

# Report of homework02 IFT6390

## Team Member

Lifeng Wan (20108546), Qiang Ye (20139927)

## Coding Environment

python 3.5.2, numpy 1.14.2 matplotlib 2.2.0

---

## 1 Linear and non-linear regularized regression (50 points)

### 1.1 Linear Regression

For training dataset  $D_n$  with  $n$  samples (input, target):

$$D_n = \{(\mathbf{x}^{(1)}, t^{(1)}), \dots, (\mathbf{x}^{(n)}, t^{(n)})\}$$

with  $\mathbf{x}^{(i)} \in \mathbb{R}^d$  and  $t^{(i)} \in \mathbb{R}$ .

The linear regression assumes a parametrized form for the function  $f$  which predicts the value of the target from a new data point  $\mathbf{x}$ . (More precisely, it seeks to predict the expectation of the target variable conditioned on the input variable  $f(\mathbf{x}) \simeq \mathbb{E}[t|\mathbf{x}]$ .)

The parametrization is a linear transformation of the input, or more precisely an *affine* transformation.

$$f(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + b$$

1. Precise this model's set of parameters  $\theta$ , as well as the nature and dimensionality of each of them.

**Answer**  $\theta = \{\mathbf{w}, b\}$ , where  $\mathbf{w} \in \mathbb{R}^d$  is a weight vector and  $b \in \mathbb{R}$  is called as a bias.

---

2. The loss function typically used for linear regression is the quadratic loss:

$$L((\mathbf{x}, t), f) = (f(\mathbf{x}) - t)^2$$

We are now defining the **empirical risk**  $\hat{R}$  on the set  $D_n$  as the **sum** of the losses on this set (instead of the average of the losses as it is sometimes defined). Give the precise mathematical formula of this risk.

**Answer**

$$\begin{aligned} \hat{R}(f_\theta, \mathbb{D}_n) &= \sum_{i=1}^n L(\mathbf{x}^{(i)}, t^{(i)}, f_\theta) \\ &= \sum_{i=1}^n (f_\theta(\mathbf{x}^{(i)}) - t^{(i)})^2 \end{aligned}$$


---

3. Following the principle of Empirical Risk Minimization (ERM), we are going to seek the parameters which yield the smallest quadratic loss. Write a mathematical formulation of this minimization problem.

**Answer** Let  $\theta^*$  be the parameters which yield the smallest quadratic loss, then:

$$\theta^* = \underset{\theta}{\operatorname{argmin}} \hat{R}(f_{\theta}, \mathbb{D}_n)$$

4. A general algorithm for solving this optimization problem is gradient descent. Give a formula for the gradient of the empirical risk with respect to each parameter.

**Answer**

$$\begin{aligned} \frac{\partial \hat{R}(f_{\theta}, \mathbb{D}_n)}{\partial \theta} &= \frac{\partial \sum_{i=1}^n L(\mathbf{x}^{(i)}, t^{(i)}, f_{\theta})}{\partial \theta} \\ &= \sum_{i=1}^n \frac{\partial (f_{\theta}(\mathbf{x}^{(i)}) - t^{(i)})^2}{\partial \theta} \\ &= 2 \sum_{i=1}^n (f_{\theta}(\mathbf{x}^{(i)}) - t^{(i)}) \frac{\partial f_{\theta}(\mathbf{x}^{(i)})}{\partial \theta} \end{aligned}$$

Since:

$$\begin{aligned} \theta &= \{\mathbf{w}, b\} = \{w_1, w_2, \dots, w_k, \dots, w_d, b\} \\ f_{\theta}(\mathbf{x}) &= \mathbf{w}^T \mathbf{x} + b = w_1 x_1 + w_2 x_2 + \dots + w_d x_d + b \end{aligned}$$

then,

$$\frac{\partial \hat{R}}{\partial w_k} = 2 \sum_{i=1}^n (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - t^{(i)}) x_k^{(i)}$$

where  $1 \leq k \leq d$ , and

$$\frac{\partial \hat{R}}{\partial b} = 2 \sum_{i=1}^n (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - t^{(i)})$$

5. Define the error of the model on a single point  $(\mathbf{x}, t)$  by  $f(\mathbf{x}) - t$ . Explain in English the relationship between the empirical risk gradient and the errors on the training set.

**Answer** According to the definition, *error* is non-zero when output of the model  $f(\mathbf{x})$  is not equal to the label  $t$ . The gradient with respect to parameter  $b$  always is the sum of the every sample's error, whereas the gradient with respect to parameter  $w_k$  is the sum of the multiplication of each sample's error and corresponding  $x_k$ . The constant 2 can be ignored.

## 1.2 Ridge Regression

Instead of  $\hat{R}$ , we will now consider a **regularized empirical risk**:  $\tilde{R} = \hat{R} + \lambda \mathcal{L}(\theta)$ . Here  $\mathcal{L}$  takes the parameters  $\theta$  and returns a scalar penalty. This penalty is smaller for parameters for which we have an a priori preference. The scalar  $\lambda \geq 0$  is an **hyperparameter** that controls how much we favor minimizing the empirical risk versus this penalty. Note that we find the unregularized empirical risk when  $\lambda = 0$ .

We will consider a regularization called *Ridge*, or *weight decay* that penalizes the squared norm ( $l^2$ ) of the weights (but not the bias):  $L(\theta) = \|\mathbf{w}\|^2 = \sum_{k=1}^d \mathbf{w}_k^2$ . We want to minimize  $\tilde{R}$  rather than  $\hat{R}$ .

