

Student: Quoc Tuan Vinh, Ngo
ID: 704526

CS-E4820 – Machine Learning: Advanced Probabilistic Methods
Homework Assignment 4

Problem 1. “*ML-II for a linear model*”

```
#####  
# Template for exercise 4.1  
import numpy as np  
import matplotlib.pyplot as plt  
  
# Load the observations  
data = np.loadtxt('ex4_1_data.txt')  
x_obs = data[:,0]  
y_obs = data[:,1]  
  
N_train = 50  
x_train = x_obs[:N_train]  
y_train = y_obs[:N_train]  
  
N_test = 10  
x_test = x_obs[N_train:N_train+N_test]  
y_test = y_obs[N_train:N_train+N_test]  
  
x_range = (-5, 5) # Possible values of x are in this range  
  
# Basis function parameters  
num_basis_functions = 11  
centers = np.linspace(x_range[0], x_range[1], num_basis_functions)  
lambdaval = 0.17  
# You can use here assume the correct basis function centers and lambda ...  
  
def rbf(x, centers, lambdaval):  
    # Radial Basis Function output for input x  
    #  
    # Inputs:  
    # x : input points (one-dimensional array)  
    # centers : basis function centers (one-dimensional array)  
    # lambdaval : basis function width (scalar)  
    #  
    # Output:  
    # Radial Basis Functions evaluated at x (two-dimensional array with len(x)  
    # rows and len(centers) columns)  
    d = x[:,np.newaxis] - centers[np.newaxis,:]  
    y = np.exp(-0.5 * (d ** 2) / lambdaval)  
    return y
```

```

def bayesian_linear_regression(phi_x, y, alpha, beta):
# Bayesian linear parameter model
#
# Inputs:
# phi_x : the basis function applied to x-data (two-dimensional array)
# y : y-data (one-dimensional array)
# alpha : the precision of the weight prior distribution (scalar)
# beta : the precision of the assumed gaussian noise (scalar)
#
# Output:
# the posterior mean, the posterior covariance, the log marginal likelihood

    N, B = phi_x.shape

    # Add here code to compute:
    # m = the posterior mean of w
    # S = the posterior covariance of w
    # S_inv = the inverse of S
    # Equation
    S_inv = (alpha * np.identity(B) + beta * np.dot(phi_x.T, phi_x))
    S = np.linalg.inv(S_inv)
    # Note: This is a corrected version of equation 18.1.19 from Barbers book
    d = beta * np.dot(phi_x.T, y)
    m = beta * S @ phi_x.T @ y
    log_likelihood = 0.5 * (-beta * np.dot(y, y) + d @ S @ d + np.log(np.linalg.det(2 *
np.pi * S)) + B * np.log(alpha) + N * np.log(beta) - N * np.log(2 * np.pi))
    return m, S, log_likelihood

# Problem 1
# Specify possible values for the alpha and beta parameters to test
alphas = np.logspace(-3, 3, 100)
betas = np.logspace(-3, 3, 100)

ll_best = 0
# Grid search over possible values of alpha and beta
for alpha in alphas:
    for beta in betas:
        # Use here functions rbf and bayesian_linear_regression to compute the
        # log marginal likelihood for given alpha and beta
        phi_x = rbf(x_train, centers, lambdaval)
        __, __, log_likelihood = bayesian_linear_regression(phi_x, y_train, alpha, beta)

        # What are the optimal values of alpha and beta, that maximize the marginal
        # likelihood?
        if log_likelihood > ll_best:
            ll_best = log_likelihood
            best_alpha = alpha
            best_beta = beta

        # Fit the model one more time using the optimal alpha and beta and the
        training

