

DS-GA 1003 HW5

Yiyan Chen

April 2018

1 Introduction

2 From Scores to Conditional Probabilities

2.1

$$\begin{aligned} E_y[l(y(f(x))|x)] &= \sum l(yf(x))p(y|x) \\ &= l(f(x))p(y=1|x) + l(-f(x))p(y=-1|x) \\ &= l(f(x))\pi(x) + l(-f(x))(1-\pi(x)) \end{aligned}$$

2.2

$$\begin{aligned} E_y[l(y(f(x))|x)] &= e^{-f(x)}\pi(x) + e^{f(x)}(1-\pi(x)) \\ &= e^{-\hat{y}}\pi(x) + e^{\hat{y}}(1-\pi(x)) \\ \frac{\partial E_y}{\partial \hat{y}} &= -e^{-\hat{y}}\pi(x) + e^{\hat{y}}(1-\pi(x)) = 0 \\ e^{\hat{y}}(1-\pi(x)) &= e^{-\hat{y}}\pi(x) \\ e^{2\hat{y}} &= \frac{\pi(x)}{1-\pi(x)} \\ \hat{y}^* &= \frac{1}{2} \ln \frac{\pi(x)}{1-\pi(x)} \end{aligned}$$

$$\begin{aligned}
e^{f^*(x)} &= \frac{\pi(x)}{1 - \pi(x)} \\
(1 - \pi(x))e^{2f^*(x)} &= \pi(x) \\
e^{2f^*(x)} - \pi(x)e^{2f^*(x)} &= \pi(x) \\
(1 + e^{2f^*(x)})\pi(x) &= e^{2f^*(x)} \\
\pi(x) &= \frac{e^{2f^*(x)}}{1 + e^{2f^*(x)}} \\
&= \frac{1}{e^{-2f^*(x)} + 1}
\end{aligned}$$

2.3

$$\begin{aligned}
E_y &= \ln(1 + e^{-f(x)})\pi(x) + \ln(1 + e^{f(x)})(1 - \pi(x)) \\
\frac{\partial E_y}{\partial \hat{y}} &= \frac{-e^{-\hat{y}}}{1 + e^{-\hat{y}}}\pi(x) + \frac{e^{\hat{y}}}{1 + e^{\hat{y}}}(1 - \pi(x)) \\
&= \frac{-1}{e^{\hat{y}} + 1}\pi(x) + \frac{e^{\hat{y}}}{1 + e^{\hat{y}}}(1 - \pi(x)) = 0 \\
-\pi(x) + e^{\hat{y}}(1 - \pi(x)) &= 0 \\
e^{\hat{y}}(1 - \pi(x)) &= \pi(x) \\
e^{\hat{y}} &= \frac{\pi(x)}{1 - \pi(x)} \\
\hat{y}^* &= \ln \frac{\pi(x)}{1 - \pi(x)}
\end{aligned}$$

$$\begin{aligned}
e^{f^*(x)} &= \frac{\pi(x)}{1 - \pi(x)} \\
(1 - \pi(x))e^{f^*(x)} &= \pi(x) \\
(1 + e^{f^*(x)})\pi(x) &= e^{f^*(x)} \\
\pi(x) &= \frac{e^{f^*(x)}}{1 + e^{f^*(x)}} \\
\pi(x) &= \frac{1}{e^{-f^*(x)} + 1}
\end{aligned}$$

2.4 [Optional]

$[-1, 1] :$

$$\begin{aligned} E_y &= (1 - f(x))\pi(x) + (1 + f(x))(1 - \pi(x)) \\ &= (1 - 2\pi(x))f(x) + 1 \end{aligned}$$

$$\begin{cases} 0 < \pi(x) < \frac{1}{2}, & f^*(x) = 1 \\ \frac{1}{2} < \pi(x) < 1, & f^*(x) = -1 \end{cases} \quad (1)$$

$$\begin{cases} \pi(x) - \frac{1}{2} < 0, & f^*(x) = 1 \\ \pi(x) - \frac{1}{2} > 0, & f^*(x) = -1 \end{cases} \quad (2)$$

