

SDS 383D, Exercises 4: Hierarchical Models: Data-analysis Problems

Jan-Michael Cabrera

April 26, 2019

Price elasticity of demand

The data in “cheese.csv” are about sales volume, price, and advertising display activity for packages of Borden sliced “cheese.” The data are taken from Rossi, Allenby, and McCulloch’s textbook on *Bayesian Statistics and Marketing*. For each of 88 stores (store) in different US cities, we have repeated observations of the weekly sales volume (vol, in terms of packages sold), unit price (price), and whether the product was advertised with an in-store display during that week (disp = 1 for display).

Your goal is to estimate, on a store-by-store basis, the effect of display ads on the demand curve for cheese. A standard form of a demand curve in economics is of the form $Q = \alpha P^\beta$, where Q is quantity demanded (i.e. sales volume), P is price, and α and β are parameters to be estimated. You’ll notice that this is linear on a log-log scale,

$$\log P = \log \alpha + \beta \log Q$$

which you should feel free to assume here. Economists would refer to β as the price elasticity of demand (PED). Notice that on a log-log scale, the errors enter multiplicatively.

There are several things for you to consider in analyzing this data set.

1. The demand curve might shift (different α) and also change shape (different β) depending on whether there is a display ad or not in the store.
2. Different stores will have very different typical volumes, and your model should account for this.
3. Do different stores have different PEDs? If so, do you really want to estimate a separate, unrelated β for each store?
4. If there is an effect on the demand curve due to showing a display ad, does this effect differ store by store, or does it look relatively stable across stores?
5. Once you build the best model you can using the log-log specification, do see you any evidence of major model mis-fit?

Propose an appropriate hierarchical model that allows you to address these issues, and use Gibbs sampling to fit your model.

We begin by specifying the form of the hierarchical linear model:

$$\bar{y}_i = X_i^T \beta_i + e_i; \quad e_i \sim N(0, \sigma_i^2 I_{n_i})$$

where the index i corresponds to a particular store. The covariates for a particular store X_i^T is a matrix with n_i rows corresponding to the number of samples taken for store i and p columns:

$$X_i^T = \begin{bmatrix} 1 & \log P_{ij} & \delta_{ij} & \delta_{ij} \log P_{ij} \\ \vdots & \vdots & \vdots & \vdots \\ 1 & \log P_{in_i} & \delta_{in_i} & \delta_{in_i} \log P_{in_i} \end{bmatrix}$$

Here δ_{ij} corresponds to whether sample j for store i had a display or not.

Each β_i has a prior of the form,

$$\beta_i \sim MVN(\theta, V).$$

For store i the likelihood is multivariate normal with

$$p(\bar{y}_i | \beta_i, \sigma_i^2) \propto \exp \left[-\frac{1}{2} (\bar{y}_i - X_i^T \beta_i)^T \sigma_i^2 I_{n_i} (\bar{y}_i - X_i^T \beta_i) \right]$$

Here we show the prior for β_i explicitly

$$p(\beta_i | \theta, V) \propto \exp \left[-\frac{1}{2} (\beta_i - \theta)^T V^{-1} (\beta_i - \theta) \right].$$

We would like to find the conditional distribution for β_i given the data given by

$$p(\beta_i | \bar{y}_i, \sigma_i^2, \theta, V) \propto p(\beta_i | \theta, V) p(\bar{y}_i | \beta_i, \sigma_i^2).$$

Combining the likelihood and prior and completing the square we find that the posterior for β_i is also multivariate normal:

$$\beta_i | \bar{y}_i \dots \sim MVN \left(\frac{V^{-1} \theta + \frac{1}{\sigma_i^2} X_i^T \bar{y}_i}{V^{-1} + X_i^T \frac{1}{\sigma_i^2} X_i}, \left[V^{-1} + X_i^T \frac{1}{\sigma_i^2} X_i \right]^{-1} \right)$$

The prior for the scaling parameter is assumed to be an improper non-informative Jeffrey's prior:

$$p(\sigma_i^2) \propto \frac{1}{\sigma_i^2}.$$

The conditional distribution is then,

$$\begin{aligned} p(\sigma_i^2 | \bar{y}_i) &\propto p(\sigma_i^2) p(\bar{y}_i | \beta_i, \sigma_i^2) \\ &\propto \frac{1}{\sigma_i^2} \left(\frac{1}{\sigma_i^2} \right)^{n_i/2} \exp \left[-\frac{1}{2} (\bar{y}_i - X_i^T \beta_i)^T \sigma_i^2 I_{n_i} (\bar{y}_i - X_i^T \beta_i) \right]. \end{aligned}$$

The conditional distribution has the form of an inverse-gamma distribution with the following parameters:

$$\sigma_i^2 | \bar{y}_i \dots \sim IG \left(\frac{n_i}{2}, \frac{1}{2} (\bar{y}_i - X_i^T \beta_i)^T (\bar{y}_i - X_i^T \beta_i) \right).$$

The prior for θ is assumed to be flat and the conditional takes the form:

$$p(\theta|\beta) \propto \exp \left[-\frac{1}{2} \sum_{i=1}^s (\beta_i - \theta)^T V^{-1} (\beta_i - \theta) \right].$$

This reduces to a multivariate normal distribution of the form:

$$\theta|\beta \dots \sim MVN \left(\frac{1}{s} \sum_{i=1}^s \beta_i, \frac{1}{s} V \right)$$

The prior for the variance for the coefficients, β_i is assumed to have the form of an inverse-Wishart prior.

$$p(V) \propto |V|^{(d+p+1)/2} \exp \left[-\frac{1}{2} \text{tr}(CV^{-1}) \right]$$

The 'likelihood' for this prior comes from the coefficients and has the form,

$$p(\beta_i|\theta, V) \propto |\det V|^{1/2} \exp \left[-\frac{1}{2} (\beta_i - \theta)^T V^{-1} (\beta_i - \theta) \right],$$

noting that the term in the determinant is necessary for finding the conditional distribution for the variance. The full 'likelihood' is then,

$$p(\beta|\theta, V) \propto |\det V|^{s/2} \exp \left[-\frac{1}{2} \sum_{i=1}^s (\beta_i - \theta)^T V^{-1} (\beta_i - \theta) \right]$$

Combining this with the inverse-Wishart prior, we can show that the condition distribution is also inverse-Wishart:

$$\begin{aligned} p(V|\beta) &\propto |V|^{(d+p+1)/2} \exp \left[-\frac{1}{2} \text{tr}(CV^{-1}) \right] |\det V|^{s/2} \exp \left[-\frac{1}{2} \sum_{i=1}^s (\beta_i - \theta)^T V^{-1} (\beta_i - \theta) \right] \\ &\propto |V|^{(d+p+s+1)/2} \exp \left[-\frac{1}{2} \text{tr} \left(V^{-1} \left[C + \sum_{i=1}^s (\beta_i - \theta)^T (\beta_i - \theta) \right] \right) \right] \\ V|\beta \dots &\sim IW \left(d + s, C + \sum_{i=1}^s (\beta_i - \theta)^T (\beta_i - \theta) \right) \end{aligned}$$

Conditionals for Gibbs Sampling

$$\beta_i | \bar{y}_i \dots \sim MVN \left(\frac{V^{-1}\theta + \frac{1}{\sigma_i^2} X_i \bar{y}_i}{V^{-1} + X_i \frac{1}{\sigma_i^2} X_i^T}, \left[V^{-1} + X_i \frac{1}{\sigma_i^2} X_i^T \right]^{-1} \right)$$

$$\sigma_i^2 | \bar{y}_i \dots \sim IG \left(\frac{n_i}{2}, \frac{1}{2} (\bar{y}_i - X_i^T \beta_i)^T (\bar{y}_i - X_i^T \beta_i) \right)$$

$$\theta | \beta \dots \sim MVN \left(\frac{1}{s} \sum_{i=1}^s \beta_i, \frac{1}{s} V \right)$$

$$V | \beta \dots \sim IW \left(d + s, C + \sum_{i=1}^s (\beta_i - \theta)^T (\beta_i - \theta) \right)$$

A Gibbs sampler was written using Python 3.7 to sample from the derived conditional distributions (see Appendix and repo). The linear fits for all the stores is shown in figure 1. Here the black dots represent data from stores where a display was not present and the red dots represent data from stores where the display was present. The black and red lines correspond the fits for display not present and present respectively.