```

```

# data to get m for the optimal model
print('Best alpha :', best_alpha)
print('Best beta :', best_beta)
best_m, _, __ = bayesian_linear_regression(phi_x, y_train, best_alpha, best_beta)

# Compute the final regression function
x_coord = np.linspace(x_range[0], x_range[1], 100)
# Compute the predicted values for inputs in x_coord using best_m
y_mean = rbf(x_coord, centers, lambdaval) @ best_m.T

# Plot the final learned regression function, together with the samples
plt.plot(x_coord, y_mean, label="learned model")
plt.plot(x_train, y_train, 'kx', label="training data")
plt.plot(x_test, y_test, 'rx', label="testing data")

# Make predictions for inputs in the test data, so that you get
# predictions 'y_pred' for inputs in x_test.
y_pred = rbf(x_test, centers, lambdaval) @ best_m.T ##WHY???????

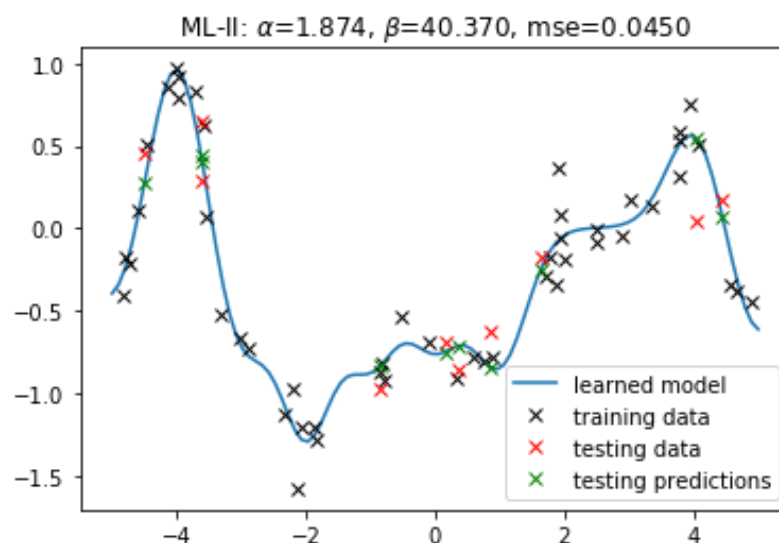
# Plot the predictions
plt.plot(x_test, y_pred, 'gx', label="testing predictions")

# Compute the mean squared prediction error for the test data.
mse_test = 1 / N_test * np.sum(y_test @ y_test - 2*y_test@y_pred +
y_pred@y_pred)

plt.legend()
plt.title("ML-II:  $\alpha$ =%.3f,  $\beta$ =%.3f, mse=%.4f" % (best_alpha,
best_beta, mse_test))
plt.show()
#####

```

The results can be summarized as the plot here:



Problem 2. “Optimizing linear model hyperparameters with validation set”

```
#####
# Template for exercise 4.2

import numpy as np
import matplotlib.pyplot as plt

# Load the observations
data = np.loadtxt('ex4_1_data.txt')
x_obs = data[:,0]
y_obs = data[:,1]

# Training data
N_train = 40
x_train = x_obs[:N_train]
y_train = y_obs[:N_train]

# Validation data
N_valid = 10
x_valid = x_obs[N_train:N_train+N_valid]
y_valid = y_obs[N_train:N_train+N_valid]

# Testing data
N_test = 10
x_test = x_obs[N_train+N_valid:]
y_test = y_obs[N_train+N_valid:]

x_range = (-5, 5) # Possible values of x are in this range

# Basis function parameters
num_basis_functions = 11
centers = np.linspace(x_range[0], x_range[1], num_basis_functions)
lambdaval = 0.17
# You can use here assume the correct basis function centers and lambda ...

def rbf(x, centers, lambdaval):
    # Radial Basis Function output for input x
    #
    # Inputs:
    # x : input points (one-dimensional array)
    # centers : basis function centers (one-dimensional array)
    # lambdaval : basis function width (scalar)
    #
    # Output:
    # Radial Basis Functions evaluated at x (two-dimensional array with len(x)
    # rows and len(centers) columns)
    d = x[:,np.newaxis] - centers[np.newaxis,:]
    y = np.exp(-0.5 * (d ** 2) / lambdaval)
    return y
```

```

def bayesian_linear_regression(phi_x, y, alpha, beta):
# Bayesian linear parameter model
#
# Inputs:
# phi_x : the basis function applied to x-data (two-dimensional array)
# y : y-data (one-dimensional array)
# alpha : the precision of the weight prior distribution (scalar)
# beta : the precision of the assumed gaussian noise (scalar)
#
# Output:
# the posterior mean, the posterior covariance, the log marginal likelihood

    N, B = phi_x.shape

    # Add here code to compute:
    # m = the posterior mean of w
    # S = the posterior covariance of w
    # S_inv = the inverse of S
    # Equation
    S_inv = (alpha * np.identity(B) + beta * np.dot(phi_x.T, phi_x))
    S = np.linalg.inv(S_inv)
    # Note: This is a corrected version of equation 18.1.19 from Barbers book
    d = beta * np.dot(phi_x.T, y)
    m = beta * S @ phi_x.T @ y
    log_likelihood = 0.5 * (-beta * np.dot(y, y) + d @ S @ d + np.log(np.linalg.det(2 *
np.pi * S)) + B * np.log(alpha) + N * np.log(beta) - N * np.log(2 * np.pi))
    return m, S, log_likelihood

# Specify possible values for the alpha and beta parameters to test
alphas = np.logspace(-3, 3, 100)
betas = np.logspace(-3, 3, 100)

min_mse_valid = float("inf")
# Grid search over possible values of alpha and beta
for alpha in alphas:
    for beta in betas:
        # Use here functions rbf and bayesian_linear_regression to fit the
        # model and compute the prediction error (using mean squared error)
        # for the validation data
        phi_x = rbf(x_train, centers, lambdaval)
        m, _, __ = bayesian_linear_regression(phi_x, y_train, alpha, beta)

        y_pred_valid = rbf(x_valid, centers, lambdaval) @ m.T
        mse_valid = 1 / N_valid * np.sum(y_valid @ y_valid - 2*y_valid@y_pred_valid
+ y_pred_valid @ y_pred_valid)
        if min_mse_valid > mse_valid:
            min_mse_valid = mse_valid
            best_alpha = alpha
            best_beta = beta

```

```

# What are the optimal values of alpha and beta, that minimize the prediction
# error in the validation data?
print('Best alpha :', best_alpha)
print('Best beta :', best_beta)

# Fit the model one more time using the optimal alpha and beta and all data
# available for model fitting (both training and validation sets)
N_train_both = 50
x_train_both = x_obs[:N_train_both]
y_train_both = y_obs[:N_train_both]

phi_x_both = rbf(x_train_both, centers, lambdaval)
best_m, _, __ = bayesian_linear_regression(phi_x_both, y_train_both, best_alpha,
best_beta)

x_coord = np.linspace(x_range[0], x_range[1], 100)
# Compute the predicted values for inputs in x_coord using best_m
y_mean = rbf(x_coord, centers, lambdaval) @ best_m.T

# Plot the final learned regression function, together with the samples
plt.plot(x_coord, y_mean, label="learned model")
plt.plot(x_train, y_train, 'kx', label="training data")
plt.plot(x_valid, y_valid, 'bx', label="validation data")
plt.plot(x_test, y_test, 'rx', label="testing data")

# Make predictions for inputs in the test data, so that you get
# predictions 'y_pred' for inputs in x_test.
y_pred = rbf(x_test, centers, lambdaval) @ best_m.T

# Plot the predictions
plt.plot(x_test, y_pred, 'gx', label="testing predictions")

# Compute the mean squared prediction error for the test data.
mse_test = 1 / N_test* np.sum(y_test @ y_test - 2*y_test@y_pred + y_pred@y_pred)

plt.legend()
plt.title("Validation:  $\alpha$ =%.3f,  $\beta$ =%.3f, mse=%.4f" %
        (best_alpha, best_beta, mse_test))
plt.show()

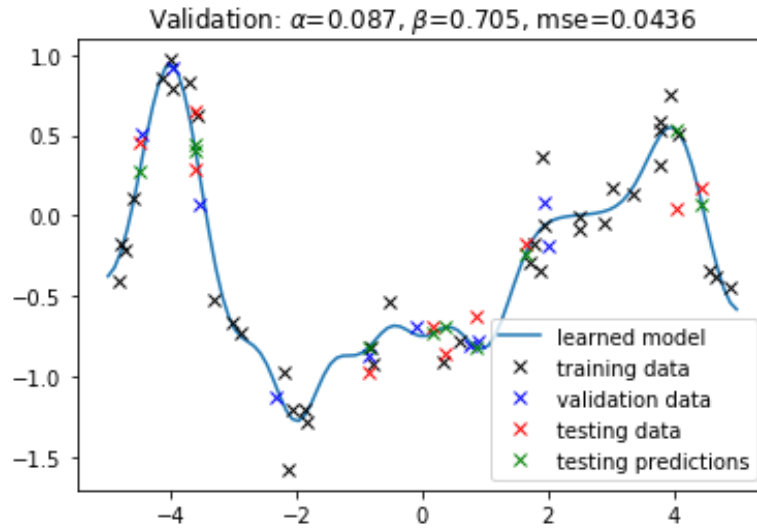
```

