

Bayesian Regression and Classification

Nikolas Tsigkas

April 2021

1 Bayesian Polynomial Regression on Linköping Temperature data

The given dataset contains daily temperature measurements in Linköping over one year. The aim is to fit a second degree polynomial to the data. More formally, the model can be formulated as:

$$T_i = \beta_0 + \beta_1 t_i + \beta_2 t_i^2 + \epsilon_i \quad \epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2) \quad i \in [1, \dots, 365]$$

Where T_i is the temperature at time i measured in $^{\circ}C$ and $t_i = \frac{i}{365}$ is the time of the year. The used prior for the β_i parameters, is the conjugate prior, i.e.:

$$\begin{aligned} \beta &\sim N(\mu_0, \sigma_0^2 \Omega_0^{-1}) \\ \sigma^2 &\sim \text{Scale-Inv-}\chi^2(\nu_0, \sigma_0^2) \end{aligned}$$

The first step is to choose suitable values for the hyperparameters μ (the mean vector), σ_0 (the scale parameter for the variance), ν_0 (the degrees of freedom for the variance) and Ω_0 (the precision matrix for the normal multivariate normal distribution). This can be done by drawing values from the prior of β and plotting the implied regression curves to see whether they agree with the data. After some iterations, the chosen suitable values were:

$$\mu = \begin{bmatrix} -15 \\ 120 \\ -110 \end{bmatrix}, \quad \sigma_0 = 1, \quad \nu_0 = 8, \quad \Omega_0 = \frac{1}{20} \mathbf{I}_3$$

The resulting plot for these prior values is shown below in figure 1, together with the observed temperature data, where one can see that the the curves cover the span of the data quite well.

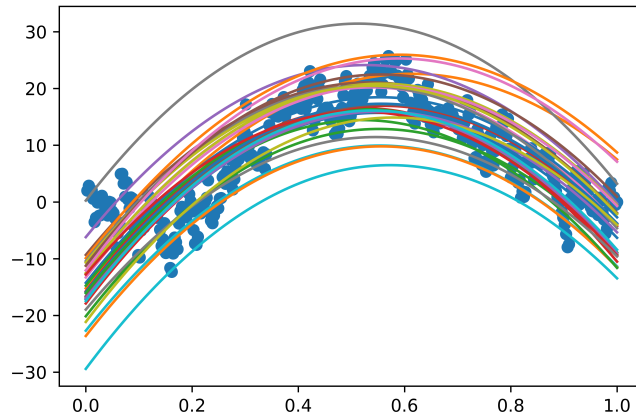


Figure 1: Regression simulated from the conjugate prior of β

Since the conjugate prior was chosen, the posterior can be evaluated analytically as:

$$\begin{aligned}\beta|\sigma, \mathbf{y}, \mathbf{X} &\sim N(\mu_n, \sigma^2 \Omega_n^{-1}) \\ \sigma^2 &\sim \text{Scale-Inv-}\chi^2(\nu_n, \sigma_n^2)\end{aligned}$$

Where:

$$\begin{aligned}\mu_n &= (\mathbf{X}^T \mathbf{X} + \Omega_0)^{-1} (\mathbf{X}^T \mathbf{X} \hat{\beta} + \Omega_0 \mu_0) \\ \Omega_n &= \mathbf{X}^T \mathbf{X} + \Omega_0 \\ \nu_n &= \nu_0 + n \\ \sigma_n^2 &= \frac{1}{\nu_n} (\mathbf{y}^T \mathbf{y} + \mu_0^T \Omega_0 \mu_0 + \mu_n^T \Omega_n \mu_n)\end{aligned}$$

And $\hat{\beta} = \mathbf{X}^\dagger \mathbf{y}$ is the ordinary least squares estimate¹.

From these expressions, one can now draw posterior samples to get a posterior model given the data. The histograms for each individual parameter are shown below in figure 2 after 10 000 simulated draws.