$f(x) \in (1, \infty) :$

$E_y = (1 + f(X))(1 - \pi(x))$ The minimizer doesn't exist

$f(x) \in (-\infty, -1) :$

$E_y = (1 - f(X))\pi(x)$ The minimizer doesn't exist

$$f^*(x) = \text{sign}(\pi(x) - \frac{1}{2})$$

3 Logistic Regression

3.1 Equivalence of ERM and probabilistic approaches

When $i \in P, y_i = 1, y'_i = 1$

When $i \in N, y_i = -1, y'_i = 0$

$$\begin{aligned} -\log \phi(w^T x_i) &= -\log \frac{1}{1 + e^{-w^T x_i}} \\ &= \log(1 + e^{-w^T x_i}) \\ -\log(1 - \phi(w^T x_i)) &= \log(1 - \frac{1}{1 + e^{-w^T x_i}}) \\ &= \log(1 - \frac{1}{1 + e^{-w^T x_i}})^{-1} \\ &= \log \frac{e^{-w^T x_i} + 1}{e^{-w^T x_i}} \\ &= \log(1 + e^{w^T x_i}) \end{aligned}$$

$$n\hat{R}_n(w) = \sum_{i \in P} \log(1 + e^{-w^T x_i}) + \sum_{i \in N} \log(1 + e^{w^T x_i})$$

$$\begin{aligned} NLL(w) &= \sum_{i \in P} [-\log \phi(w^T x_i)] + \sum_{i \in N} -\log(1 - \phi(w^T x_i)) \\ &= \sum_{i \in P} \log(1 + e^{-w^T x_i}) + \sum_{i \in N} \log(1 + e^{w^T x_i}) \\ &= n\hat{R}_n(w) \end{aligned}$$

3.2 Numerical Overflow and the log-sum-exp trick

3.2.1

$$\begin{aligned}\log(e^{x_1} + \dots + e^{x_n}) &= \log(e^{x^*}(e^{x_1-x^*} + \dots + e^{x_n-x^*})) \\ &= x^* + \log(e^{x_1-x^*} + \dots + e^{x_n-x^*})\end{aligned}$$

3.2.2

$x_i - x^* \in (-\infty, 0]$, thus, $\exp(x_i - x^*) \in (0, 1]$

3.2.3

Since at least one $x_i - x^*$ is zero when $x_i = x^*$, $e^{x_1-x^*} + \dots + e^{x_n-x^*}$ will always be bigger than 1. So the $\log(e^{x_1-x^*} + \dots + e^{x_n-x^*})$ will never go to negative infinity.

3.2.4

```
part1 = 0
part2 = -s
np.logaddexp(part1, part2)
```

3.3 Regularized Logistic Regression

3.3.1

According to Log-sum-exp rule, $\log(1 + \exp(-y_i w^T x_i)) = \log(\exp(0) + \exp(-y_i w^T x_i))$ is convex. And the sum of all convex functions is convex as well. Also, norm is convex. $\frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i w^T x_i))$ is convex, and $\lambda \|w\|^2$ is convex too. So, $J_{\text{logistic}(w)}$ is convex.

3.3.2

```
1 def f_objective(theta, X, y, l2_param=1):
2     '''
3     Args:
4         theta: 1D numpy array of size num_features
5         X: 2D numpy array of size (num_instances, num_features)
6         y: 1D numpy array of size num_instances
7         l2_param: regularization parameter
8
9     Returns:
10         objective: scalar value of objective function
11     '''
12     n = X.shape[0]
13     #x_max = np.max(np.dot(-y, np.dot(X, theta)))
14     summ = 0
15     for i in range(n):
16         summ += np.logaddexp(0, -y[i]*np.dot(theta, X[i]))
17     return summ/n + l2_param*np.sum(theta**2)
```