The results can be summarized as the plot here:

```

Best alpha : 0.0869749002618
Best beta  : 0.705480231072

```



Problem 3. “Posterior of regression weights”

Suppose $y_i = \mathbf{w}^T \mathbf{x}_i + \epsilon_i$ for $i = 1, \dots, n$ where $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$, and prior $\mathbf{w} \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I})$

To calculate the posterior of \mathbf{w} , we can follow the following steps:

$$\begin{aligned}
 & p(\mathbf{w} | \mathbf{y}, \mathbf{x}, \alpha, \beta) \\
 &= \frac{p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) p(\mathbf{w} | \alpha)}{p(\mathbf{x}, \mathbf{y} | \alpha, \beta)} \\
 &\propto p(\mathbf{w} | \alpha) \times p(\mathbf{y} | \mathbf{x}, \mathbf{w}, \beta) \\
 &= \left(\frac{\alpha}{2\pi} \right)^{\frac{B}{2}} \exp \left(-\frac{\alpha}{2} \mathbf{w}^T \mathbf{w} \right) \times \left(\frac{\beta}{2\pi} \right)^{\frac{n}{2}} \times \exp \left(-\frac{\beta}{2} \sum_{i=1}^n (y_i - \mathbf{w}^T \mathbf{x}_i)^2 \right) \\
 &= \exp \left(-\frac{\mathbf{w}^T (\alpha \mathbf{I}) \mathbf{w}}{2} - \frac{\sum_i^n (\beta \mathbf{w}^T \mathbf{x}_i \mathbf{x}_i^T \mathbf{w} - 2 \beta y_i \mathbf{w}^T \mathbf{x}_i)}{2} \right) + \text{const} \\
 &= \exp \left(-\left(\frac{1}{2} \left(\mathbf{w}^T \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \mathbf{w} \right) - \left(\beta \sum_{i=1}^n y_i \mathbf{x}_i \right) \mathbf{w} \right) \right) + \text{const}
 \end{aligned}$$

Continuing by applying completing the square trick to the transformation above, we have:

$$p(\mathbf{w} | \mathbf{y}, \mathbf{x}, \alpha, \beta)$$

$$\propto \exp \left(- \left(\frac{1}{2} \left(\mathbf{w} - \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \beta \sum_{i=1}^n y_i \mathbf{x}_i \right)^T \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right) \left(\mathbf{w} - \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \beta \sum_{i=1}^n y_i \mathbf{x}_i \right) \right) \right)$$

$$\propto \mathcal{N}(\mathbf{w} \mid \mathbf{m}, \mathbf{S})$$

with:

$$\mathbf{S} = \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1}$$

$$\mathbf{m} = \left(\alpha \mathbf{I} + \beta \sum_{i=1}^n \mathbf{x}_i \mathbf{x}_i^T \right)^{-1} \beta \sum_{i=1}^n y_i \mathbf{x}_i = \beta \mathbf{S} \sum_{i=1}^n y_i \mathbf{x}_i$$

, which is what we need to prove.