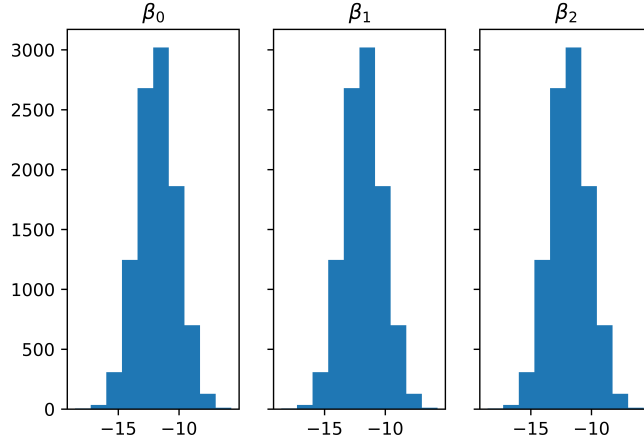


Figure 2: Histogram over the regression parameters generated from 10 000 simulated draws

From these 10 000 draws, 10 000 polynomial curves can be drawn as: $y^{(i)} = \mathbf{x}^T \beta^{(i)} + \sigma_i \xi_i$, where $\xi_i \stackrel{iid}{\sim} N(0, 1)$, $i = 1, \dots, 10000$, and $\mathbf{x} = [1, t, t^2]$ is the time values. After computing these implied temperature values from these regression curves for each day of the year, the 2.5th, 50th 97.5th percentiles were evaluated for each day. These results are shown in figure 3.

The interval between the two extreme percentiles capture the uncertainty quite well, with most of the values lying between the two. The model is however failing to capture some of the data around March and November.

From the simulations, one can also explore other features in the data, for which there might be no easy to derive probability density. One such example is the time of the year with the highest expected temperature \tilde{t} . From the model, this quantity can be found by evaluating the derivative of the temperature with regard to time, to find the unique stationary point which will yield the maximum value of the polynomial.

$$\frac{dT}{dt} = \beta_1 + 2\beta_2 \tilde{t} = 0 \implies \tilde{t} = -\frac{\beta_1}{\beta_2}$$

This random variable follows no known distribution (although the theory of quotient distributions may give some reasonable guidance for approximations), but can be approximated via simulation. Using the same random draws as before, the histogram in figure 4 was obtained.

¹† denotes the pseudo-inverse, i.e. $(\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$