3.3.3

```
1 def fit_logistic_reg(X, y, objective_function, l2_param=1):
2     '''
3     Args:
4         X: 2D numpy array of size (num_instances, num_features)
5         y: 1D numpy array of size num_instances
6         objective_function: function returning the value of the objective
7         l2_param: regularization parameter
8
9     Returns:
10         optimal_theta: 1D numpy array of size num_features
11     '''
12     theta = np.zeros(X.shape[1])
13     optimal_theta = minimize(objective_function, theta, args=(X,y,l2_param)).x
```

```

14     return optimal_theta
15
16 def load_file(file_name):
17     f_myfile = open(file_name, 'r')
18     d = StringIO(f_myfile.read())
19     data = np.loadtxt(d, delimiter=",")
20     f_myfile.close()
21     return data
22
23 X_train = load_file("/Users/cyian/Desktop/NYU/SPRING2018/DS-GA1003/hw5-probabilistic/logistic")
24 X_val = load_file("/Users/cyian/Desktop/NYU/SPRING2018/DS-GA1003/hw5-probabilistic/logistic")
25 y_train = load_file("/Users/cyian/Desktop/NYU/SPRING2018/DS-GA1003/hw5-probabilistic/logistic")
26 y_val = load_file("/Users/cyian/Desktop/NYU/SPRING2018/DS-GA1003/hw5-probabilistic/logistic")
27 train_bias = np.ones((X_train.shape[0], 1))
28 val_bias = np.ones((X_val.shape[0], 1))
29 X_train_bias = np.hstack((X_train, train_bias))
30 X_val_bias = np.hstack((X_val, val_bias))
31 X_train_bias_std = preprocessing.scale(X_train_bias)
32 X_val_bias_std = preprocessing.scale(X_val_bias)
33
34 w_1 = fit_logistic_reg(X_train_bias_std, y_train_j, f_objective)

```

```

array([[ 9.55682759e-04, -2.98411854e-04,  3.02812767e-03,
        1.05326700e-01, -3.58837262e-03, -1.35879681e-03,
       -3.85259502e-03, -7.90123362e-04, -1.14392118e-03,
       -7.17819733e-02,  6.54800235e-03, -4.51121904e-03,
        1.12491086e-02, -3.86491439e-03, -2.71224635e-03,
        1.50343327e-03, -2.78428667e-03, -9.19058606e-03,
       -6.82319847e-03, -1.02758826e-02, -1.32052354e-08])

```

3.3.4

```

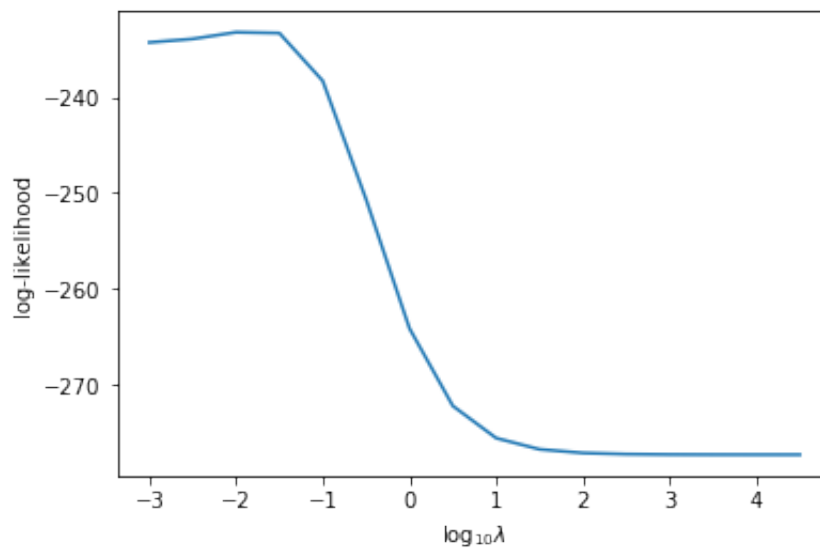
1 def get_loglikelihood(X, y, w, l2_param):
2     n = X.shape[0]
3     obj_val = f_objective(w, X, y, l2_param)
4     return -n*(obj_val - l2_param*np.sum(w**2))
5 l2_val = np.arange(-3, 5, 0.5, dtype=float)
6 loglike = []
7 for i in l2_val:
8     w = fit_logistic_reg(X_train_bias_std, y_train_j, f_objective, 10**i)
9     loglike.append(get_loglikelihood(X_val_bias_std, y_val_j, w, 10**i))
10
11 import matplotlib.pyplot as plt
12 plt.plot(l2_val, loglike)
13 plt.xlabel('$\log_{10}\lambda$')

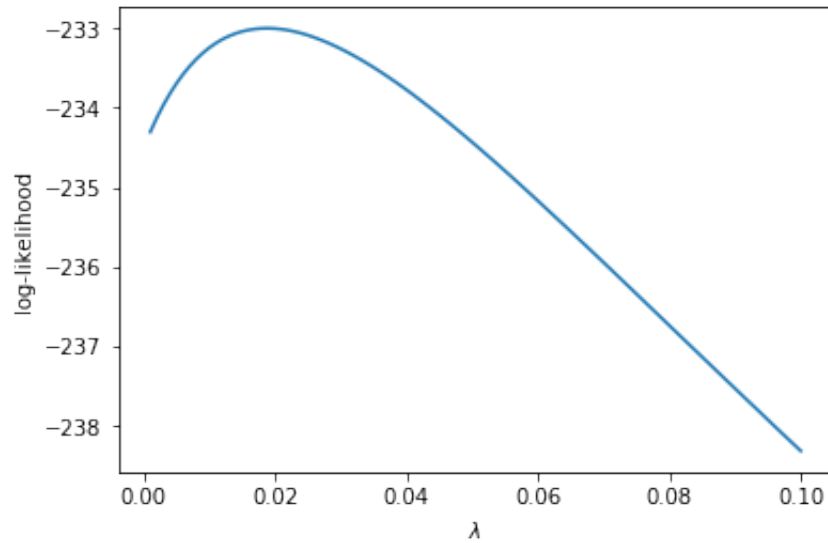
```

```

14 plt.ylabel("log-likelihood")
15 plt.show()
16
17 l2_val2 = np.linspace(0.001, 0.1, 100)
18 loglike2 = []
19 for i in l2_val2:
20     w = fit_logistic_reg(X_train_bias_std, y_train_j, f_objective, i)
21     loglike2.append(get_loglikelihood(X_val_bias_std, y_val_j, w, i))
22 plt.plot(l2_val2, loglike2)
23 plt.xlabel('$\lambda$')
24 plt.ylabel("log-likelihood")
25 plt.show()
26
27 from operator import itemgetter
28 zipped = zip(l2_val2, loglike2)
29 best_lambda = max(zipped, key = itemgetter(1))[0]
30 best_lambda

```





The l2 regularization parameter is 0.019000000000000003.

3.3.5

```

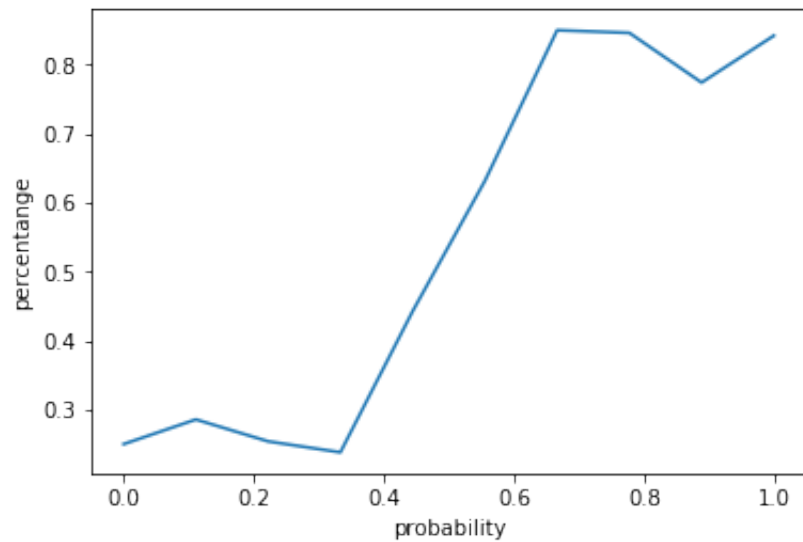
1 def slice_tuple(min_val, max_val, lst):
2     pos = []
3     for i in lst:
4         if i[1] > min_val and i[1] <= max_val:
5             pos.append(i)
6     return pos
7
8 pos1 = slice_tuple(0,0.1,lst)
9 pos2 = slice_tuple(0.1,0.2,lst)
10 pos3 = slice_tuple(0.2,0.3,lst)
11 pos4 = slice_tuple(0.3,0.4,lst)
12 pos5 = slice_tuple(0.4,0.5,lst)
13 pos6 = slice_tuple(0.5,0.6,lst)
14 pos7 = slice_tuple(0.6,0.7,lst)
15 pos8 = slice_tuple(0.7,0.8,lst)
16 pos9 = slice_tuple(0.8,0.9,lst)
17 pos10 = slice_tuple(0.9,1,lst)
18
19 def percent_tuple(lst):
20     n = len(lst)
21     count = 0
22     for i in lst:
23         if i[0] == 1:

```

```

24         count +=1
25     return count/n
26
27 per1 = percent_tuple(pos1)
28 per2 = percent_tuple(pos2)
29 per3 = percent_tuple(pos3)
30 per4 = percent_tuple(pos4)
31 per5 = percent_tuple(pos5)
32 per6 = percent_tuple(pos6)
33 per7 = percent_tuple(pos7)
34 per8 = percent_tuple(pos8)
35 per9 = percent_tuple(pos9)
36 per10 = percent_tuple(pos10)
37 percent = [per1, per2, per3, per4, per5, per6, per7, per8, per9, per10]
38 plt.plot(np.linspace(0,1,10), percent)
39 plt.xlabel("probability")
40 plt.ylabel("percentage")
41 plt.show()

```



Summa-
rize the result: When there are 10 bins, the percentage rate with respect to probability is very close to $y=x$.

4 Bayesian Logistic Regression with Gaussian Priors

4.1

$$\begin{aligned}
 P(w|D') &= \frac{P(D'|w)P(w)}{P(D')} \\
 &= k \cdot P(D'|w)P(w) \\
 &= k \cdot \exp(-NLL_{D'}(w)) \cdot P(w)
 \end{aligned}$$

4.2

$$\begin{aligned}
 -\log(p(w|D')) &= -\log k + NLL_{D'}(w) - \log P(w) \\
 &= NLL_{D'}(w) - \log P(w) + c \\
 &= -\log((2\pi^k|\Sigma|)^{\frac{1}{2}} \exp(-\frac{1}{2}w^T \Sigma^{-1}w)) + n\hat{R}_n(w) + c \\
 &= -\frac{1}{2} \log(2\pi^k|\Sigma|) + \frac{1}{2}w^T \Sigma^{-1}w + n\hat{R}_n(w) + c \\
 &= \frac{1}{2}w^T \Sigma^{-1}w + n\hat{R}_n(w) + c'
 \end{aligned}$$

In order to make MAP estimator is same as the regularized logistic regression:

$$\min(\frac{1}{2}w^T \Sigma^{-1}w + n\hat{R}_n(w)) \sim \min(n\hat{R}_n(w) + \lambda n||w||^2)$$

Then the covariance matrix is I

4.3

$$\begin{aligned}
 \frac{1}{2}w^T w + n\hat{R}_n(w) &= n\hat{R}_n(w) + \lambda n||w||^2 \\
 \frac{1}{2}||w||^2 + n\hat{R}_n(w) &= n\hat{R}_n(w) + \lambda n||w||^2 \\
 \frac{1}{2} &= \lambda n \\
 \lambda &= \frac{1}{2n}
 \end{aligned}$$

5 Bayesian Linear Regression - Implementation

5.1

```
1 def likelihood_func(w, X, y_train, likelihood_var):
2     '''
3     Implement likelihood_func. This function returns the data likelihood
4     given  $f(y_{\text{train}} | X; w) \sim \text{Normal}(Xw, \text{likelihood\_var})$ .
5
6     Args:
7         w: Weights
8         X: Training design matrix with first col all ones (np.matrix)
9         y_train: Training response vector (np.matrix)
10        likelihood_var: likelihood variance
11
12    Returns:
13        likelihood: Data likelihood (float)
14    '''
15
16    #TO DO
17    likelihood = 1
18    k = X.shape[0]
19    for i in range(k):
20        likelihood *= (2*math.pi*likelihood_var)**0.5*np.exp(-0.5*(y_train[i]-
21        np.dot(X[i], w))**2/likelihood_var)
22
23    return likelihood
```

5.2

```
1 def get_posterior_params(X, y_train, prior, likelihood_var = 0.2**2):
2     '''
3     Implement get_posterior_params. This function returns the posterior
4     mean vector \mu_p and posterior covariance matrix \Sigma_p for
5     Bayesian regression (normal likelihood and prior).
6
7     Note support_code.make_plots takes this completed function as an argument.
8
9     Args:
10        X: Training design matrix with first col all ones (np.matrix)
11        y_train: Training response vector (np.matrix)
12        prior: Prior parameters; dict with 'mean' (prior mean np.matrix)
13              and 'var' (prior covariance np.matrix)
14        likelihood_var: likelihood variance- default (0.2**2) per the lecture slides
```

```

15
16     Returns:
17         post_mean: Posterior mean (np.matrix)
18         post_var: Posterior mean (np.matrix)
19     '''
20
21     # TO DO
22     prior_var = prior['var']
23     post_mean = np.dot(np.dot((np.dot(X.T,X)+ likelihood_var*
24         prior_var.getI()).getI(), X.T), y_train)
25     post_var = (1./likelihood_var*np.dot(X.T,X)+prior_var.getI()).getI()
26     return post_mean, post_var

```

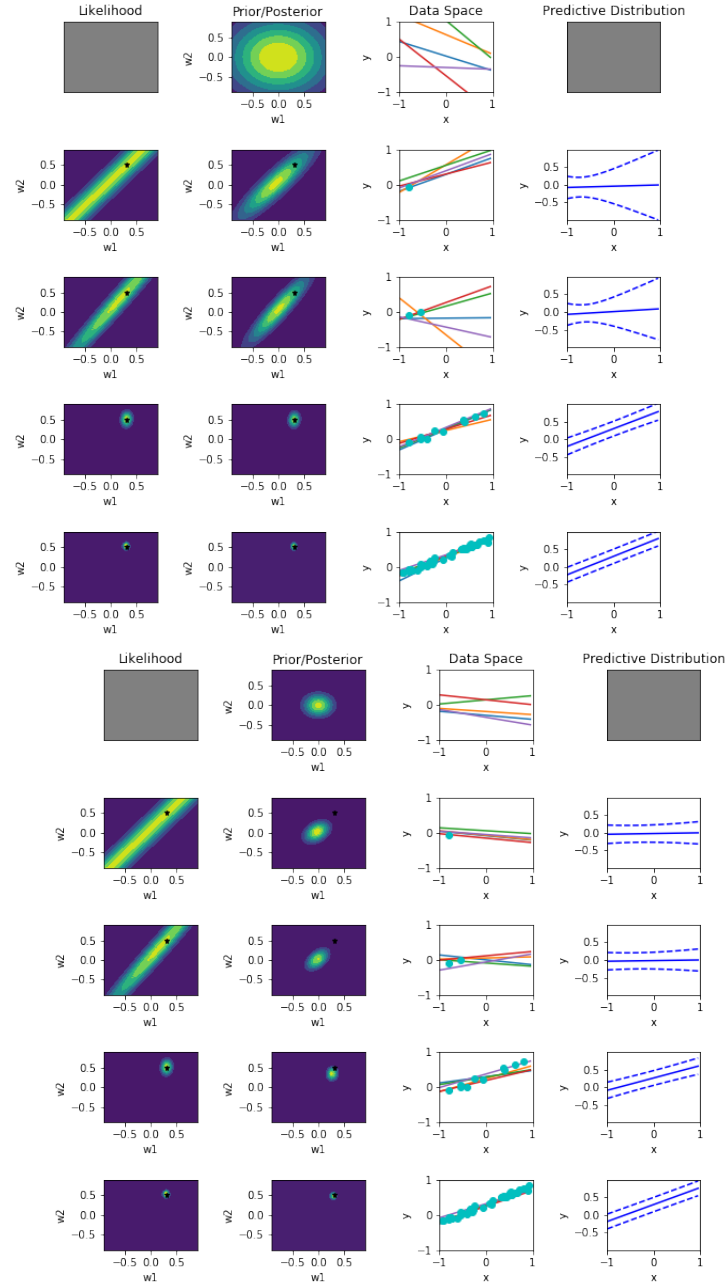
5.3

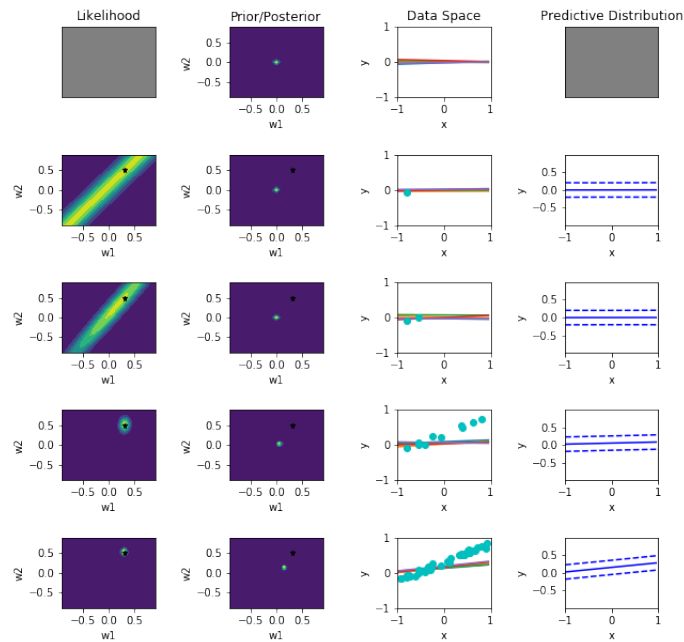
```

1 def get_predictive_params(X_new, post_mean, post_var, likelihood_var = 0.2**2):
2     '''
3     Implement get_predictive_params. This function returns the predictive
4     distribution parameters (mean and variance) given the posterior mean
5     and covariance matrix (returned from get_posterior_params) and the
6     likelihood variance (default value from lecture).
7
8     Args:
9         X_new: New observation (np.matrix object)
10        post_mean, post_var: Returned from get_posterior_params
11        likelihood_var: likelihood variance (0.2**2) per the lecture slides
12
13    Returns:
14        - pred_mean: Mean of predictive distribution
15        - pred_var: Variance of predictive distribution
16    '''
17
18    # TO DO
19    pred_mean = np.dot(post_mean.T, X_new)
20    pred_var = np.dot(X_new.T, np.dot(post_var, X_new))+likelihood_var
21    return pred_mean, pred_var

```

5.4





5.5

(i): the likelihood function shrinks as the sample size increases; when the strength of prior decreases, the likelihood function doesn't seem to get affected much (ii): the posterior distribution shrinks as either the case that sample size increases or the strength of prior decreases. (iii): As the sample size increases and the strength of prior decreases, the confidence interval get closer to the predicted value.

5.6

```

1 lamb = 0.2**2/(1/2)
2 from sklearn.linear_model import Ridge
3 ridge = Ridge(alpha = lamb)
4 ridge.fit(xtrain[:,1], ytrain)

```

first prior covariance is $\Sigma = \frac{1}{2}I$, $\lambda = \frac{\sigma^2}{\Sigma} = 0.08$.

```

ridge.coef_
array([[ 0.52399695]])

ridge.intercept_
array([ 0.30112479])

```

6 [Optional] Coin Flipping: Maximum Likelihood

6.1

$$p(D|\theta) = \theta^2(1 - \theta)$$

6.2

There are 3 different sequences of 3 coin tosses have 2 heads and 1 tail. And the probability of 2 heads and 1 tail is $C_3^2\theta^2(1 - \theta)$.

6.3

$$P(D|\theta) = \theta^{n_h}(1 - \theta)^{n_t}$$

6.4

$$L(\theta) = n_h \log \theta + n_t \log(1 - \theta)$$

$$\max(L(\theta)) = \frac{\partial \log(p(D|\theta))}{\partial \theta}$$

$$n_h \frac{1}{\theta} - n_t \frac{1}{1 - \theta} = 0$$

$$n_h - n_h \theta = n_t \theta$$

$$(n_t + n_h) \theta = n_h$$

$$\theta = \frac{n_h}{n_t + n_h}$$