1. Express the gradient of  $\tilde{R}$ . How does it differ from the unregularized empirical risk gradient?

**Answer**

$$\begin{aligned}\frac{\partial \tilde{R}}{\partial \theta} &= \frac{\partial (\hat{R} + \lambda \mathcal{L}(\theta))}{\partial \theta} \\ &= \frac{\partial \hat{R}}{\partial \theta} + \lambda \frac{\partial \mathcal{L}(\theta)}{\partial \theta}\end{aligned}$$

For parameters  $w_k$ :

$$\begin{aligned}\frac{\partial \tilde{R}}{\partial w_k} &= \frac{\partial \hat{R}}{\partial w_k} + \lambda \frac{\partial \sum_{k=1}^d w_k^2}{\partial w_k} \\ &= 2 \sum_{i=1}^n (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - t^{(i)}) x_k^{(i)} + 2\lambda w_k\end{aligned}$$

For parameter  $b$ :

$$\begin{aligned}\frac{\partial \tilde{R}}{\partial b} &= \frac{\partial \hat{R}}{\partial b} + \lambda \frac{\partial \sum_{k=1}^d w_k^2}{\partial b} \\ &= \frac{\partial \hat{R}}{\partial b} \\ &= 2 \sum_{i=1}^n (f_{\mathbf{w}, b}(\mathbf{x}^{(i)}) - t^{(i)})\end{aligned}$$

When using **regularized empirical risk** the gradient with respect to  $w_k$  has an extra item  $2\lambda w_k$  compared to empirical risk; however, both gradients with respect to parameter  $b$  are same as the regularization has no constrain to the bias  $b$ .

- 
2. Write down a detailed pseudocode for the training algorithms that finds the optimal parameters minimizing  $\tilde{R}$  by gradient descent. To keep it simple, use a constant step-size  $\eta$ .

**Answer**

```
def training(training_data, _lambda, eta, max_steps):
    ...
    Input
        training_data: np.array (n * d+1) [X, t]
        _lambda: regularization facotr
        eta: learning rate
        max_steps: end condition
    Output: optimal parameters
        w: weight vector (d * 1)
        b: bias scalar
        (optional)losses: loss for each step, list
    ...
    get X, t from training dataset
    # X, t = training_data
    get sample size n and feature numbers d from X or t
    # n, d = X.shape[0], X.shape[1]-1
    initilize weight, bias randomly
    # w = np.random.normal(0, 0.001, (d,1))
    # b = 0.00
    initialize iter_numbers to 0
    # (optional) initialize losses = []
    while iter_numbers < max_steps:
        compute error = Xw + b - t
        # (optional) compute loss from error and w
        compute gradient with respect to w and b:
        # dw = 2 * (sum(error dot X) + _lambda * w)
        # db = 2 * sum(error)
        update weight and bias:
        # w = w - eta * dw
        # b = b - eta * db
        increase iter_numbers by 1:
        # iter_numbers += 1
        # (optional) store loss to losses: losses.append(loss)
    return w, b # (optional) losses
```

3. There happens to be an analytical solution to the minimization problem coming from linear regression (regularized or not). Assuming no bias (meaning  $b=0$ ), find a matrix formulation for the empirical risk and its

gradient, with the matrix  $\mathbf{X} = \begin{pmatrix} \mathbf{x}_1^{(1)} & \cdots & \mathbf{x}_d^{(1)} \\ \vdots & \ddots & \vdots \\ \mathbf{x}_1^{(n)} & \cdots & \mathbf{x}_d^{(n)} \end{pmatrix}$  and the vector  $\mathbf{t} = \begin{pmatrix} t^{(1)} \\ \vdots \\ t^{(n)} \end{pmatrix}$ .

**Answer** Given  $\mathbf{X}$  and  $\mathbf{t}$  as above, the empirical risk  $\tilde{R}$  can be expressed by the following equation:

$$\begin{aligned} \tilde{R} &= \hat{R} + \lambda \mathcal{L} \\ &= (\mathbf{X}\mathbf{w} - \mathbf{t})^T (\mathbf{X}\mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}^T \mathbf{w} \\ &= (\mathbf{w}^T \mathbf{X}^T - \mathbf{t}^T) (\mathbf{X}\mathbf{w} - \mathbf{t}) + \lambda \mathbf{w}^T \mathbf{w} \\ &= \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{w}^T \mathbf{X}^T \mathbf{t} - \mathbf{t}^T \mathbf{X} \mathbf{w} + \mathbf{t}^T \mathbf{t} + \lambda \mathbf{w}^T \mathbf{w} \end{aligned}$$

Notice that  $(\mathbf{w}^T \mathbf{X}^T \mathbf{t})^T = \mathbf{t}^T \mathbf{X} \mathbf{w}$ . Further notice that this is a 1x1 matrix, so  $\mathbf{w}^T \mathbf{X}^T \mathbf{t} = \mathbf{t}^T \mathbf{X} \mathbf{w}$ . Thus,

$$\tilde{R} = \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{w}^T \mathbf{X}^T \mathbf{t} + \mathbf{t}^T \mathbf{t} + \lambda \mathbf{w}^T \mathbf{w}$$

The gradient of  $\tilde{R}$  with respect to  $\mathbf{w}$  is:

$$\begin{aligned}
\nabla \tilde{R} &= \nabla \mathbf{w}^T \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \nabla \mathbf{w}^T \mathbf{X}^T \mathbf{t} + \nabla \mathbf{t}^T \mathbf{t} + \nabla \lambda \mathbf{w}^T \mathbf{w} \\
&= 2 \mathbf{X}^T \mathbf{X} \mathbf{w} - 2 \mathbf{X}^T \mathbf{t} + 0 + 2 \lambda \mathbf{w} \\
&= 2(\mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{t} + \lambda \mathbf{w})
\end{aligned}$$

4. Derive a matrix formulation of the analytical solution to the ridge regression minimization problem by expressing that the gradient is null at the optimum. What happens when  $N < d$  and  $\lambda = 0$ ?

**Answer** We set  $\tilde{R}$  in last question to zero at the optimum,  $\mathbf{w}^*$ :

$$\begin{aligned}
\mathbf{X}^T \mathbf{X} \mathbf{w}^* - \mathbf{X}^T \mathbf{t} + \lambda \mathbf{w}^* &= 0 \\
(\mathbf{X}^T \mathbf{X} + \lambda I) \mathbf{w}^* &= \mathbf{X}^T \mathbf{t} \\
\mathbf{w}^* &= (\mathbf{X}^T \mathbf{X} + \lambda I)^{-1} \mathbf{X}^T \mathbf{t}
\end{aligned}$$

If  $\lambda = 0$ :

$$\mathbf{w}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

In this situation, when  $N < d$ ,  $\mathbf{X}^T \mathbf{X}$  will be non-invertible, which means  $\mathbf{w}^*$  can not be analytically figured out. It also hints that we may have set too many features or these features are not completely independent.

