

Project 2 Neural Network Optimization

(a) Using the sigmoid function defined above, compute an expression for the gradient $\nabla_w f_w(x)$. You can leave the expression in terms of ϕ and the derivative of ϕ . *

$$\nabla_w f_w(x) = [\phi(w_2x_1 + w_3x_2 + w_4x_3 + w_5)$$

$$x_1w_1\phi'(w_2x_1 + w_3x_2 + w_4x_3 + w_5)$$

$$x_2w_1\phi'(w_2x_1 + w_3x_2 + w_4x_3 + w_5)$$

$$x_3w_1\phi'(w_2x_1 + w_3x_2 + w_4x_3 + w_5)$$

$$w_1\phi'(w_2x_1 + w_3x_2 + w_4x_3 + w_5)$$

$$\phi(w_7x_1 + w_8x_2 + w_9x_3 + w_{10})$$

$$x_1w_6\phi'(w_7x_1 + w_8x_2 + w_9x_3 + w_{10})$$

$$x_2w_6\phi'(w_7x_1 + w_8x_2 + w_9x_3 + w_{10})$$

$$x_3w_6\phi'(w_7x_1 + w_8x_2 + w_9x_3 + w_{10})$$

$$w_6\phi'(w_7x_1 + w_8x_2 + w_9x_3 + w_{10})$$

$$\phi(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15})$$

$$x_1w_{11}\phi'(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15})$$

$$x_2w_{11}\phi'(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15})$$

$$x_3w_{11}\phi'(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15})$$

$$w_{11}\phi'(w_{12}x_1 + w_{13}x_2 + w_{14}x_3 + w_{15})$$

1]

(b) Compute the Derivative matrix $D_r(w)$ using the gradient from previous part (read page 383 of the textbook). Notice that $r(w) = [r_1(w), \dots, r_N(w)]$ where $r_n(w)$ is defined in (??).

```
%See function Jacob_f
w = 0.05*[1:16];
J = Jacob_f(x,w);
```

(c) We will train the neural network to approximate the non-linear function $g(x) = x_1x_2 + x_3$. To this end, generate $N = 500$ random points $x^{(n)} = [x^{(n)}_1, x^{(n)}_2, x^{(n)}_3]^T$, $1 \leq n \leq N$ in \mathbb{R}^3 and assign $y^{(n)} = x^{(n)}_1 x^{(n)}_2 + x^{(n)}_3$.

```
%generate 500 data points
%range of data is 0 to 1
[x,y] = rand_gen(500);
```

Use these $(x(n), y(n))$ N $n=1$ pairs to train the neural network by minimizing the following function (aka training loss) with respect to w :
$$l(w) = \sum_{n=1}^N r^2_n(w) + \lambda \|w\|^2_2$$

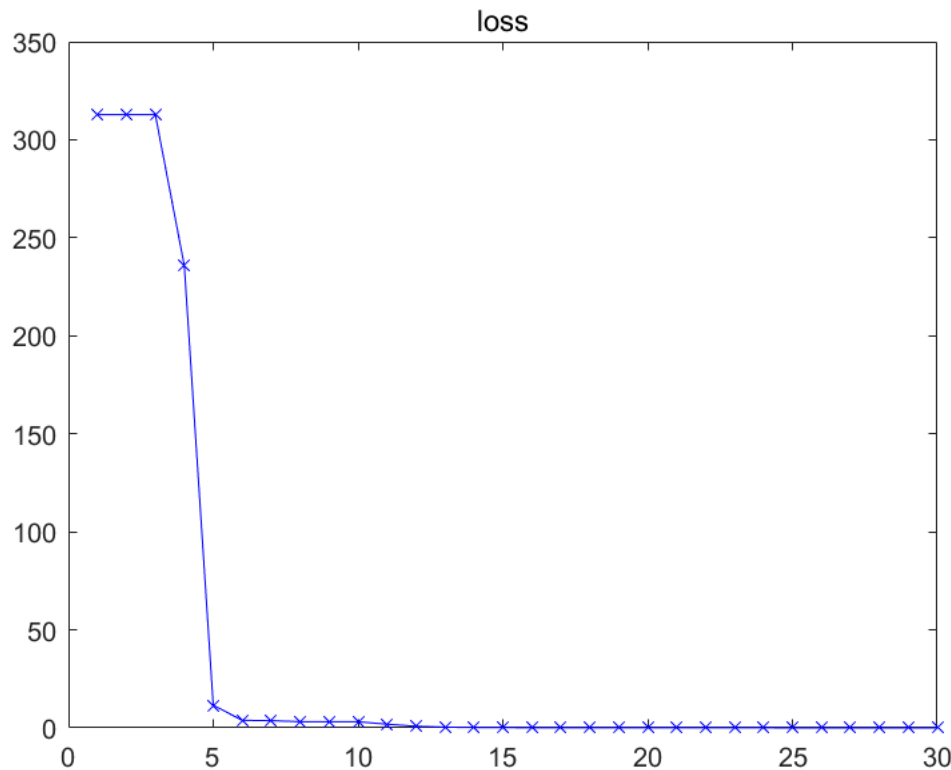
i. Write your own code implementing the Levenberg-Marquardt algorithm to approximately solve the above non-linear least squares problem. You