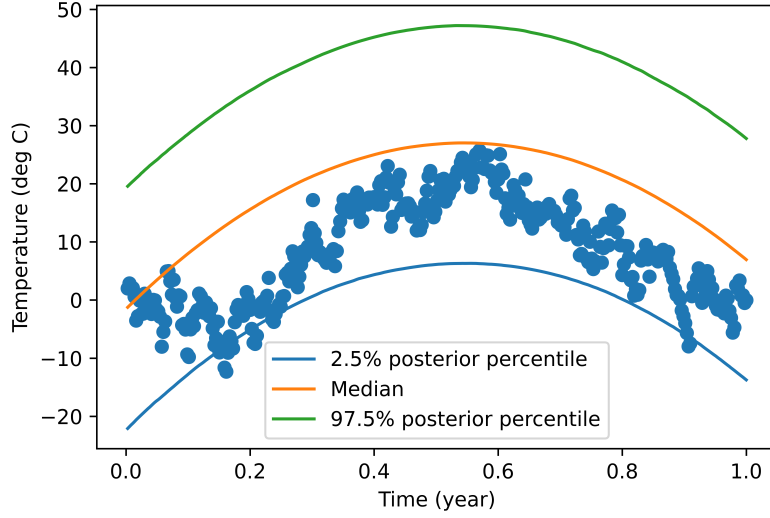


Figure 3: Posterior regression for the 2.5th, 50th and 97.5th percentiles

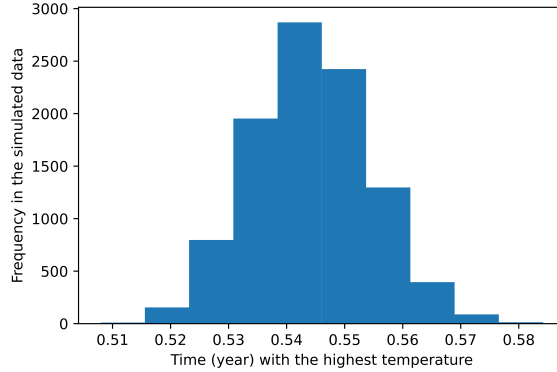


Figure 4: Histogram of the distribution of the time with the highest expected temperature

One can see that the time with the expected highest temperature is very centered around the values 0.5 – 0.6, corresponding to calendar dates between the 2nd of July and 7th of August.

This method of polynomial regression can be generalized to fit higher order polynomials, which could have been suitable given what was displayed in figure 3, where some of the data at the beginning of the year was not captured by the regression intervals. However, with such methods, one runs the risk of overfitting the model to the data which yields a model that does not generalize well. For such methods, regularization techniques such as *Lasso* or *Tikhonov Regularization* (also called *Ridge Regression*) may be suitable. These regularization methods are equivalent, in the Bayesian sense, to these prior distributions:

$$\beta_{i,Lasso}|\sigma \stackrel{iid}{\sim} \text{Laplace}(0, \frac{\sigma^2}{\lambda}) \quad \beta_{i,Ridge}|\sigma \stackrel{iid}{\sim} N(0, \frac{\sigma^2}{\lambda})$$

Where $\lambda \geq 0$ is the regularization hyperparameter, which can be interpreted as a penalty term that penalizes solutions whose β -terms have too high of a magnitude. Thus, increasing the value of λ will result in a more "smooth" polynomial. Too high values could however result in a model that underfits the data. One way of selecting the magnitude of λ is by treating it as a random variable, and letting it follow a prior distribution, such as the scaled inverse chi squared distribution.

```

1 """ Useful packages """
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from scipy import stats as st
6
7 np.random.seed(1234) # set seed for reproducibility
8
9 """ Helper functions for drawing samples and calculating posterior parameters
   """
10 # Generates num random draws from a scale-inverse-chi2 distribution
11 def draw_chi2inv(df, scale, num):
12     x = st.chi2.rvs(df-1, scale, size = num)
13     sigma_sq = (df-1)*scale/x
14     return sigma_sq
15
16 # Generates multivariate normal random draws for different sigmas
17 def draw_mvnorm(mu, Omega_inv, sigma_sq):
18     beta = np.zeros((len(sigma_sq),3))
19     for i, sig2 in enumerate(sigma_sq):
20         beta[i] = st.multivariate_normal.rvs(mean = mu, cov=sig2*Omega_inv,
21                                             size = 1)
22     return np.transpose(beta)
23
24 # Function that generates posterior parameters for the posterior
25 def post_params(X, y, mu_0, Omega_0, sigma_0):
26     n = len(y)
27
28     X_dagger = np.linalg.pinv(X)
29     XTX = X.T@X
30     beta_hat = X_dagger @ y
31
32     mu_n = np.linalg.inv(Omega_0 + XTX) @ XTX@beta_hat + Omega_0@mu_0
33     Omega_n = XTX + Omega_0
34     nu_n = nu_0 + n
35
36     yTy = y.T@y
37     mu_Omega_mu_0 = mu_0.T @ Omega_0 @ mu_0
38     mu_Omega_mu_n = mu_n.T @ Omega_0 @ mu_n
39
40     sigma_n2 = (nu_0*sigma_0**2 + yTy + mu_Omega_mu_0 + mu_Omega_mu_n )*1/nu_n
41
42     return [mu_n, Omega_n, nu_n, sigma_n2]
43
44 # Returns num random draws from the posterior of beta and sigma^2,
45 # given the posterior parameters as input
46 def post_draws(mu_n, Omega_n, nu_n, sigma_n2, num):
47     sigma_2 = draw_chi2inv(nu_n, sigma_n2, num)
48     beta_post = draw_mvnorm(mu_n, np.linalg.inv(Omega_n), sigma_2)
49     return [beta_post, sigma_2]
50
51 """ Import data """
52 data = pd.read_csv("TempLinkoping.txt", sep = "\t+")