### 1.3 Regression with a fixed non-linear pre-processing

We can make a non-linear regression algorithm by first passing the data through a fixed non-linear filter: a function  $\phi(\mathbf{x})$  that maps  $\mathbf{x}$  non-linearly to a higher dimensional  $\tilde{\mathbf{x}}$ .

For instance, if  $x \in \mathbb{R}$  is one dimensional, we can use the polynomial transformation:

$$\tilde{x} = \phi_{poly\ k}(x) = \begin{pmatrix} x \\ x^2 \\ \vdots \\ x^k \end{pmatrix}$$

We can then train a regression, not on the  $(x^{(i)}, t^{(i)})$  from the initial training set  $D_n$ , but on the transformed data  $(\phi(x^{(i)}), t^{(i)})$ . This training finds the parameters of an affine transformation  $f$ .

To predict the target for a new training point  $x$ , you won't use  $f(x)$  but  $\tilde{f}(x) = f(\phi(x))$ .

1. Write the detailed expression of  $\tilde{f}(x)$  when  $x$  is one-dimensional (uni-variate) and we use  $\phi = \phi_{poly\ k}$ .

**Answer**

$$\begin{aligned}
\tilde{f}(x) &= f(\phi_{poly\ k}(x)) \\
&= \mathbf{w}^T \phi_{poly\ k}(x) + b \\
&= w_1 x + w_2 x^2 + \dots + w_k x^k + b
\end{aligned}$$

2. Give a detailed explanation of the parameters and their dimensions.

**Answer** The parameters are  $\theta = \{\mathbf{w}, b\}$ , where  $\mathbf{w} \in \mathbb{R}^k$  is a  $k$  dimensional column vector with each  $w_k$  be a coefficient of  $x^k$ , and  $b \in \mathbb{R}$  is a scale (one dimensional vector) called bias as seen in Question 1.1.

3. If dimension  $d \geq 2$ , a polynomial transformation should include not only the individual variable exponents  $x_i^j$ , for powers  $j \leq k$ , and variables  $i \leq d$ , but also all the interaction terms of order  $k$  and less between several variables (e.g. terms like  $x_i^{j_1} x_l^{j_2}$ , for  $j_1 + j_2 \leq k$  and variables  $i, l \leq d$ ). For  $d = 2$ , write down as a

function of each of the 2 components of  $x$  the transformation  $\phi_{poly\ 1}(x)$ ,  $\phi_{poly\ 2}(x)$ , and  $\phi_{poly\ 3}(x)$ .

**Answer**

$$\begin{aligned}\phi_{poly\ 1}(x) &= \begin{pmatrix} x_1 \\ x_2 \end{pmatrix} \\ \phi_{poly\ 2}(x) &= \begin{pmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \end{pmatrix} \\ \phi_{poly\ 3}(x) &= \begin{pmatrix} x_1 \\ x_2 \\ x_1 x_2 \\ x_1^2 \\ x_2^2 \\ x_1^2 x_2 \\ x_1 x_2^2 \\ x_1^3 \\ x_2^3 \end{pmatrix}\end{aligned}$$

4. What is the dimensionality of  $\phi_{poly\ k}(x)$ , as a function of  $d$  and  $k$ ?

**Answer** Let  $dim(d, k)$  be the function of the dimensionality of  $\phi_{poly\ k}(x)$  with respect to  $d$  and  $k$ , then:

$$dim(d, k) = \sum_{i=1}^k \frac{(i + d - 1)!}{i! (d - 1)!}$$

## 2 Practical Part (50 points)

You should include all the python files you used to get your results. It should have a main file (which can be a notebook) that produces the required plots, one after another. Your results should be reproducible! Briefly explain how to run your code in the report.

1. Implement in python the ridge regression with gradient descent. We will call this algorithm `regression_gradient`. Note that we now have parameters  $\mathbf{w}$  and  $b$  we want to learn on the training set, as well an *hyper-parameter* to control the capacity of our model:  $\lambda$ . There are also *hyper-parameters* for the optimization: the step-size  $\eta$ , and potentially the number of steps.
2. Consider the function  $h(x) = \sin(x) + 0.3x - 1$ . Draw a dataset  $D_n$  of pairs  $(x, h(x))$  with  $n = 15$  points where  $x$  is drawn uniformly at random in the interval  $[-5, 5]$ . Make sure to use the same set  $D_n$  for **all** the plots below.
3. With  $\lambda = 0$ , train your model on  $D_n$  with the algorithm `regression_gradient`. Then plot on the interval  $[-10, 10]$ : the points from the training set  $D_n$ , the curve  $h(x)$ , and the curve of the function learned by your model using gradient descent. Make a clean legend. **Remark:** The solution you found with gradient descent should converge to the straight line that is closer from the  $n$  points (and also to the analytical solution). Be ready to adjust your step-size (small enough) and number of iterations (large enough) to reach this result.
4. on the same graph, add the predictions you get for intermediate value of  $\lambda$ , and for a large value of  $\lambda$ . Your plot should include the value of  $\lambda$  in the legend. It should illustrate qualitatively what happens when  $\lambda$

increases.

5. Draw another dataset  $D_{test}$  of 100 points by following the same procedure as  $D_n$ . Train your linear model on  $D_n$  for  $\lambda$  taking values in  $[0.0001, 0.001, 0.01, 0.1, 1, 10, 100]$ . For each value of  $\lambda$ , measure the average quadratic loss on  $D_{test}$ . Report these values on a graph with  $\lambda$  on the x-axis and the loss value on the y-axis.
6. Use the technique studied in problem 1.3 above to learn a non-linear function of  $x$ . Specifically, use Ridge regression with the fixed preprocessing  $\phi_{poly}^l$  described above to get a polynomial regression of order  $l$ . Apply this technique with  $\lambda = 0.01$  and different values of  $l$ . Plot a graph similar to question 2.2 with all the prediction functions you got. Don't plot too many functions to keep it readable and precise the value of  $l$  in the legend.
7. Comment on what happens when  $l$  increases. What happens to the empirical risk (loss on  $D_n$ ), and to the true risk (loss on  $D_{test}$ )?

## Answers

First, import numpy and matplotlib.pyplot

In [1]:

```
1 %matplotlib inline
2 import numpy as np
3 import matplotlib.pyplot as plt
4 np.random.seed(0)
```

### 1. Implementation of the function `regression_gradient`

Explanations are detailed given as comments in the codes.

In [2]:

```

1 def regression_gradient(data,
2                         _lambda = 0,
3                         eta = 1e-4,
4                         max_steps = 1e5,
5                         epsilon = 1e-3
6                         ):
7     '''ridge regression with gradient descent.
8     params
9         data: training dataset n rows and d+1 columns with last column the
10         true output t np.array (n, d+1)
11         _lambda: regularization factor, float
12         eta: learning rate / step size, float
13         max_steps: maximal steps of gradient descent, int
14         epsilon: condition for ending iteration, float
15     returns
16         w weight vector (d, 1)
17         b bias scalar
18     ...
19     n, d = data.shape[0], data.shape[1]-1
20     X = data[:, :-1].reshape(n, d) # (n, d)
21     t = data[:, -1].reshape(n, 1) # (n, 1)
22     w, b = np.random.normal(0, 1e-4, (d, 1)), 0.0
23     max_steps = int(max(1, max_steps))
24     losses = [] # store loss of each iteration
25     for i in range(max_steps):
26         dw, db = np.zeros_like(w), 0.0 # set dw, db to zero
27         error = np.dot(X, w) + b - t # (n, d)*(d, 1)->(n, 1)
28         if i % 1e4 == 0:
29             loss = (float(np.dot(error.T, error)) + float(np.dot(w.T, w)))/n
30             losses.append(loss)
31         if loss <= epsilon: # stop iteration
32             return w, b, losses
33         error_mul_x = np.multiply(error, X) # (n, 1) mul (n, d)
34         error_mul_x = np.clip(error_mul_x, -1e8, 1e8) # avoid overflow of value
35         # compute dw and db
36         dw = 2 * (np.sum(error_mul_x, axis = 0).reshape(d, 1) + _lambda * w)
37         db = 2 * np.float(np.sum(error, axis = 0))
38         # update weight and bias
39         w -= eta * dw
40         b -= eta * db
41     return w, b, losses

```

## 2. Create and Plot training data points

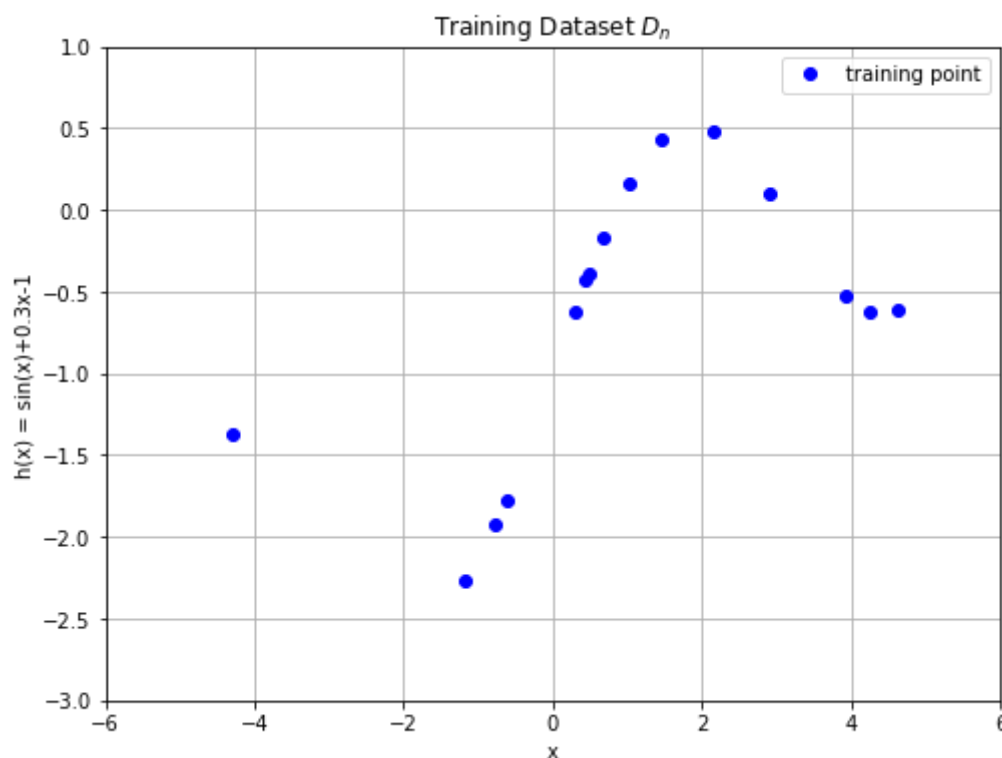


In [3]:

```

1 def h_x(x): # true mapping of x and y
2     return np.sin(x) + 0.3 * x - 1
3 # generate X uniformly at random in interval [-5, 5]
4 X = np.random.uniform(low = -5, high = 5, size = (15, 1))
5 t = h_x(X)
6 # plot training data points
7 plt.figure(figsize=(8, 6))
8 plt.plot(X, t, 'bo')
9 plt.grid(True) # add a grid
10 plt.axis([-6, 6, -3, 1]) # restriction to axes
11 plt.xlabel('x')
12 plt.ylabel('h(x) = sin(x)+0.3x-1')
13 plt.legend(['training point'], loc='upper right')
14 plt.title("Training Dataset $D_n$")
15 plt.show()

```



### 3. Plot training data points, true curve, and curve learned by `regression_descent`

We first prepared the data for the function `regression_gradient`; we executed the function to train our model and got the parameters `w` and `b`. Function `f_x` is to compute the output a linear regression with input data `X`, the parameters `w` and `b`. Based on all these data, we plotted the curves demanded. Read the comments of the code for more details.

We also wrote a function: `analytical_solution` trying to find analytical solution of a ridge regression from training data. We found that the curves drawn with points from `regression_gradient` and `analytical_solution` are completely overlapped with each other, indicating that the parameters from gradient descent have very high confidence to be correct.

