

CS299 Machine Learning: Assign#1

Due on Thursday, March 14th, 2018

Andrew Ng 6:30 am

Bryan Zhang

Contents

Problem 1	3
1(a)	3
1(b)	3
1(c)	6
Problem 2	7
2(a)	7
2(b)	7
2(c)&2(d)	7
Problem 3	8
3(a)	8
3(b) & 3(c)	9
Problem 4	10
4(a)	10
Problem 5	11
5(a)	11
5(a), i	11
5(a), ii	11
5(a), iii	12
5(b)	12
5(b), i	12
5(b), ii	13
5(b), iii	14
5(c)	18
5(c), i	18
5(c), ii	18
5(c), iii	18

Problem 1

1(a)

$$H = \nabla \cdot \nabla^T J(\theta)$$

$$\text{So } z^T H z = z^T (\nabla \cdot \nabla^T \cdot J(\theta)) z = (z^T \cdot \nabla)^2 J(\theta) = (z^T \nabla^T \cdot J(\theta))^2 \cdot \frac{1}{J(\theta)}.$$

Since the first factor of $z^T H z$ is a square term which is definitely bigger than 0, So we focus on the factor $\frac{1}{J(\theta)}$. Because $J(\theta) = -\frac{1}{m} \sum_{i=1}^m \log(h_{\theta}(y^{(i)} x^{(i)}))$ and the hypothesis function is a sigmoid function ranging from (0,1), the log of the hypothesis function must be negative and thus the cost function must be positive. So does the $z^T H z$.

q.e.d.

1(b)

The optimized θ is $[-2.62042271649454, 0.760346235045246, 1.17193037252339]$.

```

close all; clear; clc;

% the rows of X is input variables
% the rows of Y is the response variables
5 fileIDX = fopen('logistic_x.txt', 'r');
  sizeX = [2 Inf];
  formatSpec = '%f';
  X = fscanf(fileIDX, formatSpec, sizeX).';
  % append the intercept term
10 X = [ones(size(X, 1), 1) X];

  fileIDY = fopen('logistic_y.txt', 'r');
  Y = fscanf(fileIDY, formatSpec);

15 %plot the x and y
  % the sub dataset for respective y is 1
  Xp = X(1:50, :);
  %the other half of sub dataset
  Xn = X(51:size(X, 1), :);
20 sz = 25;
  x1range = [0 8];
  x2range = [-5 4];

25 %the implementation of the Newton's Method for logistic regression
  %serveral instance variables
  THETA_INITIAL = zeros(1, size(X, 2));
  ERROR_MARGINS = 0.00001;
30 %the size of sample space
  m = size(Y, 1);

35 %the sigmoid function
  sigmoid = @(z) 1./(1 + exp(-z));
  %the cost function

```

```

J = @(theta) 1 / m * sum(log(sigmoid(Y.*( X * theta '))));

40 thetaOptimized = getTheta(J, THETA_INITIAL,ERROR_MARGINS);
%the hypthesis functon
h = @(X) sigmoid(X * thetaOptimized ');

%plot the decision boundary
45 % step size for the accuracy of the boundary curve
inc = 0.01;

% generate grid coordinates
[x1, x2] = meshgrid(xlrange(1):inc:xlrange(2), x2range(1):inc:x2range(2));
50 imageSize = size(x1);

x1x2 = [x1(:) x2(:)]; % make the (x1, x2) pairs as row vectors

hypothesis = zeros(length(x1x2), 1);
55 for i = 1:length(x1x2)
    Xhypo= [1 x1x2(i,:)];
    htemp = h(Xhypo);
    if htemp > 0.5
        hypothesis(i) = 1;
60     else
        hypothesis(i) = 0;
    end
end
%reshape the hypothesis to be positioned on each grid point
65 decisionMap = reshape(hypothesis, imageSize);

% plot the decision boundary
figure
imagesc(xlrange, x2range, decisionMap);
70 hold on;
cmap = [1 0.8 0.8; 0.9 0.9 1];
colormap(cmap);

75 scatter(Xp(:,2), Xp(:, 3), sz, 'red', 'filled');
hold on;
scatter(Xn(:,2), Xn(:, 3), sz, 'blue', 'filled', 'd');
hold on;
title("Assign#1-1b: Logistic Regression Optimized with Newton's Method ");
80 xlim(xlrange); ylim(x2range);
xlabel('0 < X1 < 8');
ylabel('-5 < X2 < 4');
legend('y = 1', 'y = -1', 'Location', 'Southwest')
hold on;
85 %save the image
saveas(gcf, '1b.png')

%disp tests
90 disp(sum(sigmoid([1 2 3])));

```

```

disp(J([0 0 0]));
disp(getGradient(J, [0 0 0]));
disp(getHessian(J, [0 0 0]));
disp(getTheta(J, [0 0 0], 0.00001));
95
% get the optimized theta
function theta = getTheta(costfunc, thetaIni, errorMargins)
%% costfunc is the cost function for the logistic regression
%% thetaIni is the start point to search for the optimized theta
100 %% return the optimized theta which can make the gradient down to zero
%% thus the cost function down to the minimal
theta = thetaIni;
grad = getGradient(costfunc, theta);
while norm(grad) > errorMargins
105     grad = getGradient(costfunc, theta);
    H = getHessian(costfunc, theta);
    disp('H: '); disp(H);
    disp('grad: '); disp(grad);
    theta = theta - grad / H;
110     disp('theta '); disp(theta);
end
end

%get the hessian
115 function H = getHessian(f, x)
%% f is a function
%% x is a input varibale
%% return the hessian for the function at x
gx = getGradient(f, x);
120 H = zeros(size(x, 2));
h = 0.00001;

%iterate over al indexes in x
for i = 1: size(x, 2)
125     oldValues = x(i);
    x(i) = oldValues + h;
    gxh = getGradient(f, x); %get the grad f(x + h)
    x(i) = oldValues; % restore to previous value

130     %compute the second partial derative
    H(:, i) = (gxh - gx) ./ h;
    %iterate over to the next variable
end
end
135

% get the gradient
function grad = getGradient(f, x)
%% f is a function
140 %% x is a input varibale
%% return the gradient for the function at x

fx = f(x);