A noteworthy trend is that for many of the stores, having a display increased sales (i.e. volume) but also tended to increase price sensitivity (shown by the steeper slopes). The model also appears to do supprisingly well for stores where little or no data is available. Figures 2, 3, 4 show examples of these cases.

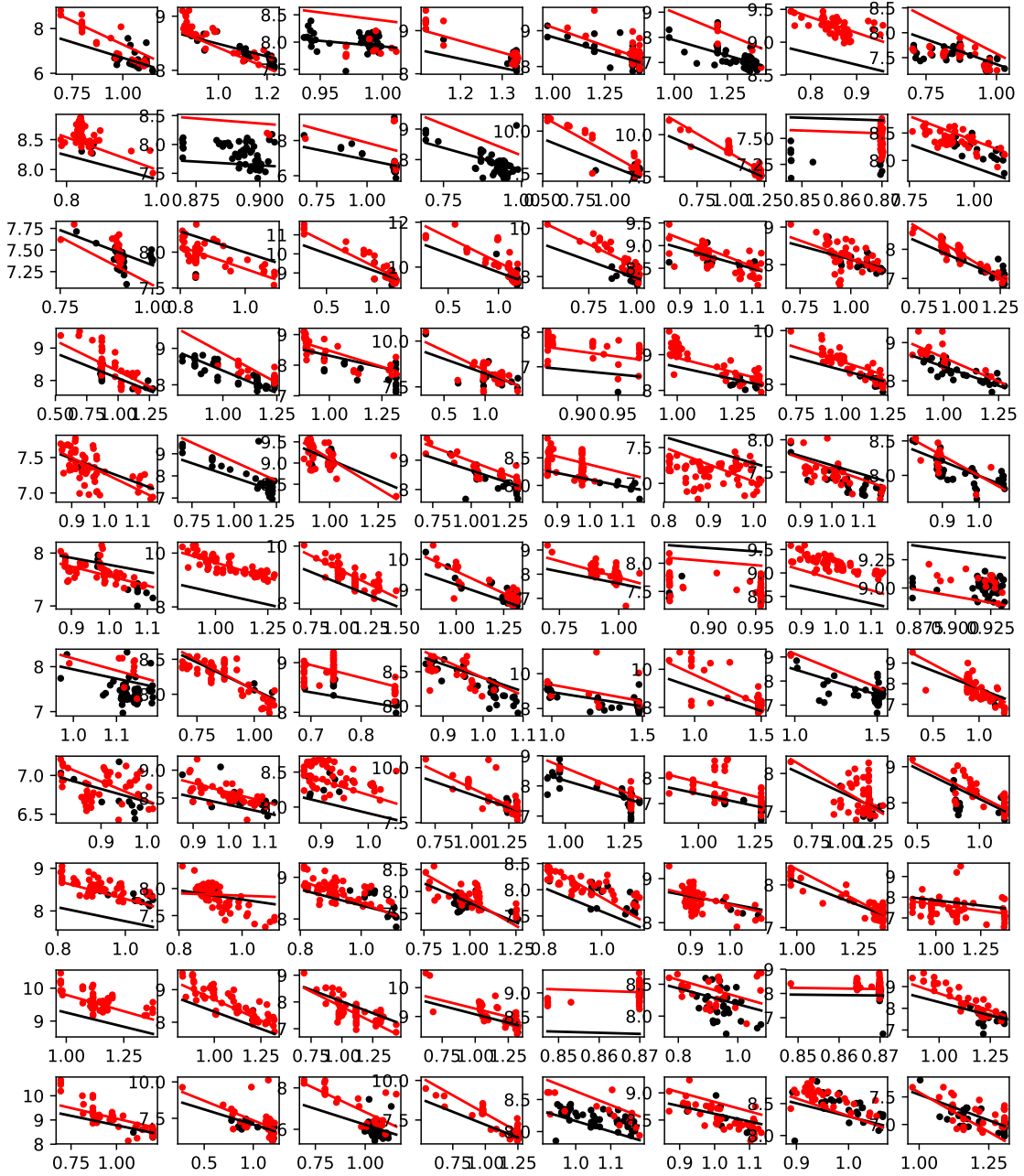


Figure 1: Data and fits of posterior means for all stores

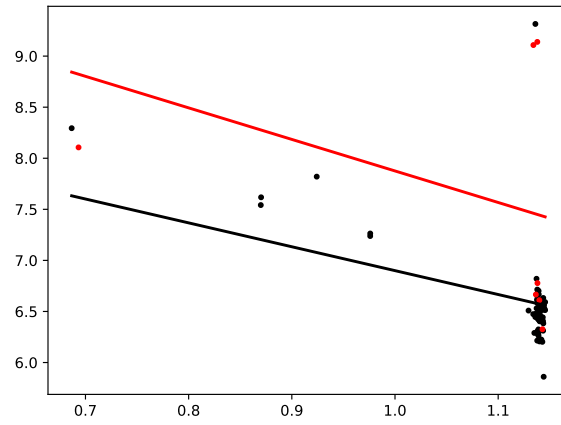


Figure 2: Data and fit for store 10

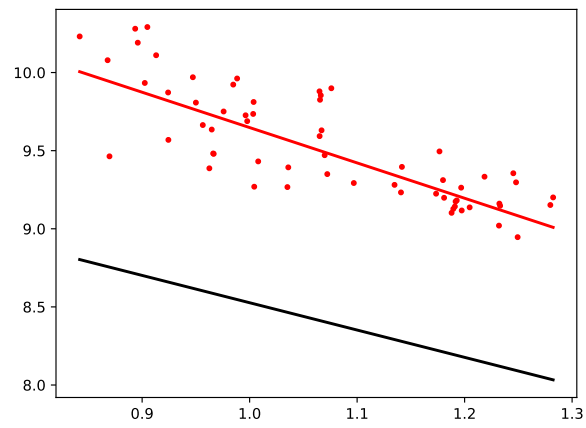


Figure 3: Data and fit for store 41

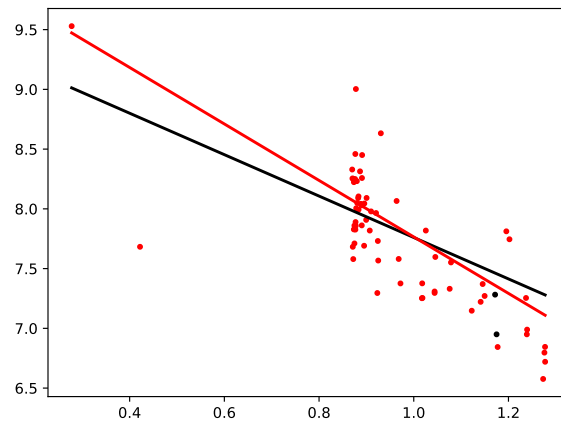


Figure 4: Data and fit for store 55

cheese.py

```
from __future__ import division
import numpy as np
import matplotlib.pyplot as plt
import pandas as pd
import sys
sys.path.append('.././scripts/')
import scipy.stats as stats
from numpy.linalg import inv

np.random.seed(3)

# Import data
df = pd.read_csv('.././data/cheese.csv', delimiter=',')

# Create an array of unique store names
stores = np.unique(df.store)

# Number of data rows
n = 5555
p = 4
d = 10

# Instantiate lists for storing sorted data
data = []
y = []
X = []

# Sort data by store
for s, store in enumerate(stores):
    data.append(df[df.store==store])

    # \bar{y}_i
    y.append(np.log(data[s].vol))
    data[s].vol = np.log(data[s].vol)
    data[s].price = np.log(data[s].price)

    # X_i^T
    X.append(np.array([np.ones(data[s].shape[0]), data[s].price, data[s].disp, data[s].price*data[s].disp]).T)