ur code should take $(x(n), y(n))$ N $n=1$ as inputs and produce the learned w as the output. Choose a very small value of λ (say $\lambda = 10^{-5}$) and experiment with a few choices of λ . Choose a suitable stopping criterion (Read Page 393 of the textbook) and specify it in your report. Plot the value of training loss $l(w)$ versus iterations to show how your training loss decreases with iterations.

My stopping criterion running 30 loops. Since I stacked my result for iterations, I have to fix the number of iterations in the multiple run algorithm.

```
%When lemda = 10^-5, lemda_k = 1
w1 = Levbg_Maqrt_alg(x,y);%multipule run algorithm
```

```
final loss=0.007660
error=0.000015
```

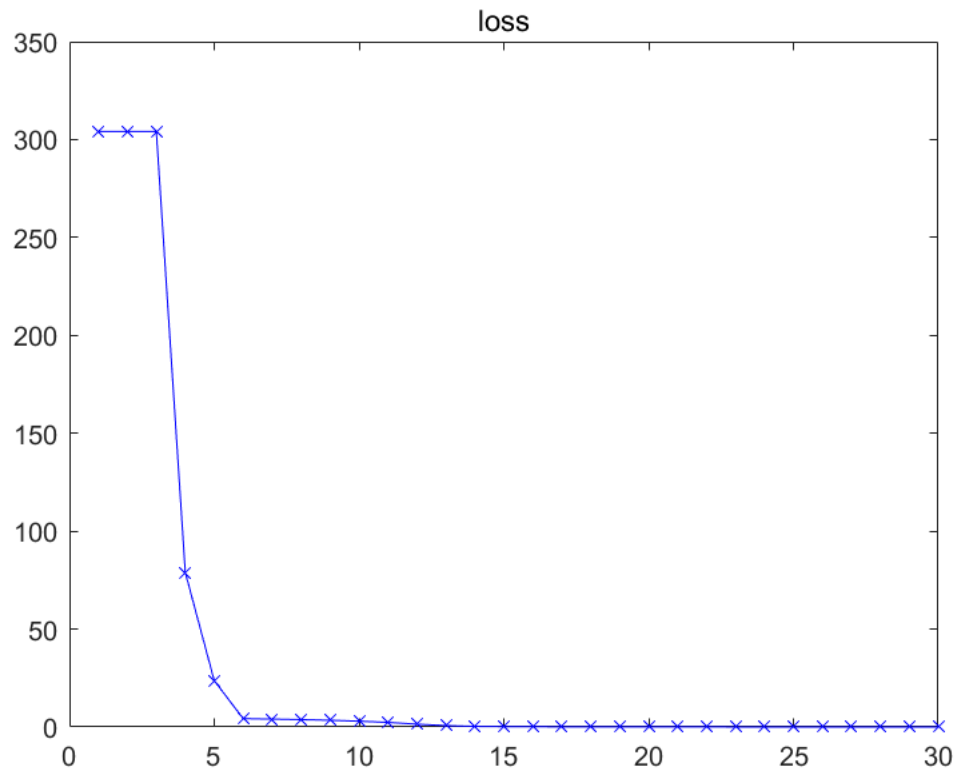


```
%error = 0.000015
```

When lemda = 10^{-3} , there is no significant impact of training loss.

```
%When lemda = 10^-3, lemda_k = 1
w2 = Levbg_Maqrt_alg(x,y);
```

```
final loss=0.024129
error=0.000048
```