```

```

52 Y = data.temp.values
53 # Data matrix consisting of 1, temp, temp^2
54 X = np.column_stack(((np.ones_like(data.time.values), data.time.values, data.
55     time.values**2)))
56
57 # Plot the data
58 plt.scatter(data.time.values, Y)
59 plt.xlabel("Time (year)")
60 plt.ylabel("Temperature (deg C)")
61 plt.show
62
63 # Prior parameters
64 mu_0 = np.array([-15, 120, -110])
65 Omega_0 = 0.05*np.eye(3)
66 Omega_0_inv = np.linalg.inv(Omega_0)
67 nu_0 = 8
68 sigma_0 = 1
69
70 # Draw 50 draws from the prior of beta
71 sigma_sq = draw_chi2inv(nu_0-1, sigma_0**2, 30)
72 beta_prior = draw_mvnorm(mu_0, Omega_0_inv, sigma_sq)
73
74 # Check visually if the priors generate a reasonable regression curve
75 zero_to_one = np.linspace(0, 1, 100)
76 X_test = np.column_stack((np.ones(100), zero_to_one, zero_to_one**2))
77
78 # By plotting the data
79 # together with the simulated regression curves from the prior
80 Y_prior = X_test@beta_prior
81 plt.plot(zero_to_one, Y_prior)
82 plt.scatter(data.time.values, Y)
83 plt.show
84 # Change the prior parameters until the results look reasonable
85
86 [mu_n, Omega_n, nu_n, sigma_n2] = post_params(X, Y, mu_0, Omega_0, sigma_0)
87 [beta_post, sigma2_post] = post_draws(mu_n, Omega_n, nu_n, sigma_n2, 10000)
88
89 # Histograms of the beta coefficients
90 fig, axs = plt.subplots(1, 3, sharey = True)
91
92 for i in range(3):
93     axs[i].hist(beta_post[0, :])
94
95 axs[0].set_title(r"$\beta_0$")
96 axs[1].set_title(r"$\beta_1$")
97 axs[2].set_title(r"$\beta_2$")
98
99 post_temp = np.zeros((365, 10000))
100 for i in range(10000):
101     post_temp[:, i] = X@beta_post[:, i] + np.sqrt(sigma2_post[i])*st.norm.rvs(1)
102
103 # Calculate the 2.5, 50 and 97.5 percentiles of the temperatures of each day
104 post_0025 = np.percentile(post_temp, 2.5, axis=1)

```

```

105 post_0500 = np.percentile(post_temp, 50, axis=1)
106 post_0975 = np.percentile(post_temp, 97.5, axis=1)
107
108 # Plot implied regression curves
109 plt.plot(X[:,1], post_0025)
110 plt.plot(X[:,1], post_0500)
111 plt.plot(X[:,1], post_0975)
112
113 plt.legend(["2.5% posterior percentile", "Median", "97.5% posterior percentile"
114            ])
115 plt.xlabel("Time (year)")
116 plt.ylabel("Temperature (deg C)")
117
118 plt.scatter(data.time.values, Y)
119 plt.show
120 plt.savefig("temp_data_posterior_curves.png", dpi = 1000)
121
122 # Find the simulated distribution of the time with maximal temperature
123 # from the posterior distribution of beta
124 time_max = -0.5*beta_post[1]/beta_post[2]
125
126 # Plot the distribution of time_max
127 plt.hist(time_max)
128 plt.xlabel("Time (year) with the highest temperature")
129 plt.ylabel("Frequency in the simulated data")

```

2 Logistic Regression on Women Worker Data

The given dataset contains 7 explanatory variables (the income of their husband, their years of education, their years of experience and the square of their years of experience, their age, and the number of small and large children respectively) on 200 women, and whether the women work or not. One way to predict whether a woman works or not, given their features, is by fitting a logistic regression model to the data.