# Number of stores
s = 88

### Instantiate priors and traces
sigma_sq = 1

theta = np.zeros(p)
V = np.eye(p)*10**6

beta = np.zeros((s,p))
sigma = np.zeros(s)

beta_trace = []
sigma_trace = []
V_trace = []

#### Iterations
iterations = 5000
burn = 500
for j in range(iterations):
    # Store variable for inverse-Wishart
    B = 0
    # Iterate over stores to calculate \beta_i's
    for store in range(s):
```

```

n = len(y[store])
# beta_cov = [V^{-1} + X_i \frac{1}{\sigma_i^2} I_{n_i} X_i^T]^{-1}
beta_cov = inv(inv(V) + X[store].T @ ((1/sigma_sq)*np.eye(n)) @ X[store])

# beta_mean = beta_cov [V^{-1} \theta + \frac{1}{\sigma_i^2} X_i \bar{y}_i]
beta_mean = beta_cov @ (inv(V) @ theta + (1/sigma_sq) * (X[store] @ y[store]))

# Sample from multivariate normal for each store
beta[store] = stats.multivariate_normal.rvs(mean=beta_mean, cov=beta_cov)

# Shape and scale parameters for sigma
a = n/2
b = (y[store] - X[store] @ beta[store]) @ (y[store] - X[store] @ beta[store])/2

# Sample from inverse-gamma distribution
sigma[store] = stats.invgamma.rvs(a, 1/b)

# Sum variable for inverse-Wishart
B += np.tensordot((beta[store] - theta), (beta[store] - theta).T, axes=0)

# Sample from multivariate normal, theta_mean = 1/s \sum_{i=1}^s \beta_i, theta_cov = V/s
theta = stats.multivariate_normal.rvs(mean=(1/s)*beta.sum(axis=0), cov=V/s)

# Sample from inverse-Wishart distribution
V = stats.invwishart.rvs(d+s, np.eye(p) + B)

# Append samples to trace
beta_trace.append(beta)
sigma_trace.append(sigma)

