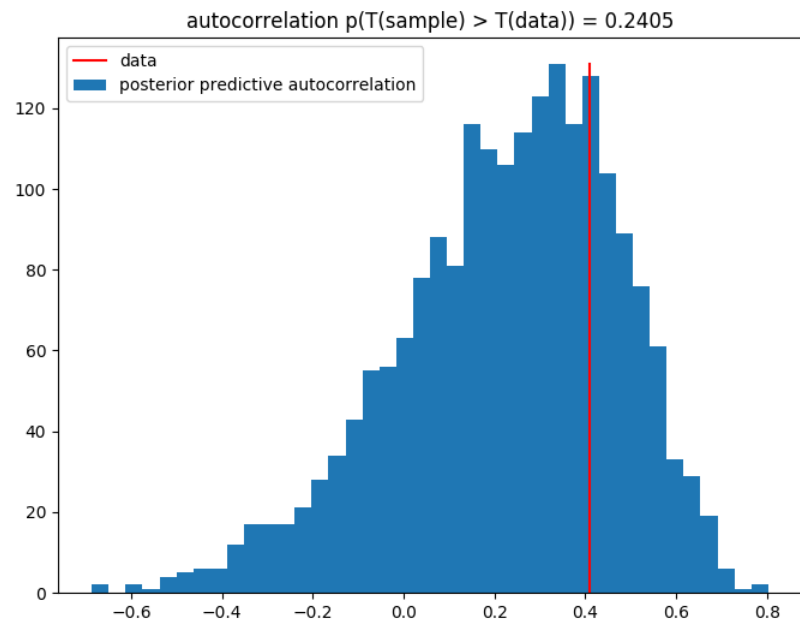
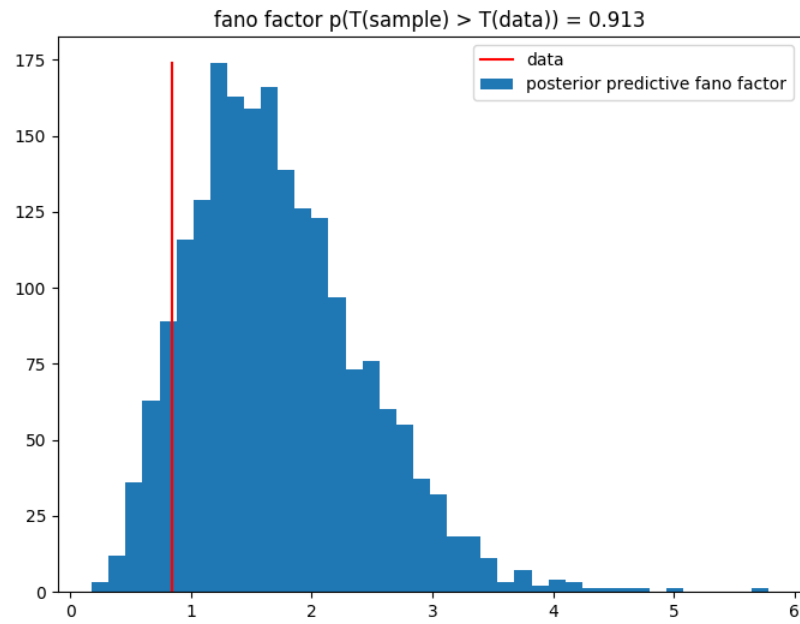


Figure 3: Posterior predictive Fano factor (top) and autocorrelation distribution (bottom) for model 2