In [4]:

```

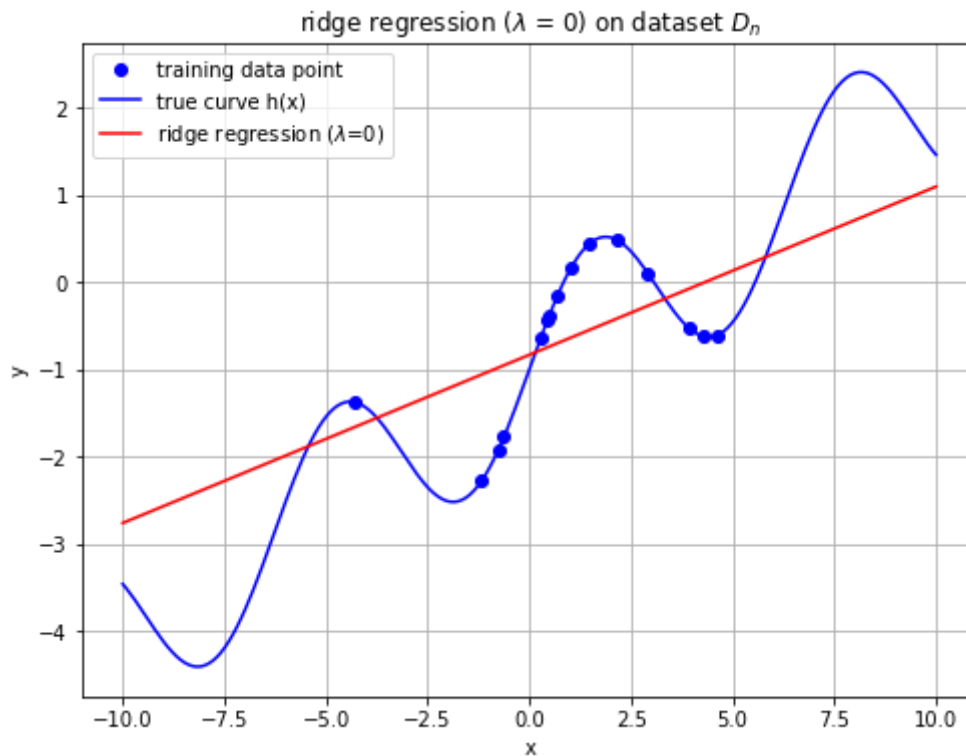
1 data = np.concatenate((X, t), axis = 1) # prepare training data
2 # training model to get parameters
3 w, b, step_losses = regression_gradient(data,
4                                     _lambda = 0.0, # no weight penalty
5                                     eta = 1e-4, # learning rate
6                                     max_steps = 1e4 # max numbers of iterat
7                                     )
8
9 def analytical_solution(data, _lambda = 0.0):
10     """find analytical solution of ridge regression from data
11     params
12         data: training dataset n rows and d+1 columns with last column the true
13         output t np.array (n, d+1)
14         _lambda: regularization factor, float
15     returns
16         w weight vector (d, 1)
17         b bias scalar
18     """
19     n, d = data.shape
20     d -= 1 # the last column is label t
21     X = data[:, :-1].reshape(n, d) # (n, d)
22     t = data[:, -1].reshape(n, 1) # (n, 1)
23     X = np.concatenate((np.ones((n, 1))), X, axis = 1)
24     I = np.identity(d+1)
25     temp = np.mat(np.dot(X.T, X) + _lambda * I) # (d+1, d+1)
26     temp = np.linalg.inv(temp) # inverse of temp
27     theta = np.dot(temp, X.T) # (d+1, d+1) dot (d+1, n) -> (d+1, n)
28     theta = np.dot(theta, t).reshape(-1, 1) # (d+1, n) dot (n, 1) -> (d+1, 1)
29     b = float(theta[0, 0])
30     w = theta[1:, :].reshape(-1, 1) # (d, 1)
31     return w, b
32
33 # w_1, b_1 = analytical_solution(data)
34
35 def f_x(X, w, b):
36     """generate output of a linear regression model by parameters w and b
37     params
38         X: variables, np.array (n, d)
39         w: weights, np.array (d, 1)
40         b: bias, float scalar
41     returns
42         y: np.dot(X, w) + b
43     """
44     return np.dot(X, w) + b
45
46 support = np.linspace(-10, 10, 200).reshape(-1, 1)
47 y = f_x(support, w, b) # ridge regression with _lambda = 0.00
48 h = h_x(support) # true output
49 # analytical_y = f_x(support, w_1, b_1)
50
51 plt.figure(figsize=(8, 6))
52 plt.grid(True) # add a grid
53 plt.plot(X, t, 'bo')
54 plt.plot(support, h, 'b-')
55 # plt.plot(support, analytical_y, 'g-')
56 plt.plot(support, y, 'r-')
57 plt.xlabel('x')
58 plt.ylabel('y')
59 plt.legend(['training data point',

```

```

60     'true curve h(x)',
61     # 'analytical solution',
62     'ridge regression ($\lambda=0$)'
63 ], loc='upper left')
64 plt.title("ridge regression ($\lambda = 0$) on dataset $D_n$")
65 plt.show()

```



In [5]:

```

1  def draw_step_loss(step_losses):
2      steps = np.array([i for i in range(len(step_losses))])
3      step_losses = np.array(step_losses).reshape(-1, 1)
4      plt.figure(figsize=(8, 6))
5      plt.grid(True) # add a grid
6      plt.plot(steps, step_losses, 'b-')
7      plt.xlabel('steps')
8      plt.ylabel('loss')
9      plt.legend(['loss'], loc='upper left')
10     plt.title("loss of ridge regression ($\lambda = 0$) on dataset $D_n$")
11     plt.show()
12
13 # to see the loss trends trough the iteration, uncomment next line
14 # draw_step_loss(step_losses)

```

#### 4. Plot curves from ridge regression with an intermediate and a large $\lambda$

In our codes below, we use 100 and 1000 for intermediate and large  $\lambda$  respectively. Run `regression_gradient` for each  $\lambda$  to get corresponding optimum parameters. With these parameters and the function `f_x` implemented in last question, we can achieve different outputs of the model for plotting.