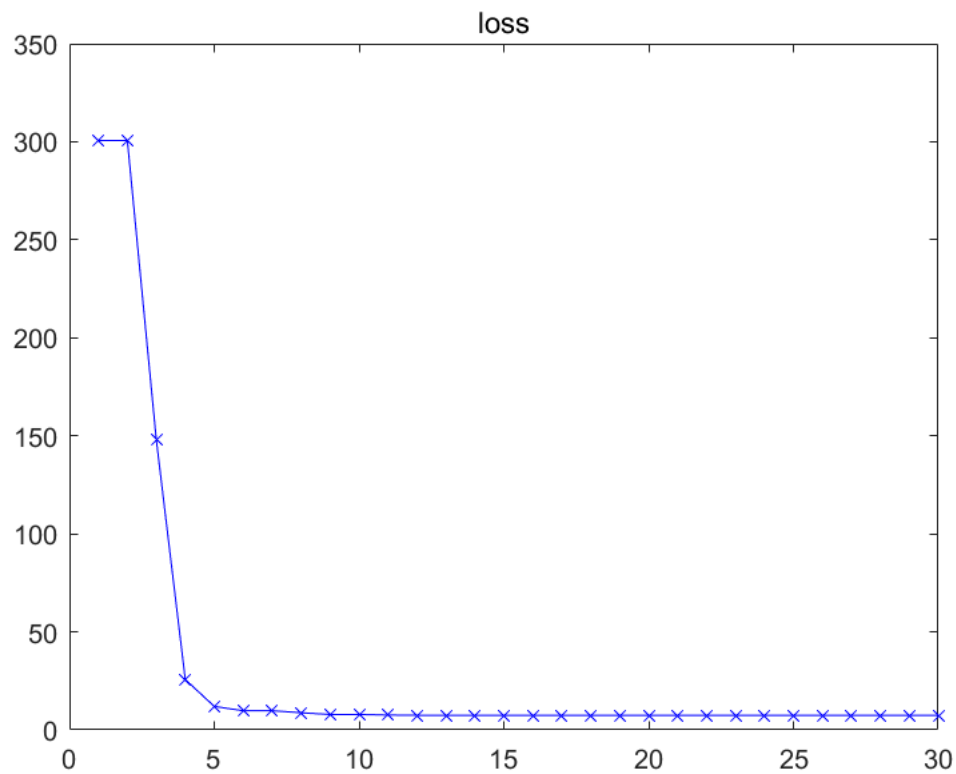


```
%error = 0.000048
```

When $\text{lemda} = 1$, the final training loss is larger than before since the training loss put more weight on $\text{norm}(w)$ instead of $\text{norm}(r)$.

```
%When lemda = 1, lemda_k = 1
w3 = Levbg_Maqrt_alg(x,y);
```

```
final loss=7.373758
error=0.014748
```

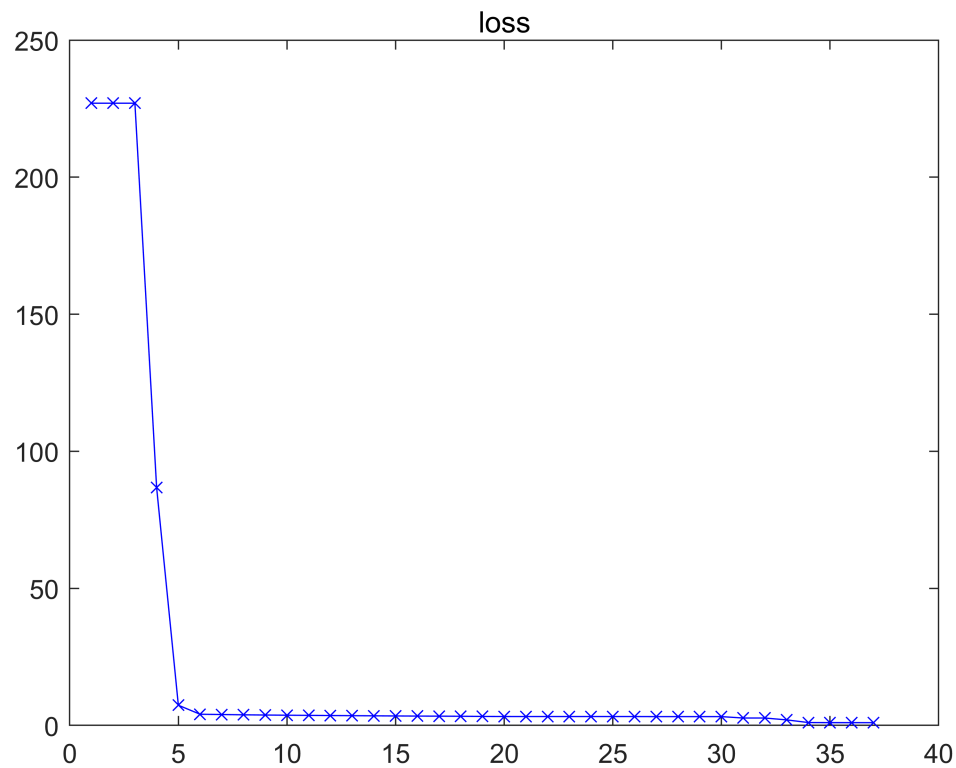


```
%error = 0.014748
```

Try different initializations for the Levenberg-Marquardt algorithm and see the effect on final model training loss error. **Now $\text{lemda} = 10^{-5}$ not changing. Let w equals to 0.1. My stopping criteria for single run algorithm is loss smaller than 1 or 50 loops if loss not smaller than 1.**

```
%When lemda =  $10^{-5}$ ,  $w = 0.1$ 
w4 = LevgMaqdt_alg(x,y);%single run algorithm
```

```
final loss=0.998053
```

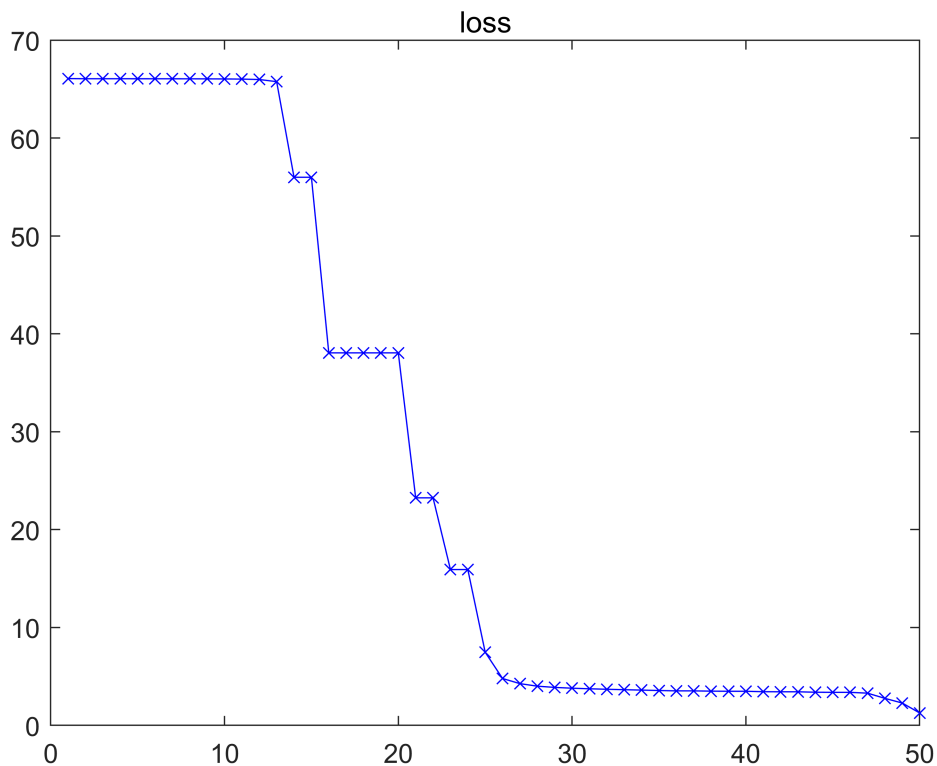


error=0.001996

When $\text{lemda} = 10^{-5}$ and initial w much larger than 1. We the loss needs more iterations to drop since initial w is far from optimized w . Also the loss is larger when than $w = 0.1$

```
%When lemda = 10^-5,w = 3
w5 = LevgbMaqdt_alg(x,y);%single run algorithm
```