Formally, the model is:

$$P[y_i = 1|\mathbf{x}_i] = \frac{1}{1 + e^{-\mathbf{x}_i^T \beta}} = \Lambda(\mathbf{x}_i^T \beta) \quad \mathbf{x}_i, \beta \in \mathbb{R}^8$$

Where $y_i = 1$ represents a working woman with features \mathbf{x}_i and $\Lambda(\cdot)$ is the cdf of the logistic distribution.

This model comes with the difficulty of choosing a suitable prior for the β -vector, as the posterior will not come in the form of a nice analytical expression from a known density kernel. However, by using a second order Taylor expansion of the log-posterior around its mode, an approximate posterior with a nice form can be derived.

$$\ln p(\beta|y) \approx \ln p(\tilde{\beta}|y) + (\beta - \tilde{\beta})^T \nabla_{\beta}(\ln p(\tilde{\beta}|y)) + \frac{1}{2}(\beta - \tilde{\beta})^T H_{\ln p}(\tilde{\beta})(\beta - \tilde{\beta})$$

From the definition of the mode, the first order term in the expression vanishes, as the function is maximized at this point. Using this, the above equation can be solved for the posterior density, resulting in.

$$\ln p(\beta|y) \approx \ln p(\tilde{\beta}|y) \exp \left[\frac{1}{2}(\beta - \tilde{\beta})^T H_{\ln p}(\tilde{\beta})(\beta - \tilde{\beta}) \right] \propto \exp \left[\frac{1}{2}(\beta - \tilde{\beta})^T H_{\ln p}(\tilde{\beta})(\beta - \tilde{\beta}) \right]$$

From this expression, it can be seen that the posterior is approximately proportional to the multivariate Gaussian kernel, with mean vector equal to the mode, and covariance matrix equal to the negative inverse hessian of the log-posterior (which can be interpreted as the observed information) i.e.:

$$\beta|y \stackrel{appr.}{\sim} N(\tilde{\beta}, -H_{\ln p}^{-1}(\tilde{\beta}))$$

The parameters $\tilde{\beta}$ and $-H_{\ln p}^{-1}(\tilde{\beta})$ can be evaluated numerically, by running a suitable optimization algorithm on the problem:

$$\max_{\beta} \ln p(\beta|y) \iff \max_{\beta} \ln \sum_1^n p(y_i|\beta) + \ln p(\beta)$$

The first term comes from the logistic function defined above, and the second from the prior. The used prior for was $\beta \sim N(0, 100I)$. As the objective of the optimization was to get estimates of both the mean and covariance matrix of the approximate posterior, a quasi-Newton optimization algorithm is suitable, as it uses an estimator for the hessian to find stationary points. The Broyden-Fletcher-Goldfarb-Shanno algorithm provided the following estimators:

$$\tilde{\beta} = [0.628 \quad -0.012 \quad 0.18 \quad 0.168 \quad -0.145 \quad -0.082 \quad -1.359 \quad -0.0247]^T$$