# Reduce trace
beta_trace = np.asarray(beta_trace)
beta_trace_mean = np.mean(beta_trace[burn:], axis=0)

# Plot all stores together into one plot
fig, ax = plt.subplots(11, 8, figsize=(10,12))
fig.subplots_adjust(hspace=0.6)
n = 0
for i in range(11):
    for j in range(8):
        x_hat = X[n][X[n][:,1].argsort()]
        ax[i,j].plot(data[n][data[n].disp==0].price, data[n][data[n].disp==0].vol, '.k', linewidth=0.1)
        ax[i,j].plot(data[n][data[n].disp==1].price, data[n][data[n].disp==1].vol, '.r', linewidth=0.1)
        ax[i,j].plot(x_hat[:,1], beta_trace_mean[n][0]+beta_trace_mean[n][1]*x_hat[:,1], '-k')
        ax[i,j].plot(x_hat[:,1], (beta_trace_mean[n][0]+beta_trace_mean[n][2])+(beta_trace_mean[n][1]+beta_trace_mean[n][3])*x_hat[:,1], '-r')
        n+=1

plt.savefig('figures/cheese_all_plots.pdf')

# Plot selected stores
i = 10, 41, 55
for p, plots in enumerate(i):
    plt.figure()
    x_hat = X[plots][X[plots][:,1].argsort()]
    plt.plot(data[plots][data[plots].disp==0].price, data[plots][data[plots].disp==0].vol, '.k', linewidth=2)
    plt.plot(data[plots][data[plots].disp==1].price, data[plots][data[plots].disp==1].vol, '.r', linewidth=2)
    plt.plot(x_hat[:,1], beta_trace_mean[plots][0]+beta_trace_mean[plots][1]*x_hat[:,1], '-k', linewidth=2)
    plt.plot(x_hat[:,1], (beta_trace_mean[plots][0]+beta_trace_mean[plots][2])+(beta_trace_mean[plots][1]+beta_trace_mean[plots][3])*x_hat[:,1], '-r', linewidth=2)

plt.savefig('figures/cheese_plots_'+str(plots)+'.pdf')

```