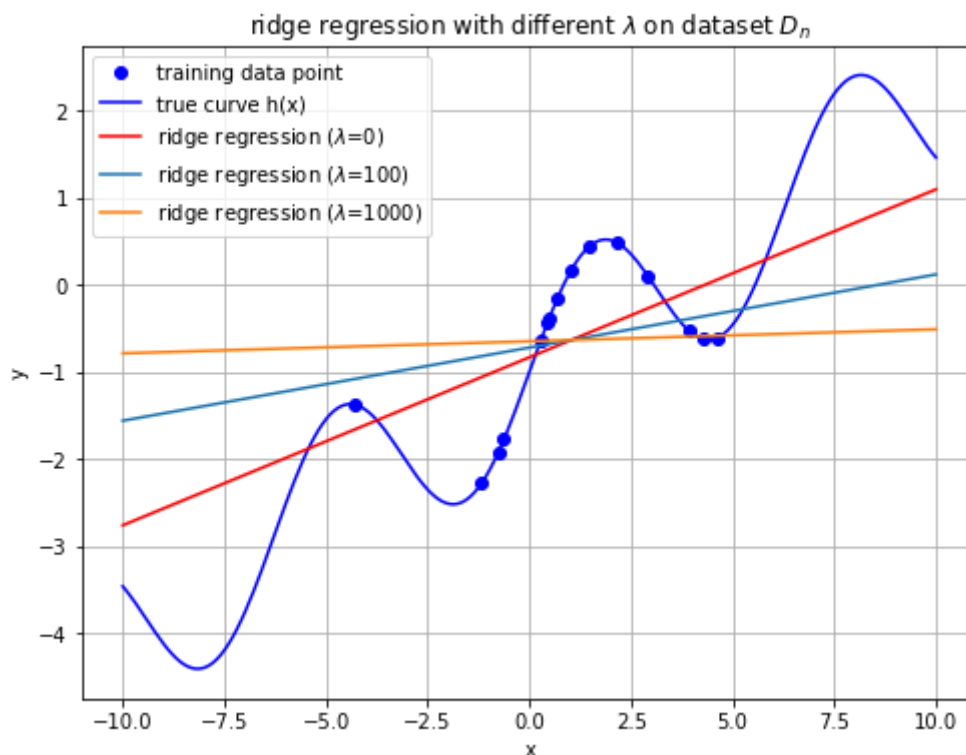
See the comments in the codes for more detail.

In [6]:

```

1 _lambdas = [100, 1000]
2 params = [] # keep parameters (w, b) for each _lambdas
3 for _lambda in _lambdas:
4     w, b, _ = regression_gradient(data, _lambda = _lambda)
5     params.append((w, b))
6
7 plt.figure(figsize=(8, 6))
8 plt.grid(True) # add a grid
9 plt.plot(X, t, 'bo') # training data points
10 plt.plot(support, h, 'b-') # true curve
11 plt.plot(support, y, 'r-') # ridge regression with _lambda = 0.
12
13 legends = ['training data point',
14            'true curve h(x)',
15            'ridge regression ($\lambda=0$)']
16
17 for i in range(len(params)): # plot curves with different _lambda
18     plt.plot(support, f_x(support, params[i][0], params[i][1]))
19     legends.append('ridge regression ($\lambda=${})'.format(_lambdas[i]))
20
21 plt.xlabel('x')
22 plt.ylabel('y')
23 plt.legend(legends, loc='upper left')
24 plt.title("ridge regression with different $\lambda$ on dataset $D_n$")
25 plt.show()

```



## 5. Plot loss of different $\lambda$ s

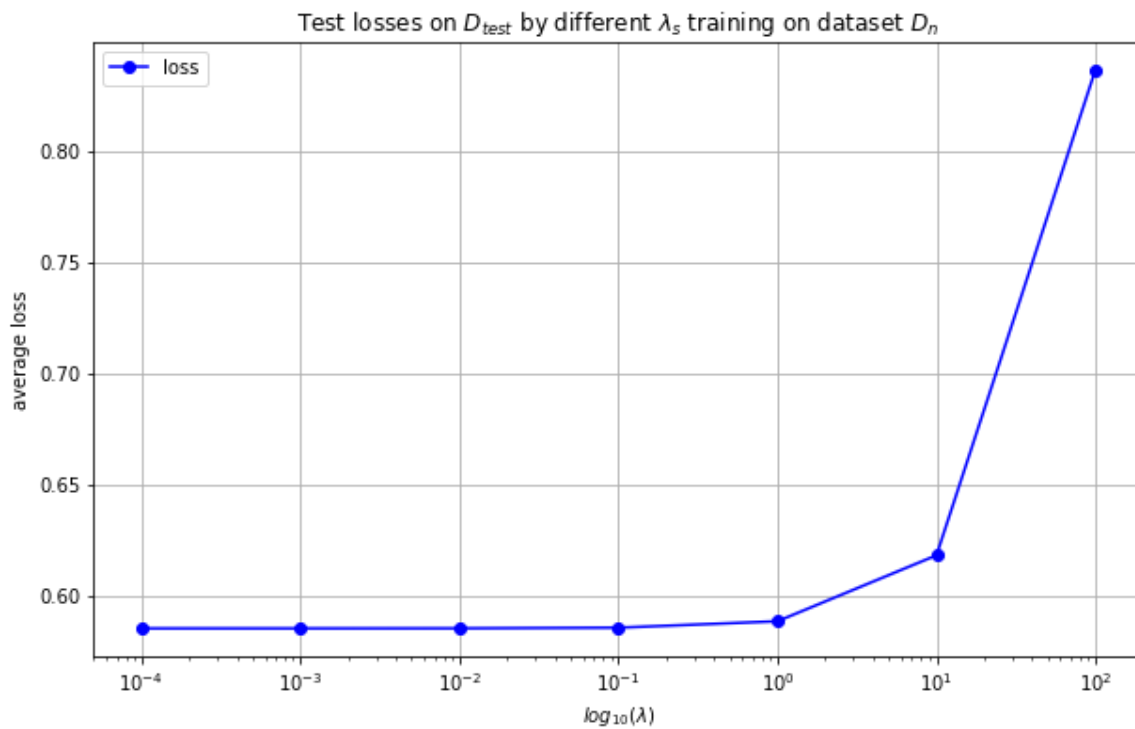
In this section, test dataset is generated based on the requirement. A function to compute the average loss receiving output of a model and true output as main parameters is also implemented. After completing each training on certain  $\lambda$ . We also calculated and collected the average loss on testing data set and stored them in a list `losses`. See the comments in the codes for more detail.