$$-H_{\ln p}^{-1}(\tilde{\beta}) = \begin{bmatrix} 0.169 & 0.007 & -0.004 & 0.024 & -0.071 & -0.004 & 0.072 & -0.176 \\ 0.007 & 0.001 & -0.001 & 0.001 & -0.003 & 0 & 0.004 & -0.007 \\ -0.004 & -0.001 & 0.005 & -0.001 & 0.005 & -0.001 & -0.015 & 0.004 \\ 0.024 & 0.001 & -0.001 & 0.008 & -0.025 & -0.001 & 0.009 & -0.025 \\ -0.071 & -0.003 & 0.005 & -0.025 & 0.082 & 0.002 & -0.027 & 0.074 \\ -0.004 & 0 & -0.001 & -0.001 & 0.002 & 0 & 0.001 & 0.004 \\ 0.072 & 0.004 & -0.015 & 0.009 & -0.027 & 0.001 & 0.174 & -0.076 \\ -0.176 & -0.007 & 0.004 & -0.025 & 0.074 & 0.004 & -0.076 & 0.183 \end{bmatrix}$$

From this, one can evaluate individual features. For instance, the coefficient corresponding to the number of small children β_6 follows the distribution $\beta_6 \sim N(-1.359, 0.417)$ (the second parameter here is the standard deviation). From this, the 95% equal tail credible interval was determined as: $[-1.9599, -0.54127]$. This indicates that the number of small children that a woman has, with a high probability, correlates negatively with the probability that she works.

One can also do inference on out of sample data. For instance, a 37 year old woman with:

1. a husband whose income is 13
2. 8 years of education
3. 11 years of work experience
4. 2 small children (under the age of 7)

can be represented, in the model, as:

$$\mathbf{x} = [1 \quad 13 \quad 8 \quad 11 \quad 1.21 \quad 37 \quad 2 \quad 0]^T$$

Given that the approximate posterior distribution of the coefficients is known, The distribution of the probability that she works can be evaluated through simulations, by drawing β -vectors, multiplying them with her data, and passing the product through the logistic cdf. The resulting histogram for this, when 10 000 draws from β was performed, is shown in figure 5.

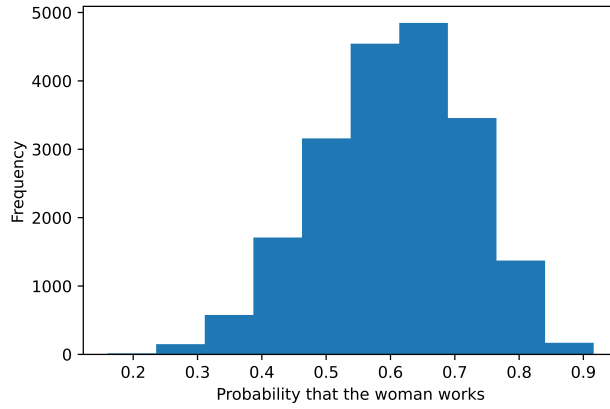


Figure 5: Histogram over the probability that the 37 year old woman works

Similarly, one can examine the distribution of the number of women working, given a population of eight women with the same features as the one mentioned above. This is equivalent of drawing samples from a $\text{Bin}(8, \Lambda(\mathbf{x}_i^T \beta))$, with $\beta \sim N(\tilde{\beta}, -H_{\ln p}^{-1}(\tilde{\beta}))$. The resulting histogram, for 10 000 draws from this distribution, is shown in figure 6.