```

```

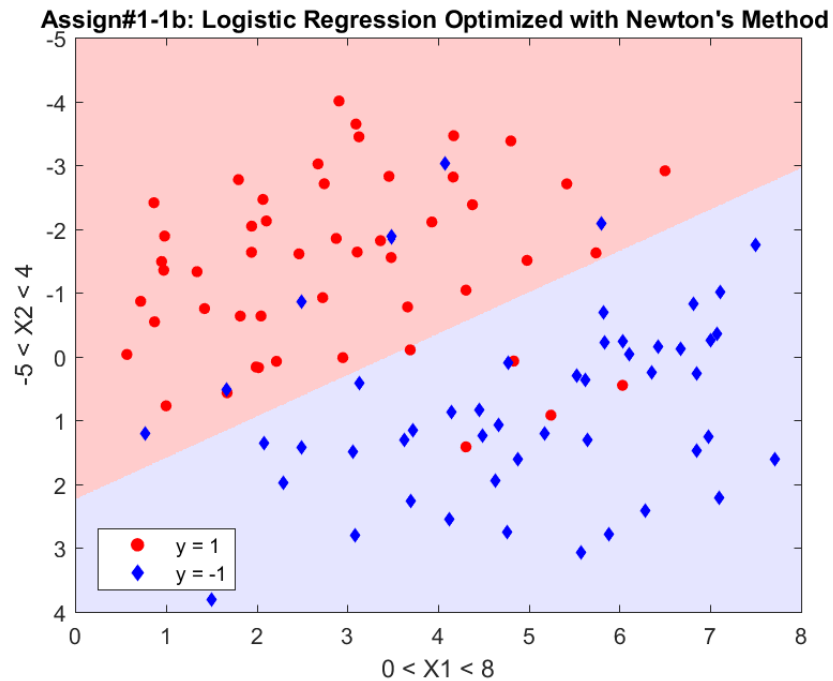
grad = zeros(size(x));
145 h = 0.00001;

%iterate over all indexes in x
for i = 1:size(x, 2)
    oldValues = x(i);
150 x(i) = oldValues + h; %increment by h
    fxh = f(x); % evaluate f(x + h)
    x(i) = oldValues; %restore to the previous value for x(i)

%compute the partial derative
155 grad(i) = (fxh - fx) / h; %the slop
%iterate to the next index
end
end

```

1(c)



Problem 2

2(a)

The poisson distribution, $p(y; \lambda) = \frac{e^{-\lambda} \lambda^y}{y!} = \frac{1}{y!} e^{(\log(\lambda)y - \lambda)}$. Compared with the exponential family $p(y; \eta) = b(y) \exp(n^T T(y) - a(\eta))$, we can get

$$\begin{aligned} b(y) &= \frac{1}{y!} \\ \eta &= \log(\lambda) \\ T(y) &= y \\ \alpha(\eta) &= \lambda = e^\eta \end{aligned}$$

2(b)

The canonical response function gives the mean of the distribution, which in this case, a Poisson distribution, is λ , in terms of the natural parameter η . So the canonical response function for the Poisson distribution is $g(\eta) = e^\lambda$.

2(c)&2(d)

Here we directly look at all GLMs instead of the Poisson case. Our goal is first to derive the $p(y^{(i)} | x^{(i)}; \theta)$ with respect to θ_j .

For exponential family, $p(y; \eta) = b(y) \exp(n^T T(y) - a(\eta))$. After using assumption 3, we substitute the natural parameter for $\theta^T X$. For training set

$$\{(X^{(i)}, Y^{(i)}); i = 1, \dots, m\},$$

the distribution is

$$p(y^{(i)} | x^{(i)}; \theta) = b(y) \exp(\theta^T X^{(i)} - a(\theta^T X^{(i)})). \quad (1)$$

Then we take derivative of the log-likelihood above with respect to θ_j , which shows,

$$\begin{aligned} & \frac{d}{d\theta_j} \log(p(y^{(i)} | x^{(i)}; \theta)) \\ &= \frac{1}{p(y^{(i)} | x^{(i)}; \theta)} \cdot \frac{d}{d\theta_j} p(y^{(i)} | x^{(i)}; \theta) \\ &= \frac{1}{p(y^{(i)} | x^{(i)}; \theta)} \cdot (p(y^{(i)} | x^{(i)}; \theta) [T(Y^{(i)}) - \frac{d}{d\eta} a(\eta)] \cdot x_j) \\ &= (T(Y^{(i)}) - \frac{d}{d\eta} a(\eta)) \cdot x_j \end{aligned} \quad (2)$$

Using the assumption 2, we can get $T(Y^{(i)}) = Y^{(i)}$. Because of this assumption and α being only a hyper-parameter, we need only to prove $\frac{d}{d\eta} a(\eta)$ is the response function when natural parameter η is replaced by $\theta^T X^{(i)}$ to prove the stochastic gradient ascent has the update rule

$$\theta_j = \theta_j - \alpha(h(x) - y)x_j$$

for every j and looping i until m .

$a(\eta)$ is the log partition function, which is to make sure the probability function can integrate to unity at last. So

$$a(\eta) = \log \left[\int b(Y) \exp(\eta^T T(Y)) v(dY) \right]$$

$$\begin{aligned}
\frac{d}{d\eta^T} a(\eta) &= \frac{\int b(Y)T(Y)\exp(\eta^T T(Y)v(dY)}{\int b(Y)\exp(\eta^T T(Y)v(dY)} \\
&= \int T(Y)[b(Y)\exp(\eta^T T(Y)v(dY) \cdot \exp(-\log[\int \exp(\eta^T T(Y)b(Y)v(dY)]) \\
&= \int T(Y)[b(Y)\exp(\eta^T T(Y)v(dY) \cdot \exp(-a(\eta)) \\
&= \int T(Y)[b(y)\exp(\eta^T T(Y) - a(\eta))]v(dY) \\
&= \int T(Y)p(Y|X; \theta)V(dY) \\
&= E(T(Y)|X)
\end{aligned} \tag{3}$$

Since the response function give the estimation of Y (the $T(Y)$ in most cases) in terms of input features X , so $\frac{d}{d\eta^T} a(\eta)$ is $h_\theta(X)$. Thus, we can conclude that the derivative of log-likelihood given $(X^{(i)}, Y^{(i)})$ with respect to θ_j is the update rule for stochastic gradient ascent. Since 2(c) is only a special case for 2(d), q.e.d.

Problem 3

3(a)

Our goal ¹ is to prove $p(y = 1|x; \phi, \Sigma, \mu_{-1}, \mu_1) = \frac{1}{1 + \exp(-\theta^T x - \theta_0)}$. and the similiar form for $y = -1$. We can use Bayes's Formula to prove

$$\begin{aligned}
&p(y = 1|x; \phi, \Sigma, \mu_{-1}, \mu_1) \\
&= \frac{p(x|y = 1) \cdot p(y = 1)}{p(x)} \\
&= \frac{p(x|y = 1) \cdot p(y = 1)}{p(x|y = 1) \cdot p(y = 1) + p(x|y = -1) \cdot p(y = -1)} \\
&= \frac{1}{1 + \frac{p(x|y=-1) \cdot p(y=-1)}{p(x|y=1) \cdot p(y=1)}}
\end{aligned} \tag{4}$$

The form is already similiar to our goal. Next we only need to prove the fraction $\frac{p(x|y=-1) \cdot p(y=-1)}{p(x|y=1) \cdot p(y=1)} = \exp(\theta^T x + \theta_0)$.

$$\begin{aligned}
&\frac{p(x|y = -1) \cdot p(y = -1)}{p(x|y = 1) \cdot p(y = 1)} \\
&= \exp[-1/2(x - \mu_{-1})^T \Sigma^{-1}(x - \mu_{-1}) + 1/2(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)] \cdot \frac{1 - \phi}{\phi} \\
&= \exp[-1/2(x - \mu_{-1})^T \Sigma^{-1}(x - \mu_{-1}) + 1/2(x - \mu_1)^T \Sigma^{-1}(x - \mu_1)] \cdot \exp(\log(\frac{1 - \phi}{\phi})) \\
&= \exp[\frac{(\mu_1 - \mu_{-1})^T x}{\Sigma} - \frac{\mu_{-1}^2 - \mu_1^2}{2\Sigma} + \log(\frac{1 - \phi}{\phi})]
\end{aligned} \tag{5}$$

Campared with $\exp(-\theta^T x + \theta_0)$, we can get

$$\theta = \frac{\mu_1 - \mu_{-1}}{\Sigma},$$

¹I got this section down by referencing

<https://duphan.wordpress.com/2016/10/27/gaussian-discriminant-analysis-and-logistic-regression/>

$$\theta_0 = \log\left(\frac{1-\phi}{\phi}\right) - \frac{\mu_{-1}^2 - \mu_1^2}{2\Sigma}$$

Since the fraction is just the inverse in case $y = -1$, so the form satisfies in both case; q.e.d.

3(b) & 3(c)

Our goal ² is to perform log-likelihood maximum estimation on the log-likelihood. By Bayes's rule³,

$$\begin{aligned} l(\phi, \mu_1, \mu_{-1}, \Sigma) &= \log \prod_{i=1}^m p(x^{(i)}|y^{(i)}; \mu_{-1}, \mu_1, \Sigma) \cdot p(y^{(i)}) \\ &= \sum_{i=1}^m \log \{ [\phi \cdot N(x^{(i)}|\mu_1, \Sigma)]^{1\{y^{(i)}=1\}} \cdot [(1-\phi) \cdot N(x^{(i)}|\mu_{-1}, \Sigma)]^{1\{y^{(i)}=-1\}} \} \\ &= \sum_{i=1}^m \{ 1\{y^{(i)}=1\} \log \phi + 1\{y^{(i)}=1\} \log N(x^{(i)}|\mu_1, \Sigma) + 1\{y^{(i)}=-1\} \log(1-\phi) + 1\{y^{(i)}=-1\} \log N(x^{(i)}|\mu_{-1}, \Sigma) \} \end{aligned} \quad (6)$$

Let's first estimate ϕ . Consider the parts that contains ϕ , we have:

$$\sum_{i=1}^m [1\{y^{(i)}=1\} \log \phi + 1\{y^{(i)}=-1\} \log(1-\phi)]$$

Take the derivative over ϕ :

$$\sum_{i=1}^m [1\{y^{(i)}=1\} \frac{1}{\phi} - 1\{y^{(i)}=-1\} \frac{1}{1-\phi}]$$

Setting it to zero, we have

$$\phi = \frac{\sum_{i=1}^m 1\{y^{(i)}=1\}}{\sum_{i=1}^m 1\{y^{(i)}=1\} + \sum_{i=1}^m 1\{y^{(i)}=-1\}}$$

Since $y =$ either 1 or -1,

$$\phi = \frac{1}{m} \sum_{i=1}^m 1\{y^{(i)}=1\}$$

We now move onto estimating μ_1 or μ_{-1} . Since both have identical form, we only need to consider one of mean vectors. Considering the parts contain μ_1 , we have:

$$\begin{aligned} &\sum_{i=1}^m 1\{y^{(i)}=1\} \log N(x^{(i)}|\mu_1, \Sigma) \\ &= \sum_{i=1}^m 1\{y^{(i)}=1\} \frac{-1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) + \text{const} \end{aligned} \quad (7)$$

Take the derivative over μ_1 and set it to zero, we have:

$$\begin{aligned} &\sum_{i=1}^m 1\{y^{(i)}=1\} \Sigma^{-1} (x^{(i)} - \mu_1) = 0 \\ \mu_1 &= \frac{\sum_{i=1}^m 1\{y^{(i)}=1\} x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)}=1\}} \end{aligned}$$

²I got this section down by referencing <http://web.engr.oregonstate.edu/~xfern/classes/cs534/notes/LDA.pdf>

³ $1\{y=i\}$ is the notation introduced in the lecture notes 1 page 27.

Similarly we have:

$$\mu_{-1} = \frac{\sum_{i=1}^m 1\{y^{(i)} = -1\}x^{(i)}}{\sum_{i=1}^m 1\{y^{(i)} = -1\}}$$

Finally, we will estimate Σ . Taking the part that contains Σ we have ⁴ :

$$\begin{aligned} & \sum_{i=1}^m [1\{y^{(i)} = 1\} \log N(x^{(i)} | \mu_1, \Sigma) + 1\{y^{(i)} = -1\} \log N(x^{(i)} | \mu_{-1}, \Sigma)] \\ &= \sum_{i=1}^m [1\{y^{(i)} = 1\} \left[\frac{(x - \mu_1)^T \Sigma^{-1} (x - \mu_1)}{-2} - \frac{1}{2} \log |\Sigma| \right] + 1\{y^{(i)} = -1\} \left[\frac{(x - \mu_{-1})^T \Sigma^{-1} (x - \mu_{-1})}{-2} - \frac{1}{2} \log |\Sigma| \right]] \\ &= -\frac{m}{2} \log |\Sigma| - \frac{1}{2} \sum_{i=1}^m 1\{y^{(i)} = 1\} \text{Tr}(\Sigma^{-1} ((x^{(i)} - \mu_1)^T (x^{(i)} - \mu_1))) - \frac{1}{2} \sum_{i=1}^m 1\{y^{(i)} = -1\} \text{Tr}(\Sigma^{-1} ((x^{(i)} - \mu_{-1})^T (x^{(i)} - \mu_{-1}))) \\ &= -\frac{m}{2} \log |\Sigma| - \frac{m}{2} \text{Tr} \left(\frac{1}{m} \sum_{i=1}^m 1\{y^{(i)} = 1\} \Sigma^{-1} ((x^{(i)} - \mu_1)^T (x^{(i)} - \mu_1)) \right) - \frac{m}{2} \text{Tr} \left(\sum_{i=1}^m 1\{y^{(i)} = -1\} \Sigma^{-1} ((x^{(i)} - \mu_{-1})^T (x^{(i)} - \mu_{-1}))) \right) \\ &= -\frac{m}{2} \log |\Sigma| - \frac{m}{2} \text{Tr}(\Sigma^{-1} \frac{1}{m} (\sum_{i=1}^m 1\{y^{(i)} = 1\} ((x^{(i)} - \mu_1)^T (x^{(i)} - \mu_1)) + \sum_{i=1}^m 1\{y^{(i)} = -1\} ((x^{(i)} - \mu_{-1})^T (x^{(i)} - \mu_{-1})))) \end{aligned} \quad (8)$$

Taking derivative over Σ ⁵ and set it to zero, we have

$$(\Sigma^{-1})^T + \frac{1}{m} \left(\sum_{i=1}^m 1\{y^{(i)} = 1\} ((x^{(i)} - \mu_1)^T (x^{(i)} - \mu_1)) + \sum_{i=1}^m 1\{y^{(i)} = -1\} ((x^{(i)} - \mu_{-1})^T (x^{(i)} - \mu_{-1}))) \right)^T (-\Sigma^{-1})^2 = 0$$

Thus, we can get

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T.$$

Problem 4

4(a)

We will use induction to prove linear invariance of Newton's method. We already have $z^{(0)} = A^{-1}x^{(0)}$. Applying Newton's method to function f , we can have

$$x^{(i+1)} = x^{(i)} - H^{-1} \nabla_x f(x^{(i)})$$

Applying Newton's method to function g , we can have

$$z^{(i+1)} = z^{(i)} - H^{-1} \nabla_z g(z^{(i)})$$

Assuming and substituting $z^{(i)} = A^{-1}x^{(i)}$, we have

$$z^{(i+1)} = A^{-1}x^{(i)} - H^{-1} \nabla_x (A^{-1}x^{(i)})$$

⁴To get $A^T B A = \text{Tr}(B A^T A)$, in which A is a $n \times 1$ vector and B is a $n \times n$ square matrix, $A^T B A = \sum_{i=1}^n \sum_{j=1}^n A_i B_{ij} A_j = \sum_{q=1}^n \sum_{p=1}^n B_{pq} A_p A_q = \text{Tr}(B(A^T A))$.

⁵To figure out how to take derivative of a determinant, I read the lecture notes, Matrix Cookbook, <https://www.ics.uci.edu/~welling/teaching/KernelsICS273B/MatrixCookBook.pdf>.

Here we prove that $\nabla_x(Ax) = A\nabla_x x$, as long as A is independent with x . Since $\nabla_x(Ax)_p = \frac{\partial}{\partial x_p} \sum_{q=1}^n A_q A_p x_p = \sum_{q=1}^n A_q A_p \frac{\partial}{\partial x_p} x_p = (A\nabla_x x)_p$, we can prove our conclusion. The similar way for the hessian matrix can make it eligible to move the constant matrix above the hessian.

$$\begin{aligned} z^{(i+1)} &= A^{-1}x^{(i)} - H^{-1}\nabla_x(A^{-1}x^{(i)}) \\ &= A^{-1}(x^{(i)} - H^{-1}\nabla_x(x^{(i)})) \\ &= A^{-1}x^{(i+1)} \end{aligned} \tag{9}$$

q.e.d. Similar procedure can be applied to gradient optimization to prove its linear invariance.

Problem 5

5(a)

5(a), i

We want to prove $\frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2 = (X\theta - \vec{y})^T W (X\theta - \vec{y})$, we choose W is a diagonal matrix by $m \times m$,

$$W = \begin{bmatrix} \frac{1}{2}w_1 & 0 & 0 & \dots & 0 \\ 0 & \frac{1}{2}w_2 & 0 & \dots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & \frac{1}{2}w_m \end{bmatrix}$$

$$\begin{aligned} LHS &= J(\theta) \\ &= \frac{1}{2} \sum_{i=1}^m w^{(i)} (\theta^T x^{(i)} - y^{(i)})^2 \\ RHS &= (X\theta - \vec{y})^T W (X\theta - \vec{y}) \\ &= [\dots \frac{1}{2} ((x^{(i)})^T \theta - y^{(i)}) w^{(i)} \dots] \cdot \begin{bmatrix} \ddots \\ ((x^{(i)})^T \theta - y^{(i)}) \\ \dots \end{bmatrix} \\ &= LHS \end{aligned} \tag{10}$$

q.e.d

5(a), ii

From i, we know that $J(\theta) = (X\theta - \vec{y})^T W (X\theta - \vec{y})$ under this weight setting. First, we will prove $\nabla_A(XA - Y)^T W (XA - Y) = 2W(X\theta - Y)X^T$, is W is symmetric.