Problem 4. “Poisson regression with Laplace approximation”

Suppose we have y_i as observed counts, $\mathbf{x}_i = \{x_1, \dots, x_n\}$, $\boldsymbol{\theta}$ as regression weights. We also have:

$$y_i \mid \boldsymbol{\theta} \sim \text{Poisson}(\exp(\boldsymbol{\theta}^T \mathbf{x}_i))$$

$$\boldsymbol{\theta} \sim \mathcal{N}(0, \alpha^{-1} \mathbf{I})$$

a. Derive the gradient $\nabla \log p(\boldsymbol{\theta} \mid \mathbf{y})$ and Hessian $\mathbf{H} = \nabla \nabla \log p(\boldsymbol{\theta} \mid \mathbf{y})$:

First, we need to compute $\log p(\boldsymbol{\theta} \mid \mathbf{y}, \alpha)$ as below:

$$\log p(\boldsymbol{\theta} \mid \mathbf{y}, \alpha)$$

$$= \log p(\mathbf{y} \mid \boldsymbol{\theta}, \alpha) + \log p(\boldsymbol{\theta} \mid \alpha) - \log p(\mathbf{y} \mid \alpha)$$

$$= \log p(\mathbf{y} \mid \boldsymbol{\theta}, \alpha) + \log p(\boldsymbol{\theta} \mid \alpha) + \text{const}$$

$$= \log \prod_{i=1}^n \frac{e^{-e^{\boldsymbol{\theta}^T \mathbf{x}_i}} \times e^{y_i \boldsymbol{\theta}^T \mathbf{x}_i}}{y_i!} + \log \left(\left(\frac{\alpha}{2\pi} \right)^{\frac{B}{2}} \times e^{-\frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta}} \right) + \text{const}$$

$$= \sum_{i=1}^n (-e^{\boldsymbol{\theta}^T \mathbf{x}_i} + y_i \boldsymbol{\theta}^T \mathbf{x}_i) - \frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \text{const}$$

Therefore, we can derive the negative gradient as below:

$$\begin{aligned}
 & -\nabla \log p(\boldsymbol{\theta}|\mathbf{y}) \\
 &= - \frac{\delta \left(\sum_{i=1}^n (-e^{\boldsymbol{\theta}^T \mathbf{x}_i} + y_i \boldsymbol{\theta}^T \mathbf{x}_i) - \frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \text{const} \right)}{\delta \boldsymbol{\theta}} \\
 &= \sum_{i=1}^n (e^{\boldsymbol{\theta}^T \mathbf{x}_i} \mathbf{x}_i^T - y_i \mathbf{x}_i^T) + \alpha \boldsymbol{\theta}^T
 \end{aligned}$$

And the Hessian can be derived as below:

$$\begin{aligned}
 \mathbf{H} &= -\nabla \nabla \log p(\boldsymbol{\theta}|\mathbf{y}, \alpha) \\
 &= - \frac{\delta \left(\sum_{i=1}^n (-e^{\boldsymbol{\theta}^T \mathbf{x}_i} \mathbf{x}_i^T + y_i \mathbf{x}_i^T) - \alpha \boldsymbol{\theta}^T \right)}{\delta \boldsymbol{\theta}} \\
 &= \sum_{i=1}^n (\mathbf{x}_i e^{\boldsymbol{\theta}^T \mathbf{x}_i} \mathbf{x}_i^T) + \alpha \mathbf{I}
 \end{aligned}$$

b. “Laplace Approximation”

According to the lecture slide, we can obtain the Laplace approximation as:

$$\tilde{E}(\boldsymbol{\theta}) = E(\bar{\boldsymbol{\theta}}) + \frac{1}{2} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})^T \mathbf{H}_{\bar{\boldsymbol{\theta}}} (\boldsymbol{\theta} - \bar{\boldsymbol{\theta}})$$

