

# CS 6190: Probabilistic Machine Learning Spring 2023

## Midterm Report

Handed out: 30 March, 2023

Due: 11:59pm, 05 May, 2023

### 1. Problem definition and motivation

The final term report is about earthquake prediction. I expect to use the dataset from Kaggle and knowledge from the class to predict the earthquake magnitude. Many countries are engaging in earthquake prediction. Especially at the beginning of this year, Turkey experienced several enormous disasters which caused many casualties. I believe there are relationships between magnitude and the dataset from the observation stops. Hence, this project aims to use machine learning to find the connection. The mid-term report also shows that the dataset and magnitude have a non-linear relationship, which motivates me to use machine learning skills to deal with the prediction. There are many ways to predict earthquakes, but none only use the dataset from the observation stops but not from seismicity.

### 2. My solution

In this project, I built three models to solve the problem. They are linear model (Built in the mid-term report), Gaussian distribution, and Dirichlet distribution. On another side, I used Laplace approximation to get the result to make things easy because Gaussian and Dirichlet are hard to compute the evidence. Then, use mean square error as my evaluation. I used the linear model as my baseline and tried to build other models to compete with it. Overall, Dirichlet has the best performance among competitors. (Mean square error: Dirichlet 2.75/ Gaussian 9.07/ linear 9.45)

In the linear model, I used Bayesian linear regression taught in the class. Equations are shown below:

$$p(\omega|D) \propto p(D|\omega)p(\omega)$$

$$\Rightarrow N(\omega|m_N, S_N) \propto N(t|\Phi\omega, \beta^{-1}I)N(\omega|m_0, S_0)$$

$$\text{Then, we can get } m_N = S_N(S_0^{-1}m_0 + \beta\Phi^T t) \text{ and } S_N^{-1} = S_0^{-1} + \beta\Phi^T\Phi$$

$$\text{Or, simplify it as } m_N = \beta S_N\Phi^T t \text{ and } S_N = \alpha I + \beta\Phi^T\Phi$$

The codes are shown in fig 1

```
40      """
41      SN = inv(inv(S0) + beta * Phi.T @ Phi)
42      mN = SN @ (inv(S0) @ m0 + beta * Phi @ t)
43      return SN, mN
```

Figure 1: Conjugate linear regression

In Gaussian distribution, I used factorized model taught in the class. Equations are shown below:

$$p(\alpha) = \text{Gam}(a_0, b_0)$$

$$\text{Prior: } p(\omega|\alpha) = N(\omega|0, \alpha I)$$

$$\text{Likelihood: } p(t|\omega) = \prod_{n=1}^N N(t_n|\omega\Phi_n, \beta^{-1})$$

$$\text{Joint: } p(t, \omega, \alpha) = p(t|\omega)p(\omega|\alpha)p(\alpha)$$

The task I want is to get the posterior  $p(\omega|t, \alpha)$ . The most significant difference from the linear model is that the gamma function influences the prior. I used Laplace approximation to get the features. Here are the equations:

$f(\theta) = \log(p(\theta, D))$   $f(\theta) \approx f(\theta_0) + \nabla f(\theta_0)^T(\theta - \theta_0) + \frac{1}{2}(\theta - \theta_0)^T A(\theta - \theta_0)$  Hence,  $p(\theta|D) = N(\theta|\theta_0, A)$   
In this Gaussian model, the logarithm joint is shown below:  
 $\log(p(w, t, \alpha)) = \log(p(t|\omega, \alpha)) + \log(p(\omega|\alpha)) + \log(p(\alpha))$   
 $\Rightarrow \log(\text{gamma}(\alpha)) + \frac{-1}{2}\omega\alpha I\omega^T + \frac{-1}{2}(t_n - \Phi_n\omega^T)\beta(t_n - \Phi_n\omega^T)$ . The codes are shown in fig 2

```

24 def compute_log_gamma(alpha):
25     """ Compute the expected value of every log pi_k, i.e., compute  $\sum_k E[\log \pi_k]$  """
26     #return digamma(alpha) - digamma(alpha.sum())
27     return np.log(gamma(alpha))
28 def compute_log_pw(alpha, w):
29     w = w.reshape(1,8)
30     #return -0.5 * w @ (alpha*np.eye(8)) @ w.T
31     #return np.diagonal(0.5 * (w @ (alpha*np.eye(8))).T @ w)
32     return 0.5 * w @ (alpha*np.eye(8)) @ w.T
33
34 def compute_log_pt_w(w, phi, t, beta):
35     return np.sum(0.5 * np.diagonal((t - phi @ w.T) * beta * ((t - phi @ w.T).T)))
36
37 def log_posterior(w, phi, t, alpha, beta):
38     #optimize w[0:w[i]]
39     #print("compute_log_gamma: ", compute_log_gamma(alpha))
40     #print("compute_log_pw: ", compute_log_pw(alpha, w).T)
41     #print("compute_log_pt_w: ", compute_log_pt_w(w, phi, t, beta))
42     result = -1 * compute_log_gamma(alpha) + compute_log_pw(alpha, w).T + compute_log_pt_w(w, phi, t, beta)
43     return np.sum(result)

```

Figure 2: Logarithm joint prob. in factorized Gaussian

Then, we need to find the argmax of  $\omega$ . Because it is just a derivative of the logarithm toward  $\omega$ , I use Scipy optimize minimize library to find the mode.(shown in fig 3)

```

84 solution = scipy.optimize.minimize(log_posterior, w0_laplace, args=(phi_laplace, tn_laplace, alpha, beta,0),method='BFGS')
85 w_map = solution.x
86 print("w_map")
87 print(w_map)
88

```

Figure 3: Find Gaussian logarithm joint mode

In Dirichlet distribution, I changed the prior from Gaussian distribution to Dirichlet distribution. The main goal is to see the performance difference from Gaussian.