Taking the derivative of the cost function with respect to θ ,⁶

$$\begin{aligned} \nabla_{\theta} J(\theta) &= \nabla_{\theta} (X\theta - \vec{y})^T W (X\theta - \vec{y}) \\ &= \nabla_{\theta} (\theta^T X^T W X\theta - \theta^T X^T W \vec{y} - \vec{y}^T W X\theta + \vec{y}^T W \vec{y}) \\ &= \nabla_{\theta} Tr[\theta^T X^T W X\theta - \theta^T X^T W \vec{y} - \vec{y}^T W X\theta + \vec{y}^T W \vec{y}] \\ &= \nabla_{\theta} (Tr[\theta^T X^T W X\theta] - 2Tr[\vec{y}^T W X\theta]) \\ &= 2[X^T W^T X\theta - X^T W^T \vec{y}] \end{aligned} \tag{11}$$

Setting it to zero, we have

$$\theta = (X^T W^T X)^{-1} X^T W^T \vec{y}$$

⁶ $\nabla_A(XA - Y)^T W (XA - Y) = 2W(X\theta - Y)X^T$, I refer to Matrix Cookbook at <https://www.ics.uci.edu/~welling/teaching/KernelsICS273B/MatrixCookBook.pdf>

5(a), iii

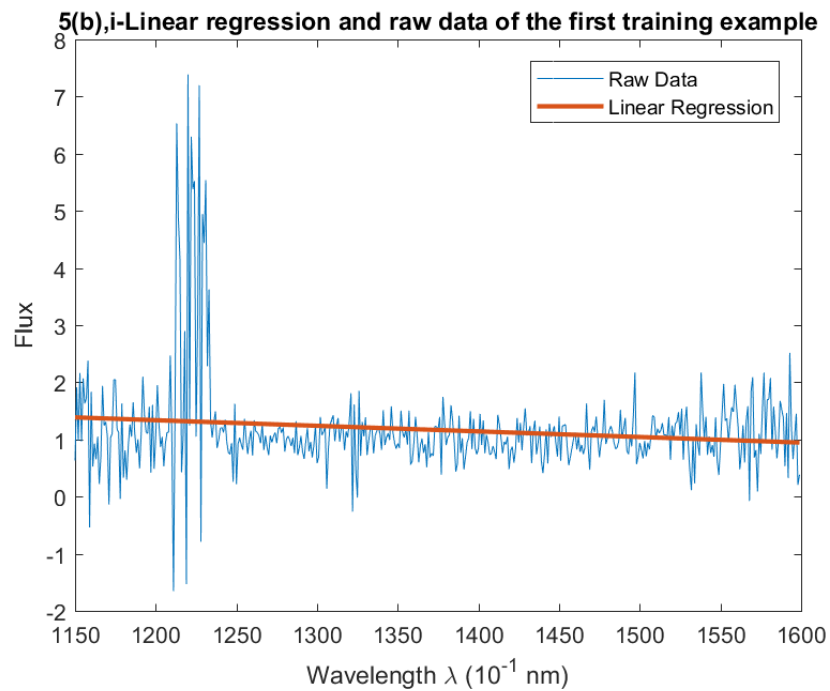
$$\begin{aligned}
 l(\theta) &= \log \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma^{(i)}} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2(\sigma^{(i)})^2}\right) \\
 &= -\sum_{i=1}^m \frac{1}{2(\sigma^{(i)})^2} (y^{(i)} - \theta^T x^{(i)})^2 + z
 \end{aligned} \tag{12}$$

So we only need to optimize the first term. Compared with cost function, we can get $w^{(i)} = \frac{1}{2(\sigma^{(i)})^2}$.

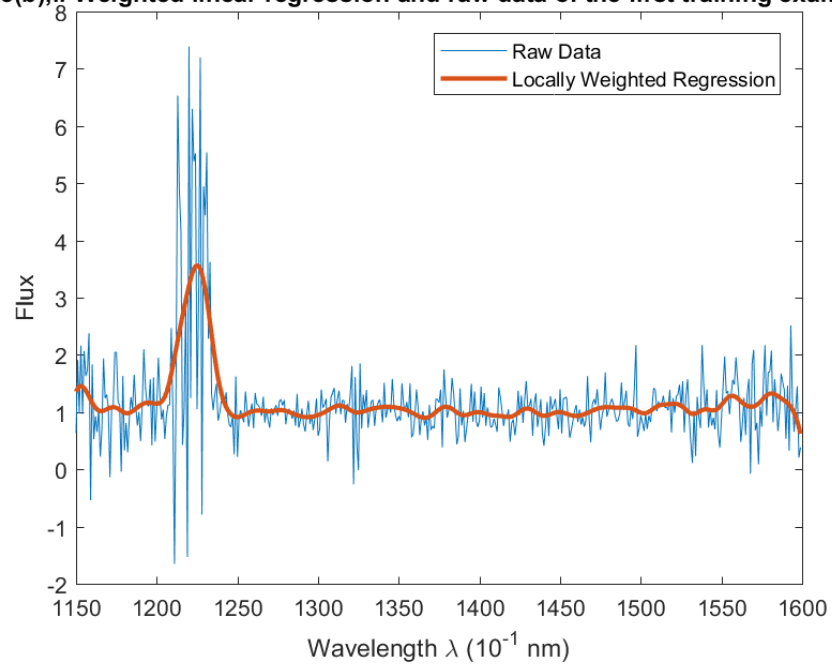
5(b)