final loss=1.271375



error=0.002543

ii. Test the model on another set of NT randomly generated points $\tilde{x}(n)$, $n = 1, 2, \dots, NT$ (which are different from those used to train the network). Show how you compute the error for test points (try different values of NT between 50 and 300) and report the values. Summarize your training and testing error metrics for different initializations, different choices of λ and NT. Comment on your results.

Generate 50 data points as test set.

```
%generate 50 data points as test set
[x,y] = rand_gen(50);
```

Use multiple run LM algorithm to calculate the error.

```
%When lemda = 10^-5
w6 = Levgb_Maqrt_alg(x,y);
```

```
final loss=0.000842
error=0.000017
```

```
%error = 0.000017
```

```
%When lemda = 10^-3
w7 = Levgb_Maqrt_alg(x,y);
```

```
final loss=0.013910
error=0.000278
```

```
%error = 0.000278
```

```
%When lemda = 1  
w8 = Levgb_Maqrt_alg(x,y);
```

```
final loss=2.848993  
error=0.056980
```

```
%error = 0.05698
```

Generate 100 data points as test set.

```
%generate 100 data points as test set  
[x,y] = rand_gen(100);
```

Use multiple run LM algorithm to calculate the error.

```
%When lemda = 10^-5  
w15 = Levgb_Maqrt_alg(x,y);
```

```
final loss=0.001735  
error=0.000017
```

```
%error = 0.000017
```

```
%When lemda = 10^-3  
w16 = Levgb_Maqrt_alg(x,y);
```

```
final loss=0.015339  
error=0.000153
```

```
%error = 0.000153
```

```
%When lemda = 1  
w17 = Levgb_Maqrt_alg(x,y);
```

```
final loss=3.519118  
error=0.035191
```

```
%error = 0.035191
```

Generate 200 data points as test set.

```
[x,y] = rand_gen(200);
```

```
%When lemda = 10^-5  
w9 = Levgb_Maqrt_alg(x,y);
```

```
final loss=0.001886  
error=0.000009
```

```
%error = 0.000009
```

```
%When lemda = 10^-3  
w10 = Levbg_Maqrt_alg(x,y);
```

```
final loss=0.018087  
error=0.000090
```

```
%error = 0.00009
```

```
%When lemda = 1  
w11 = Levbg_Maqrt_alg(x,y);
```

```
final loss=4.214546  
error=0.021073
```

```
%error = 0.021073
```

Generate 300 data points as test set.

```
[x,y] = rand_gen(300);
```

```
%When lemda = 10^-5  
w12 = Levbg_Maqrt_alg(x,y);
```

```
final loss=0.003997  
error=0.000013
```

```
%error = 0.0000013
```

```
%When lemda = 10^-3  
w13 = Levbg_Maqrt_alg(x,y);
```

```
final loss=0.020934  
error=0.000070
```

```
%error = 0.00007
```

```
%When lemda = 1  
w14 = Levbg_Maqrt_alg(x,y);
```

```
final loss=5.205768  
error=0.017353
```

```
%error = 0.017353
```

```
%first row is NT  
%first column is lemnda  
%others are the test error matrix of corresponding NT and lemnda  
err_matx = [0 50 100 200 300;
```



```
10^-5 0.0017 0.0017 0.0009 0.00013;  
10^-3 0.0278 0.0153 0.009 0.007;  
1 5.698 3.5191 2.1073 1.7353]
```

```
err_matx = 4x5  
    0    50.0000   100.0000   200.0000   300.0000  
0.0000    0.0017    0.0017    0.0009    0.0001  
0.0010    0.0278    0.0153    0.0090    0.0070  
1.0000    5.6980    3.5191    2.1073    1.7353
```

```
%To see the trend easier, I multiple all the error rate with 100
```

The error matrix shows that as test NT approaches training NT(500), the error rates decrease. Also, as lambda grows, the error rates also increase.