Prior:  $Dir(\omega, \alpha) = \frac{\Gamma(\alpha_0)}{\Gamma(\alpha_1)\Gamma(\alpha_2)\dots\Gamma(\alpha_k)} \prod_{k=1}^K \omega_k^{\alpha_k-1}$

Likelihood:  $p(t|\omega) = \prod_{n=1}^N N(t_n|\omega\Phi_n, \beta^{-1})$ .

The logarithm joint:  $\log(p(t|\omega)) + \log(Dir(\omega, \alpha))$ . I divided it into a couple of items:

$\log(p(t|\omega)) = \frac{-1}{2}(t_n - \Phi_n\omega^T)\beta(t_n - \Phi_n\omega^T)$  and

$\log(Dir(\omega, \alpha)) = \log(\Gamma(\alpha_0)) - \log(\Gamma(\alpha_1)) - \log(\Gamma(\alpha_2))\dots - \log(\Gamma(\alpha_k)) + (\alpha_0 - 1)\log(\omega_0) + (\alpha_1 - 1)\log(\omega_1) + \dots(\alpha_k - 1)\log(\omega_k)$

The codes are shown in fig 4 (prior) and fig 5 (likelihood). Afterward, I use Scipy library to get the mode of the logarithm.(shown in fig 6)

### 3. Experimental evaluation

To evaluate the model performance, I use mean square error on the test dataset to confirm how far the model is from the ground truth. The result is shown in fig 7. The result, linear model(9.44), Gaussian model(9.07), and Dirchlet model(2.74), shows how the different models influence performance. For example, the linear and Gaussian models are similar, except the prior is limited by another distribution or not. With the restriction, the performance starts improving (from 9.44 to 9.07). On the other side, we can tell how important a model is in regression. By change the prior from Gaussian to Dirichlet, the performance improves enormously (from 9.07 to 2.74). Though there still is some improvement space (2.74 is too large), I believe trying more distributions can improve the performance to an acceptable score.

### 4. Future plan

We already saw how much important a model influences performance. The next step would be to find the most proper model(try other models, e.g., Poisson or Wishart). Then, using more accurate computing approximation ways to get the result (e.g., Variational inference).

The source code is put on the GitHub: [link](#)

```

2 def compute_log_pw2(alpha, w): #dirichlet prior
3     #left side (not log yet): gamma(a0)/(gamma(a1)*gamma(a2)...(gamma(ak))
4     #right side (not log yet): w0^(a0-1) * w1^(a1-1) * w2^(a2-1) ... * wk^(ak-1)
5
6     #left side log:
7     log_gamma_alpha = np.log(gamma(alpha))
8     log_gamma_alpha[0] = -1 * log_gamma_alpha[0]
9     log_gamma_alpha = -1 * log_gamma_alpha
10    leftside = log_gamma_alpha
11
12    #right side:
13    alpha_minus_1 = (alpha -1).reshape(1,8)
14    #print("alpha_minus_1 shape: ",alpha_minus_1.shape)
15    w_normalize = w/np.sum(w)
16    #print("w: ",w)
17    #print("w_normalize: ",w_normalize)
18    log_w = np.log(w_normalize)
19    #print("log_w shape: ",log_w.shape)
20
21    rightside = alpha_minus_1 @ log_w.T
22    result = leftside + rightside
23    return np.sum(result)
24

```

Figure 4: Dirichlet prior

```

24
25 def compute_log_pt_w2(w, phi, t, beta):
26     return np.sum(0.5 * np.diagonal((t- phi @ w.T) * beta * ((t- phi @ w.T)).T))
27
28 def log_posterior2(w, phi, t, alpha, beta):
29     #print("compute_log_pw2: ",compute_log_pw2(alpha, w))
30     #print("compute_log_pt_w2: ",compute_log_pt_w2(w, phi, t, beta))
31     result = compute_log_pw2(alpha, w) + compute_log_pt_w2(w, phi, t, beta)
32     return result
33

```

Figure 5: Dirichlet likelihood and posterior

```

34
35 solution = scipy.optimize.minimize(log_posterior2, w0_laplace, args=(phi_laplace, tn_laplace, alpha, beta),method='BFGS')
36 w_map = solution.x
37 print("w_map")
38 print(w_map)

```

Figure 6: Find logarithm Dirichlet distribution mode

```

1 #validation
2 def Get_predic(coef,input):
3     input_numpy=input.to_numpy()
4     output=np.arange(len(input_numpy), dtype=float)
5     for idx in range(len(input_numpy)):
6         output[idx] = f(input_numpy[idx], coef)
7     return output
8 test_prediction = testY
9 test_prediction=Get_predic(mN,testX)
10 print("Baysian linear Mean Square Error:")
11 print(np.square(np.subtract(testY,test_prediction)).mean())
12 #validation end

```

Figure 7: Evaluation

## 5. Reference to literature

- Laplace approximation: [https://james-brennan.github.io/posts/laplace\\_approximation/](https://james-brennan.github.io/posts/laplace_approximation/)
- Kaggle data set: <https://www.kaggle.com/datasets/grigol1/earthquakes-2000-2023>
- Columans meaning: <https://earthquake.usgs.gov/data/comcat/data-eventterms.php>
- Baysien linear learning: <https://www.microsoft.com/en-us/research/uploads/prod/2006/01/Bishop-Pattern-Recognition-and-Machine-Learning-2006.pdf>
- Regular linear learning: [https://scikit-learn.org/stable/modules/generated/sklearn.linear\\_model.LinearRegression.html](https://scikit-learn.org/stable/modules/generated/sklearn.linear_model.LinearRegression.html)