5(b), i

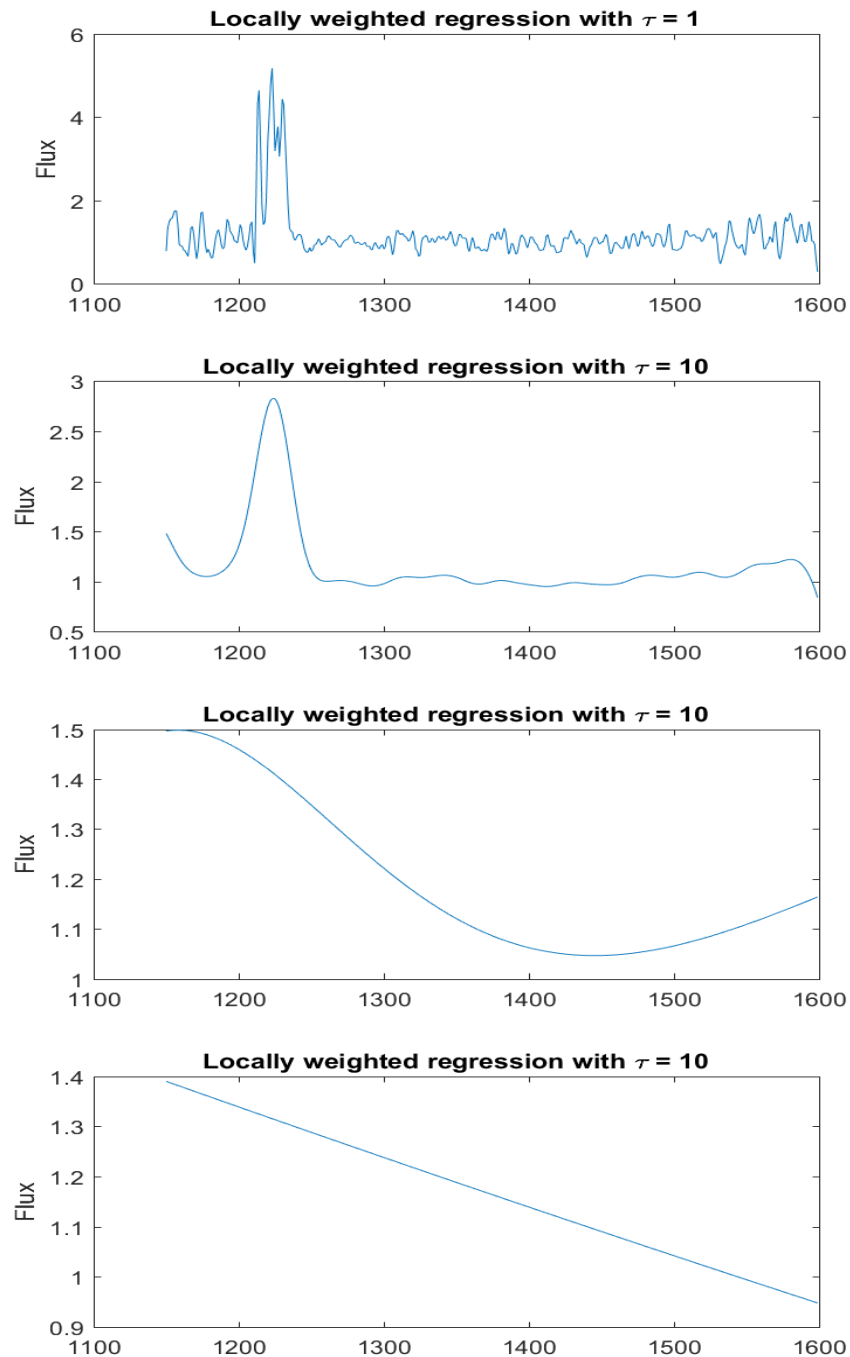
The optimized θ is $[2.513 - 0.0009]^T$.



5(b), ii

5(b),ii-Weighted linear regression and raw data of the first training example

5(b), iii



Below is the code for the section, 5(b),i,ii and iii. The function are seperated into another file for convienience.

```

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%% 5(b) visualize the data
%% i. demands a linear regression with
%% visualization and optimized parameter;
5 %% ii. visualize raw data and weighted
%% linear regression;
%% iii. vary the bandwidth parameter
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
close all; clear; clc;

10 % define the query point
run 'load_quasar_data.m';
%load the first training example
trainEx1 = train_qso(1,:)';

15 %%%%%%%%%% linear regression for the first training example %%%%%%%%%%
%plot the raw data
figure
plot( lambdas, trainEx1);
20 title('5(b),i-Linear regression and raw data of the first training example');
xlabel('Wavelength \lambda (10^{-1} mm)');
ylabel('Flux');
hold on;

25 %compute theta using normal equation
X = [ones(size(lambdas)) lambdas];
theta = ( X' * X)^(-1) * X' * trainEx1;
disp(theta);

30 %plot the linear regression line
YHypoLin = X * theta;
p1 = plot(lambdas, YHypoLin);
legend('Raw Data', 'Linear Regression');
p1(1).LineWidth = 2;

35 saveas(gcf, '5b(i).png')

%%%%%%%%%% weighted linear regression %%%%%%%%%%
40 %plot the raw data
figure
plot( lambdas, trainEx1);
title('5(b),ii-Weighted linear regression and raw data of the first training example')
;
xlabel('Wavelength \lambda (10^{-1} mm)');
45 ylabel('Flux');
hold on;

% for each lambda compute the predict value
LWRX = X;
50 %the demension of sample space
m = size(LWRX, 1);
% the bandwidth parameter

```

```
tau = 5;
LWRY = smoothLWR(X, trainEx1, LWRX, tau);
55 %plot the locally weighted regression
p2 = plot(lambdas, LWRY);
ylim([-2 8]);
legend('Raw Data', 'Locally Weighted Regression')
60 p2(1).LineWidth = 2;
saveas(gcf, '5b(ii).png')

65 %%%%% vary the bandwidthparameter %%%%%

% when tau is 1
tau = 1;
LWRY = smoothLWR(X, trainEx1, LWRX, tau);
70 variedTau = figure(3);
set(variedTau, 'Position', [100, 100, 512, 1200]);
ax1 = subplot(4, 1, 1);
plot(ax1, lambdas, LWRY);
75 title(ax1, 'Locally weighted regression with \tau = 1');
ylabel('Flux');

% when tau is 10
tau = 10;
80 LWRY = smoothLWR(X, trainEx1, LWRX, tau);

ax2 = subplot(4, 1, 2);
plot(ax2, lambdas, LWRY);
title(ax2, 'Locally weighted regression with \tau = 10');
85 ylabel('Flux');

% when tau is 100
tau = 100;
LWRY = smoothLWR(X, trainEx1, LWRX, tau);
90 ax3 = subplot(4, 1, 3);
plot(ax3, lambdas, LWRY);
title(ax3, 'Locally weighted regression with \tau = 10');
ylabel('Flux');

95 % when tau is 1000
tau = 1000;
LWRY = smoothLWR(X, trainEx1, LWRX, tau);

100 ax4 = subplot(4, 1, 4);
plot(ax4, lambdas, LWRY);
title(ax4, 'Locally weighted regression with \tau = 10');
ylabel('Flux');

105 saveas(gca, '5b(iii).png')
```



```

%%% comments on varied tau  %%%
% When tau get bigger, the weight decrease quickly. Thus no input will have
% much difference on the regression model.
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