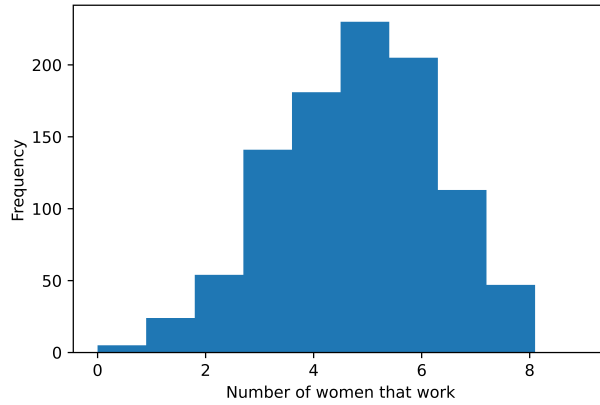


Figure 6: Histogram over the number of women out of the eight that work


```

1 """ Useful packages """
2 import pandas as pd
3 import numpy as np
4 import matplotlib.pyplot as plt
5 from scipy import stats as st
6 from scipy.optimize import minimize as fmin
7
8 np.random.seed(1234) # set seed for reproducibility
9
10 """ Helper functions """
11 # Returns the log posterior of the logistic function
12 def log_posterior(beta, y, X, mu, Cov):
13     log_prior = np.log( st.multivariate_normal.pdf(beta, mean=mu, cov=Cov) )
14     regr = X@beta
15     logL = np.sum( regr * y - np.log(1+np.exp(regr) ) )
16     return (logL + log_prior)
17
18 # Returns the posterior distribution of the probability that a particular sample
    works
19 def posterior_prob(x, mu, Cov, num):
20     beta = st.multivariate_normal.rvs(mean = mu, cov = Cov, size = num)
21     return st.logistic.cdf( np.matmul( x , np.transpose( beta ) ) )
22
23 # Returns the posterior distribution of the number of samples that work
24 def posterior_num(x, mu, Cov, num):
25     n = len(x)
26     prob = posterior_prob(x, mu, Cov, num)
27     mult_pred = np.zeros(num)
28     for i, p in enumerate(prob):
29         mult_pred[i] = st.binom.rvs(n, p)
30     return mult_pred
31
32 """ Import the data """
33 data = pd.read_csv("WomenWork.dat", sep = " ")
34
35 # Binary response variable
36 y = data.Work.values
37
38 # Covariates
39 X = data.iloc[:, 1:].values
40
41 # Prior parameters
42 mu_0 = np.zeros(8)
43 Cov_0 = 100*np.eye(8)
44
45 # Define the log posterior as a function of beta only, for the optimization
46 # Also, multiply with -1, as we seek the maximum
47 # and the algorithm performs minimization
48 def obj_fun(beta):
49     return -log_posterior(beta, y, X, mu_0, Cov_0)
50
51 # Find the max of the log posterior using a quasi-Newton method
52 # Which also returns a numerical approximation of  $H^{-1}$ 
53 x0 = 1*np.ones(8) # Initial guess
54 opt_res = fmin(obj_fun, x0, method = "BFGS", options={'disp' : True})
55
56 # The inverse of J, no minus sign
57 # as the objective function already has "switched signs"

```

```

58 J_inv = opt_res["hess_inv"]
59 # The mode, i.e. the variables at optimum
60 post_mode = opt_res["x"]
61
62 # Extract parameters for the coefficient of "NSmallChild",
63 # i.e. the second to last covariate
64 mu_SC = post_mode[-2]
65 sd_SC = J_inv[-2,-2]**0.5
66
67 eq_tail_CI_SC = [mu_SC-st.norm.ppf(0.925)*sd_SC, mu_SC+st.norm.ppf(0.975)*sd_SC]
68
69 """ Data for a woman with:
70 - a husband with income 13
71 - 8 years of education
72 - 11 years of work experience
73 - 37 years of age
74 - 2 small children (under the age of 7)
75 """
76 x_lw = np.array([1, 13, 8, 11, (11.0/10.0)**2, 37, 0, 0])
77
78 # The posterior probability that this particular woman works
79 post_prob = posterior_prob(x_lw, post_mode, J_inv, 20000)
80
81 plt.hist(post_prob)
82 plt.xlabel("Probability that the woman works")
83 plt.ylabel("Frequency")
84 plt.savefig("work_prob_hist.png", dpi = 1500)
85 plt.show
86
87 # The posterior distribution of the number of women working
88 post_num = posterior_num(x_lw, post_mode, J_inv, 1000)
89
90 plt.hist(post_num, range=(0,9))
91 plt.xlabel("Number of women that work")
92 plt.ylabel("Frequency")
93 plt.savefig("work_num_hist.png", dpi = 1500)
94 plt.show

```