In [7]:

```

1  # Prepare test data (100 data points) from interval [-5, 5]
2  X_test = np.random.uniform(-5, 5, (100, 1))
3  t_test = h_x(X_test)
4  data_test = np.concatenate((X_test, t_test), axis = 1)
5  # using given different lambdas
6  _lambdas = [0.0001, 0.001, 0.01, 0.1, 1, 10, 100]
7
8  def compute_average_loss(y_from_model, true_y, _lambda = 0.0, w = None):
9      """compute loss, if _lambda is not zero and w is not None, the weight
10         penalty will be counted to the total loss.
11         params
12             y_from_model: output generated by a ridge regression model
13             true_y: true output
14             _lambda: factor of weight penalty
15             w: weight matrix(vector)
16         returns
17             average loss of the model
18         """
19         assert y_from_model.shape == true_y.shape, "shape mismatch."
20         loss = np.sum(np.power(y_from_model - true_y, 2))
21         if not (_lambda == 0.0 or w is None):
22             loss += _lambda * np.dot(w.T, w)
23         return float(loss) / len(true_y)
24
25     losses = [] # store losses with different _lambda
26     for _lambda in _lambdas:
27         # training model with different _lambda
28         w, b, _ = regression_gradient(data,
29                                     _lambda = _lambda,
30                                     eta = 1e-5,
31                                     max_steps = 1e5
32                                     )
33         y_from_model = f_x(X_test, w, b)
34         loss = compute_average_loss(y_from_model, t_test)
35         losses.append(loss)
36
37     # plot loss
38     plt.figure(figsize=(10, 6))
39     plt.grid(True) # add a grid
40     plt.xscale('log')
41     plt.plot(_lambdas, losses, 'bo-')
42     plt.xlabel('$\lambda$')
43     plt.ylabel('average loss')
44     plt.legend(['loss'], loc='upper left')
45     plt.title("Test losses on $D_{test}$ by different $\lambda_s$ training on datas")
46     plt.show()

```



## 6. Ridge regression on different polynomial pre-processed training datasets

Before combining the polynomial data pre-processing and ridge regression, we first implemented a function to generate dataset with  $k$  features from original data  $X$ ; each of the features is  $i$  ( $i \leq k$ ) times power of the original data  $X$ .

After that, function `regression_gradient` is executed for training dataset with different  $k$ . Training loss and optimal parameters are stored for further computation or analysis.

It may take several seconds or minutes to run the following codes as the `max_steps` here is set to  $1e6$ .

See comments in the following codes for more detail.

In [8]:

```

1 def polynomialize(data, k):
2     """generate a np.array with shape:(n, k). the value of ith column is i time
3     the first column value in the same row
4     param
5         data: original data, np.array (n, 1)
6         k: max power k >= 1
7     returns
8         new_data np.array (n, k)
9     """
10    data = data.reshape(-1, 1)
11    if k < 2:
12        return data
13    new_data = data.copy()
14    for i in range(2, k+1):
15        new_data = np.concatenate((new_data, np.power(data, i)), axis = 1)
16    return new_data
17
18    _lambda, k = 0.01, 5
19    X_polynomial = polynomialize(X, k) # polynomialized training dataset
20    support_polynomial = polynomialize(support, k) # for plotting
21
22    parameters = []
23    training_losses = []
24
25    for l_order in range(1, k+1): # different l_order
26        data = np.concatenate((X_polynomial[:, :l_order], t), axis = 1)
27        eta = pow(10, -3 - l_order) # different learning rate for different order
28        w, b, losses = regression_gradient(data,
29                                          _lambda = _lambda,
30                                          eta = eta,
31                                          max_steps = 1e6)
32        parameters.append((w, b))
33        training_losses.append(losses[-1]) # store the most recent loss

```

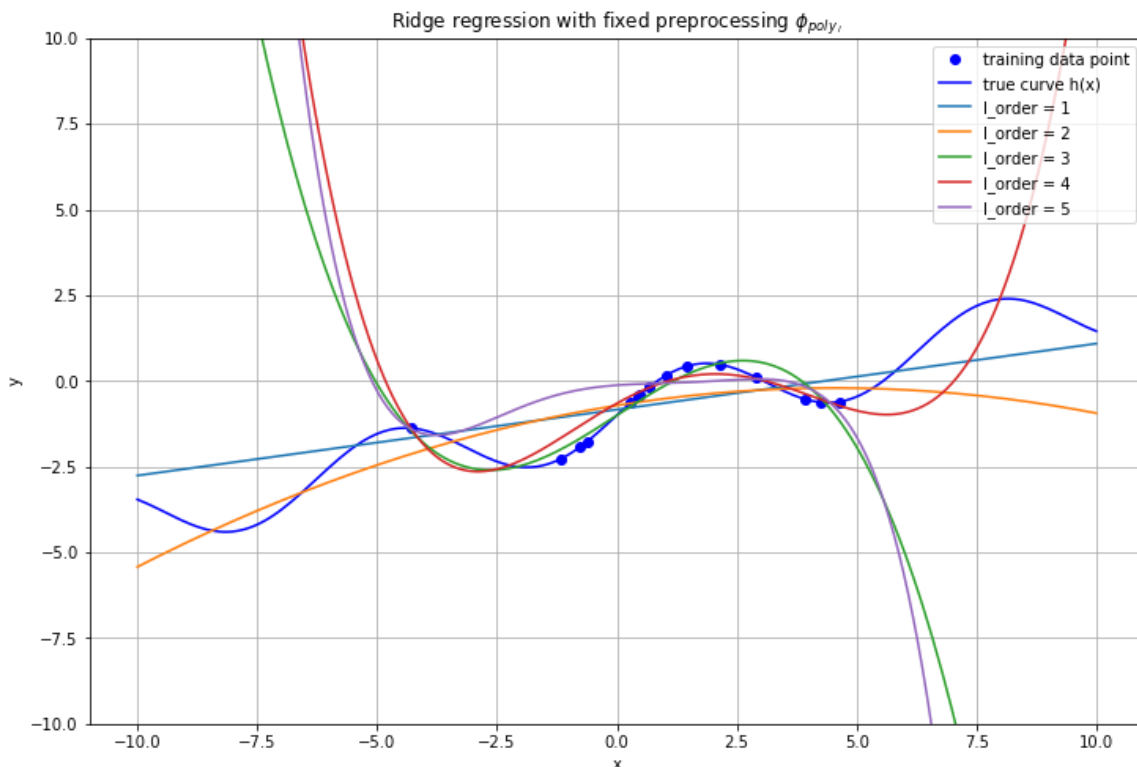
Plot ridge regression on different polynomial pre-processed dataset:

In [9]:

```

1 # plotting figure
2 plt.figure(figsize=(12, 8))
3 plt.grid(True) # add a grid
4 plt.plot(X, t, 'bo')
5 plt.plot(support, h_x(support), 'b-')
6
7 legends = ['training data point',
8           'true curve h(x)']
9
10 for l_order in range(1, k+1):
11     w, b = parameters[l_order-1]
12     plt.plot(support, f_x(support_polynomial[:, :l_order], w, b))
13     legends.append('l_order = {}'.format(l_order))
14
15 plt.xlabel('x')
16 plt.ylabel('y')
17 plt.axis([None, None, -10, 10])
18 plt.legend(legends, loc='upper right')
19 plt.title("Ridge regression with fixed preprocessing  $\phi_{poly_l}$ ")
20 plt.show()

```



## 7. Comment on what happens when $l$ increases.

Usually, When  $l$  increases, the ridge regression with polynomial pre-processing will have higher power components of the original data  $\mathbf{X}$ ; thus, the capacity of the model will grow, and it will try its best to fit the training data in the interval  $[-5, 5]$  rather than find the inner true mapping of  $h(x)$  in both  $[-5, 5]$  and  $(-\infty, \infty)$ . Therefore, as  $l$  increase, the trend is that losses on training set will decrease as long as the model is well tuned and fully trained (with adequate iteration), whereas on testing set (even if data is sampled in  $[-5, 5]$ ), loss will have very very high opportunity to increase.

In practice, due to non-perfect training and randomized sampling of dataset on interval  $[-5, 5]$ , losses on training set may not always decrease as  $l$  increases; loss on testing set may have a local minimum at certain  $l$ .



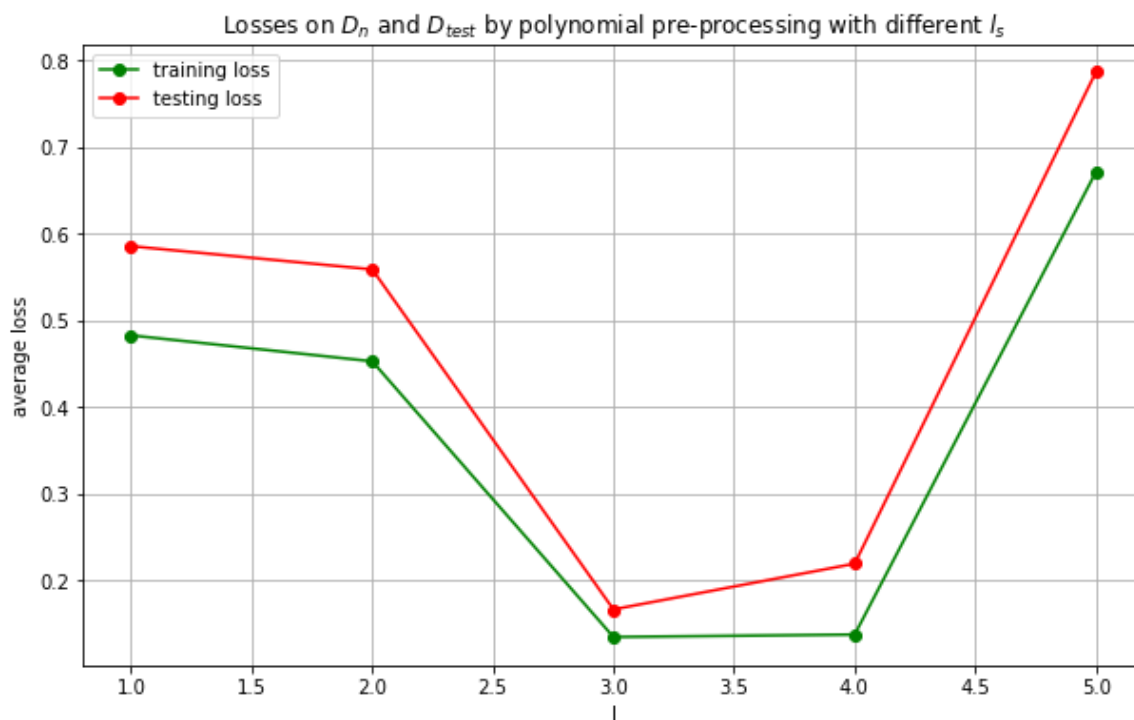
The following code computed and plotted the losses of our trained model on training set and testing set.

In [10]:

```

1 X_test_polynomial = polynomialize(X_test, k)
2 testing_losses = []
3 for l_order in range(1, k+1):
4     w, b = parameters[l_order - 1]
5     y_test = f_x(X_test_polynomial[:, :l_order], w, b)
6     test_loss = compute_average_loss(y_test, t_test)
7     testing_losses.append(test_loss)
8
9 l_orders = [i for i in range(1, k+1)]
10
11 plt.figure(figsize=(10, 6))
12 plt.grid(True) # add a grid
13 plt.plot(l_orders, training_losses, 'go-')
14 plt.plot(l_orders, testing_losses, 'ro-')
15 plt.xlabel('l')
16 plt.ylabel('average loss')
17 plt.legend(['training loss', 'testing loss'], loc='upper left')
18 plt.title("Losses on  $D_n$  and  $D_{test}$  by polynomial pre-processing with dif
19 plt.show()

```



**The end of the report.**

### Team Member

Lifeng Wan (20108546), Qiang Ye (20139927)

### Coding Environment

python 3.5.2, numpy 1.14.2 matplotlib 2.2.0