```

The function file for the section, 5(b),i,ii and iii is here.

```

function outputs = smoothLWR(trainInputs, trainOutputs, queryPoints, tau)
%%
% return smoothed curve
5 %-'trainInputs': training input data
  %-'trainOutputs': training output data
  %-'tau': also called as bandwidth parameter for the gaussian kernel
%%
  outputs = zeros([size(queryPoints, 1) 1]);
10 % the demension of the sample size
  m = size(trainInputs, 1);
  for i = 1:m
    outputs(i) = LWRPredict(trainInputs, trainOutputs, queryPoints(i,:), tau);
  end
15 end

function Output = LWRPredict(trainInputs, trainOutputs, queryPoint, tau)
%%
% return the predicted value for the query point
20 %-'trainInputs': training input data
  %-'trainOutputs': training output data
  %-'queryPoints': the data points we want to predict on
  %-'tau': also called as bandwidth parameter for the gaussian kernel
%%
25 W = getWeight(trainInputs, queryPoint, tau);
  theta = (trainInputs' * W * trainInputs)^(-1) * trainInputs' * W * trainOutputs;
  Output = queryPoint * theta;
end

30 function W = getWeight(trainInputs, queryPoint, tau)
%%
% get the diagonal weighted matrix
%-'trainInputs': training input data
%-'queryPoints': the data points we want to predict on
35 %-'tau': also called as bandwidth parameter for the gaussian kernel
%%

% the demension of the training input
m = size(trainInputs, 1);
40 W = eye(m);
  for i = 1: m
    distance = trace((trainInputs(i,:) - queryPoint)' * (trainInputs(i,:) - queryPoint
    ));
    W(i, i) = exp(- distance / (2* tau^2));
  end
45 end

```

5(c)**5(c), i**

Below is the matlab code for the section 5c, i. The function file are listed in section 5(b).

```

close all; clear; clc;
run 'load_quasar_data.m';

        %%%%%% 5(c),i-smooth the entire dataset %%%%%%%%%%
5 % smooth the test and trian data set to get rid of random noise
smoothed_qso_train = smoothTrainSet(train_qso , lambdas);
smoothed_qso_test = smoothTrainSet(test_qso , lambdas);

save('smoothed');

10
function smoothed_qso_train = smoothTrainSet(train_qso , lambdas)
%%
% return the smoothed training data set
%-'train_qso': the original training data set that contains random noised
15 %-'lambdas': the training input wavelengths
%%
train_inputs = [ones(size(lambdas)) lambdas];
% bandwidth parameter
tau = 5;
20 %sample size
m = size(train_qso , 1);
smoothed_qso_train = zeros(size(train_qso));

for i = 1:m
25     smoothed_qso_train(i,:) = smoothLWR(train_inputs , train_qso(i,:) , train_inputs ,
        tau)';
end
end

```

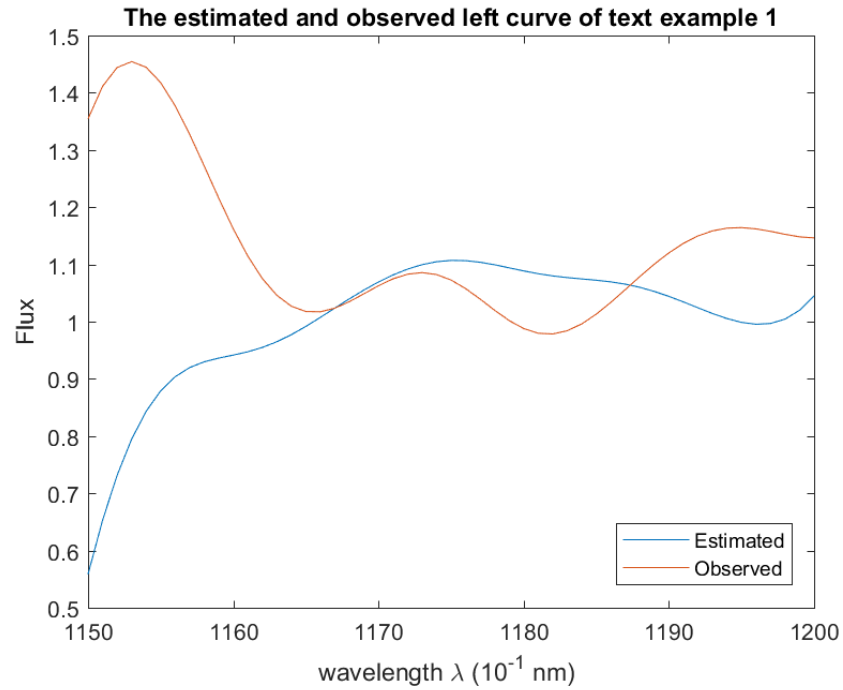
5(c), ii

The average training error is 2.3084.

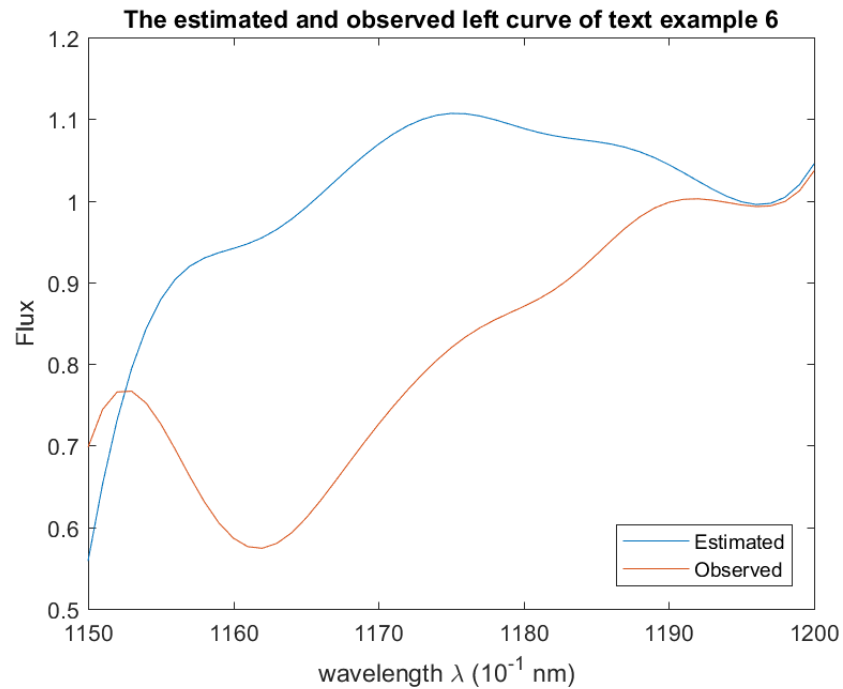
The average test error is 2.1080.

5(c), iii

Plot for test example 1:



Plot for test example 6:



Below is the matlab code for the section 5c,ii and iii.

```
close all; clear; clc;
load('smoothed');

%      %%% 5(c), ii - Function Estimator %%%
```

```

%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%TEST FOR FUNCTION%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%
%% test for getDistance method
% distance = getDistance(smoothed_qso_test(2,:), smoothed_qso_test(3,:), true);
% disp(distance);
10 %
%% test for getNeighbors and getDistances
% distances = getDistances(smoothed_qso_test, smoothed_qso_test(2,:));
% disp(distances);
% neighbors = getNeighbors(distances, 3);
15 % disp(neighbors);
%
%% test for the kerl function
% temp = ker(2);
% disp(temp);
20 %
% temp = ker(0.45);
% disp(temp);
%
%% test for estimateLeft method
25 % lambdasLeft = lambdas(1:51);
% lambdasRight = lambdas(151:end);
% f = smoothed_qso_train(3,:);
% fright = f(:,151:end);
% fleft = f(:,1: 51);
30 % fleftEstimated = estimateLeft(smoothed_qso_train, fright, 3);
% plot(lambdasLeft, fleftEstimated);
% hold on;
% plot(lambdasLeft, fleft);
% disp(getDistance(fleftEstimated, fleft, false));
35 %
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%END OF TEST%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%

% compute the training error
averageTrainError = getAverageTrainError(smoothed_qso_train, 3);
40 disp(averageTrainError);

% compute the test error
averageTestError = getAverageTestError(smoothed_qso_train, smoothed_qso_test, 3);
45 disp(averageTestError);

% plot for test example 1
figure
lambdasLeft = lambdas(1:51);
lambdasRight = lambdas(151:end);
50 f = smoothed_qso_train(1,:);
fright = f(:,151:end);
fleft = f(:,1: 51);
fleftEstimated = estimateLeft(smoothed_qso_train, fright, 3);
plot(lambdasLeft, fleftEstimated);
55 hold on;
plot(lambdasLeft, fleft);
title('The estimated and observed left curve of text example 1');
xlabel('wavelength \lambda (10^{-1} mm) ');

```

```

ylabel('Flux');
60 legend('Estimated','Observed','Location','southeast');
saveas(gcf,'5c(iii)-testEx1.png');

%plot for the test example 6
figure
65 lambdasLeft = lambdas(1:51);
lambdasRight = lambdas(151:end);
f = smoothed_qso_train(6,:);
fright = f(:,151:end);
fleft = f(:,1: 51);
70 fleftEstimated = estimateLeft(smoothed_qso_train, fright, 3);
plot(lambdasLeft, fleftEstimated);
hold on;
plot(lambdasLeft, fleft);
title('The estimated and observed left curve of text example 6');
75 xlabel('wavelength \lambda (10^{-1} mm) ');
ylabel('Flux');
legend('Estimated','Observed','Location','southeast');
saveas(gcf,'5c(iii)-testEx6.png');

80 function averageTestError = getAverageTestError(trainSet, testSet, k)
%%
% return the average test error for the entire test data set
% -'trainSet': the training data set
% -'testSet': the test data set
85 % -'k': the number of nearest neighbors we want to focus on
%%
% the test size
m = size(testSet, 1);
testSetRight = trainSet(:,151:end);
90 testSetLeft = trainSet(:,1: 51);
testEstimated = zeros(size(testSetLeft));
testErrors = zeros(size(testSet, 1), 1);
for i = 1: m
    testEstimated(i,:) = estimateLeft(trainSet, testSetRight(i,:), k);
95    testErrors(i,:) = getDistance(testSetLeft(i,:), testEstimated(i,:), false);
end
averageTestError = sum(testErrors) / m;
end

100 function averageTrainError = getAverageTrainError(trainSet, k)
%%
% return the estimated left function for the entire training right function;
% -'trainSet': the training data set;
% -'k': the number of nearest neighbors we want to focus on.
105 %%
trainSetRight = trainSet(:,151: end);
trainSetLeft = trainSet(:,1: 51);
trainEstimated = zeros(size(trainSetLeft));
errors = zeros(size(trainSet, 1), 1);
110 % the size of training set
m = size(trainSet, 1);

```

```

for i = 1:m
    trainEstimated(i,:) = estimateLeft(trainSet, trainSetRight(i,:), k);
    errors(i) = getDistance(trainSetLeft(i,:), trainEstimated(i,:), false);
115 end
averageTrainError = sum(errors) / m;
end

function fleft = estimateLeft(trainSet, fright, k)
120 %%
% return the estimated left function
% -'trainSet': the training data set
% -'fright': the observed right function
% -'k': the number of nearest neighbors
125 %%
% cut the training set into left and right part
trainSetRight = trainSet(:,151: end);
trainSetLeft = trainSet(:,1: 51);

130 distances = getDistances(trainSetRight, fright);
neighbors = getNeighbors(distances, k);
h = max(distances);

%compute the numerator and xc
135 numerator = zeros(1, size(trainSetLeft, 2));
denominator = 0;
for i = 1:k
    neighbor = neighbors(k);
    numerator = numerator + ker(distances(neighbor) / h) * trainSetLeft(i,:);
140 denominator = denominator + ker(distances(neighbor) / h);

end

fleft = numerator / denominator;
145 end

function temp = ker(t)
temp = max([1-t; 0]);
end

150 function neighbors = getNeighbors(distances, k)
%%
% return the K neighbors that is closet to our query function
%-'distances': is the vector containing distance of each training example
155 %-'K': the numbers of neighbors we want to find, must be a positive
%integers
%%
neighbors = zeros(k, 1);
for i = 1:k
160     [~, neighbors(i)] = min(distances, [], 'omitnan');
    % update the the min to NaN for the next closet neighbor
    distances(neighbors(i)) = NaN;
end
end

```

```
165     function distance = getDistance(f1, f2, right)
diff = (f1 - f2).^2;
if right == true
    % 151 is index of the wavelength 1300 among the lambdas vector
170     distance = sum(diff(151:end));
else
    % 51 is the index of the wavelength 1200 among the lambdas vector
    distance = sum(diff(1:51));
end
175 end

function distances = getDistances(trainSet, f)
%%
% return the all samples' distances relative to the function f.
180 %-'trainSet': the training set that contains both left and right examples
%-'f': the smoothed function
%%
diff = trainSet - f;
% only cumulate the right function distance
185 diff = diff(:,151:end).^2;
temp = cumsum(diff, 2);
distances = temp(:, end);
end
```