, in which $E(\boldsymbol{\theta}) = -\log p(\boldsymbol{\theta}|\mathbf{y}, \alpha) = \sum_{i=1}^n (e^{\boldsymbol{\theta}^T \mathbf{x}_i} - y_i \boldsymbol{\theta}^T \mathbf{x}_i) + \frac{\alpha}{2} \boldsymbol{\theta}^T \boldsymbol{\theta} + \text{const}$

The function above follows a Gaussian approximation $q(\boldsymbol{\theta}|\mathbf{y}, \alpha) \propto \exp(-\tilde{E}(\boldsymbol{\theta}))$ to $q(\boldsymbol{\theta}|\mathbf{y}, \alpha)$ as:

$$q(\boldsymbol{\theta}|\mathbf{y}, \alpha) \sim \mathcal{N}(\boldsymbol{\theta}|\mathbf{m}, \mathbf{S})$$

With:

$$\mathbf{S} = \mathbf{H}_{\bar{\boldsymbol{\theta}}}^{-1} \text{ and } \mathbf{m} = \bar{\boldsymbol{\theta}}$$

Where:

- $\bar{\boldsymbol{\theta}}$ is the minimum value of $E(\boldsymbol{\theta})$ that can be found analytically (root of the derivative) or by numerical optimization (such as Newton’s method).
- $\mathbf{H}_{\bar{\boldsymbol{\theta}}}^{-1}$ is the Hessian $\mathbf{H}_{\bar{\boldsymbol{\theta}}}$ of $E(\boldsymbol{\theta})$ at $\bar{\boldsymbol{\theta}}$.

c. “Compare between Laplace and true posterior”

Please refer to the code below:

#####

ML: Advanced Probabilistic Methods

```

# Template for exercise 4.4
import numpy as np
import matplotlib.pyplot as plt

# get some data
data = np.loadtxt('ex4_4_data.txt')
x = data[:,0]
y = data[:,1]

theta_true = np.pi / 4 # true parameter used to generate the data
alpha = 1e-2 # prior's parameter

# compute Laplace approximation
theta_lapl = 0.5 # initial

# iterate to optimum with newton's method to find the MAP estimate for theta
for iter in range(100):
    grad = - y @ x + np.sum(np.exp(theta_lapl * x)*x) + alpha * theta_lapl
    H = np.sum(x * np.exp(theta_lapl * x) * x) + alpha
    theta_lapl = theta_lapl - grad / H # do newton step

# compute Hessian at optimum
H = np.sum(x * np.exp(theta_lapl * x) * x) + alpha
difference = theta_lapl - theta_true

# plot posterior densities
theta = np.linspace(0.55, 0.95, 1000)
post_true = np.zeros(len(theta))

for i in range(len(theta)):
    # log posterior:
    from scipy.misc import factorial
    post_true[i] = (np.dot(y, x * theta[i]) - np.sum(np.exp(x * theta[i]) -
        np.log(factorial(y)))) - 0.5*alpha*np.dot(theta[i], theta[i]))
M = np.amax(post_true)
post_true = np.exp(post_true-M) / np.sum(np.exp(post_true-M)) / (theta[1]-
theta[0]) # normalize

import scipy.stats
post_laplace = np.zeros(len(theta))
for i in range(len(theta)):
    post_laplace[i] = scipy.stats.norm(theta_lapl, 1/np.sqrt(H)).pdf(theta[i])
    # compute approximative density at the points 'theta'
    # Hint: you can use norm.pdf from scipy.stats

plt.figure(1)
plt.plot(theta, post_true, '-k', label="True posterior")
plt.plot(theta, post_laplace, '-.r', label="Laplace approximation")
plt.plot(theta_true, 0, 'o', label="True value")
plt.xlim(0.55, 0.95)

```

```

plt.xlabel('$\\theta$')
plt.title('Posterior $p(\\theta | y)$')
plt.legend()

plt.figure(2)
plt.plot(x, y, 'o', x, np.exp(theta_lapl*x), '-r')
plt.xlabel('x')
plt.ylabel('y')
plt.title('Data')
plt.show()
#####

```

We have the plots as below:

