

1.2 Find the Minimum of the Rosenbrock Function [Code and Report]

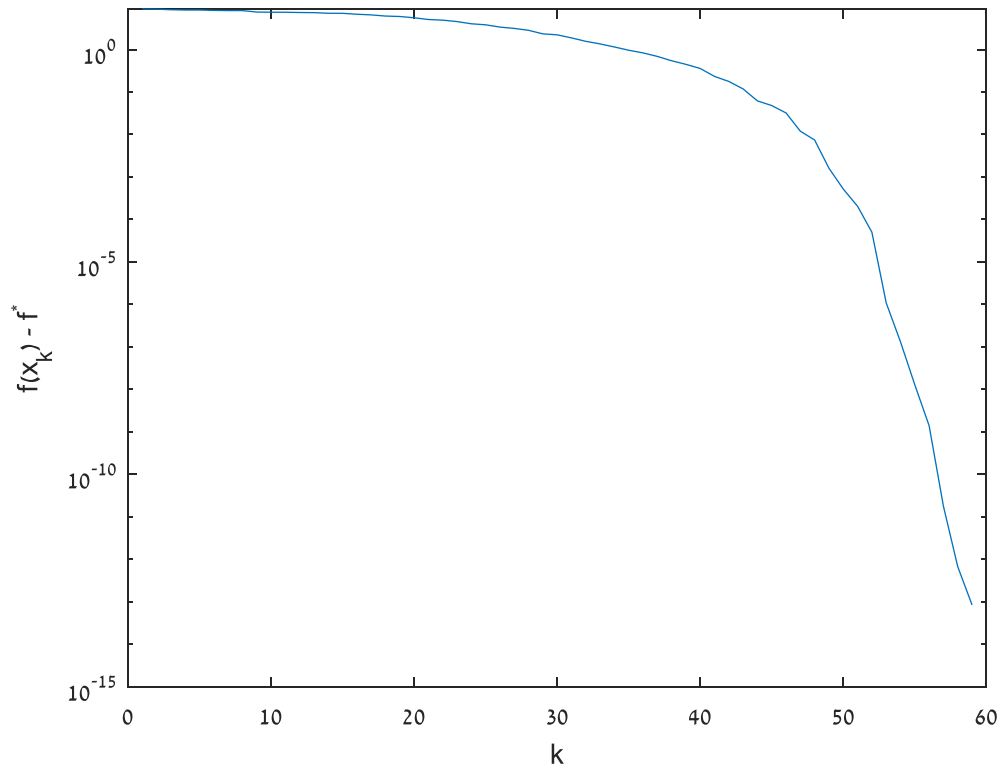
Use the BFGS method to find the minimum point of the [Rosenbrock function](#) for the case of $n = 10$ and starting point $x_0 = (0, 0, \dots, 0) \in \mathbb{R}^{10}$. As a reminder, the Rosenbrock function is given as follows:

$$f(x_1, x_2, \dots, x_n) = \sum_{i=1}^{n-1} \left((1 - x_i)^2 + 100 (x_{i+1} - x_i^2)^2 \right) \quad (1)$$

plot the convergence curve $f(x_k) - f^*$ (the y-axis) as function of the iteration number k (the x-axis), where f^* is the optimal/minimal value of the Rosenbrock function.

Use logarithmic scale for the y-axis.

From the following graph, we can see that the convergence happens swiftly using BFGS, where the number of iterations is relatively small and thus the running time for rosenbrock function is significantly less

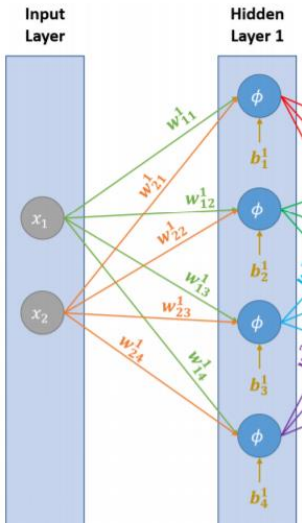


1.3.4 Explicit Expression of the Neural Network Model [Report]

Complete the following:

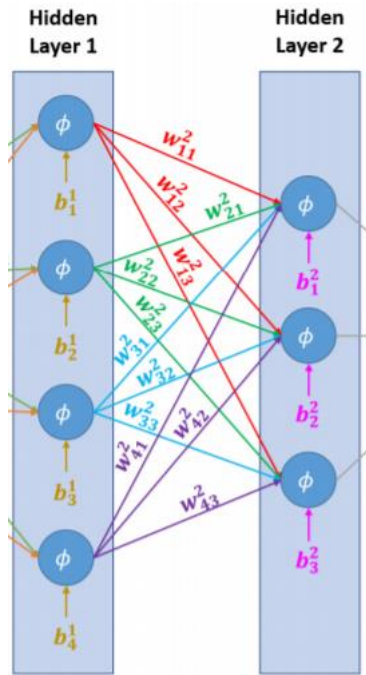
- Write an explicit vectorized expression for the neural network given in section 1.3.1. In other words, you should write the explicit expression of $F(x|\mathcal{W})$. Please note, \mathcal{W} represent the set of model parameters. That is, in our feed-forward neural network context, \mathcal{W} consists of all weight matrices $\{W_i\}$ and all bias vectors $\{b_i\}$ across all layers. Even though your explicit expression for $F(x|\mathcal{W})$ should be written in terms of x (the input vector), $\{W_i\}$, $\{b_i\}$ and the activation function ϕ , you can also think of \mathcal{W} , for implementation purposes, as a single long vector which is made of stacking the columns/rows of the matrices $\{W_i\}$ and the vectors $\{b_i\}$ one on top of the other. For guidance, please refer to figures 5 and 6.

the output of the following first hidden layer would be:



$$\phi \left(\begin{bmatrix} w_{11}^1 & w_{21}^1 \\ w_{12}^1 & w_{22}^1 \\ w_{13}^1 & w_{23}^1 \\ w_{14}^1 & w_{24}^1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} b_1^1 \\ b_2^1 \\ b_3^1 \\ b_4^1 \end{bmatrix} \right) = \begin{bmatrix} \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_1^1) \\ \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_2^1) \\ \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_3^1) \\ \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_4^1) \end{bmatrix} = \phi(w^{1T} x_i + b_i^1)$$

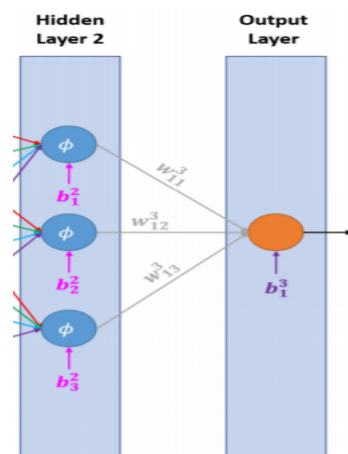
the output of the following second hidden layer would be:



$$\phi \left(\begin{bmatrix} w_{11}^2 w_{21}^2 w_{31}^2 w_{41}^2 \\ w_{12}^2 w_{22}^2 w_{32}^2 w_{42}^2 \\ w_{13}^2 w_{23}^2 w_{33}^2 w_{43}^2 \end{bmatrix} \begin{bmatrix} \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) \\ \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{21}^1) \\ \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{31}^1) \\ \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{41}^1) \end{bmatrix} + \begin{bmatrix} b_1^2 \\ b_2^2 \\ b_3^2 \end{bmatrix} \right) = \phi(w^{2T} \phi(w^{1T} x_i + b_i^1) + b_i^2)$$

$$\begin{bmatrix} \phi(w_{11}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{21}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{21}^1) + w_{31}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{31}^1) + w_{41}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{41}^1) + b_1^2) \\ \phi(w_{12}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{22}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{21}^1) + w_{32}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{31}^1) + w_{42}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{41}^1) + b_2^2) \\ \phi(w_{13}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{23}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{21}^1) + w_{33}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{31}^1) + w_{43}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{41}^1) + b_3^2) \end{bmatrix}$$

the output of the following output layer would be:



$$\left(\begin{matrix} w_{i1}^3 & w_{i2}^3 & w_{i3}^3 \end{matrix} \begin{bmatrix} \phi(w_{11}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{21}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{31}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) + w_{41}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{11}^2) \\ \phi(w_{12}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{22}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{32}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) + w_{42}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{22}^2) \\ \phi(w_{13}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{23}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{33}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) + w_{43}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{32}^2) \end{bmatrix} \right) + b_{i1}^3 = \boxed{w^3{}^T \phi \left(w^2{}^T \phi \left(w^1{}^T x_i + b_i^1 \right) + b_i^2 \right) + b_i^3} =$$

$$\begin{aligned} & [w_{11}^3 \cdot \phi(w_{11}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{21}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{31}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) + w_{41}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{11}^2) + \\ & w_{12}^3 \cdot \phi(w_{12}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{22}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{32}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) \\ & + w_{42}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{22}^2) + \\ & w_{13}^3 \cdot \phi(w_{13}^2 \cdot \phi(w_{11}^1 x_1 + w_{21}^1 x_2 + b_{11}^1) + w_{23}^2 \cdot \phi(w_{12}^1 x_1 + w_{22}^1 x_2 + b_{12}^1) + w_{33}^2 \phi(w_{13}^1 x_1 + w_{23}^1 x_2 + b_{13}^1) + w_{43}^2 \phi(w_{14}^1 x_1 + w_{24}^1 x_2 + b_{14}^1) + b_{32}^2) + b_{11}^3] \end{aligned}$$

To summarize:

- $W^1: 2 \times 4$; $b^1: 4 \times 1$
- $W^2: 4 \times 3$; $b^2: 3 \times 1$
- $W^3: 3 \times 1$; $b^3: 1 \times 1$
- Our output function expression would be:

$$F(x|w) = w^3{}^T \phi_2 \left(w^2{}^T \phi_1 \left(w^1{}^T x_i + b_i^1 \right) + b_i^2 \right) + b_i^3$$

When:

$$\phi(x) = \tanh(x); \phi'(x) = \frac{4e^{-2x}}{(1 + e^{-2x})^2}$$

$$\phi_1 \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = \begin{pmatrix} \frac{e^{u_1} - e^{-u_1}}{e^{u_1} + e^{-u_1}} \\ \frac{e^{u_2} - e^{-u_2}}{e^{u_2} + e^{-u_2}} \\ \frac{e^{u_3} - e^{-u_3}}{e^{u_3} + e^{-u_3}} \end{pmatrix} \quad \phi_2 \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} \frac{e^{u_1} - e^{-u_1}}{e^{u_1} + e^{-u_1}} \\ \frac{e^{u_2} - e^{-u_2}}{e^{u_2} + e^{-u_2}} \\ \frac{e^{u_3} - e^{-u_3}}{e^{u_3} + e^{-u_3}} \end{pmatrix}$$

$$\phi_1' \begin{pmatrix} u_1 \\ u_2 \\ u_3 \\ u_4 \end{pmatrix} = \begin{pmatrix} \phi'(u_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \phi'(u_4) \end{pmatrix}; \quad \phi_2' \begin{pmatrix} u_1 \\ u_2 \\ u_3 \end{pmatrix} = \begin{pmatrix} \phi'(u_1) & \cdots & 0 \\ \vdots & \ddots & \vdots \\ 0 & \cdots & \phi'(u_3) \end{pmatrix}$$

1.3.7 Evaluating the Loss Function's Derivative with Respect to Network's Output of Single Training Example [Code] (Updated: 29/05/2021)

In this section we will treat the loss function as a standalone module. Write a routine that calculates the derivative of the loss function with respect to the neural network's

output. That is, you should derive and implement $\frac{\partial L}{\partial F(x_i, \mathcal{W})}$.

$$\begin{aligned} L &= \frac{1}{n} \sum_{i=1}^n \left((F(x^i, W) - y_i)^2 \right); \nabla_w L = \frac{1}{n} \sum_{i=1}^n \nabla_w \left((F(x^i, W) - y_i)^2 \right) \\ \frac{\partial L}{\partial F(x_i, w)} &= \frac{\partial \frac{1}{n} \sum_{i=1}^n \left((F(x^i, W) - y_i)^2 \right)}{\partial F(x_i, w)} = \frac{\frac{1}{n} \sum_{i=1}^n \partial \left((F(x^i, W) - y_i)^2 \right)}{\partial F(x_i, w)} \\ &= \frac{\frac{1}{n} \sum_{i=1}^n \partial(\psi(r))}{\partial r_i} = \frac{1}{n} \sum_{i=1}^n \frac{\partial(\psi(r_i))}{\partial r_i} \end{aligned}$$

We notice we want output of a single training example:

$$F(x^i|w) = w^3{}^T \phi_2 \left(w^2{}^T \phi_1 \left(w^1{}^T x_i + b_i^1 \right) + b_i^2 \right) + b^3$$

$$r_i = F(x^i, W) - y_i$$

$$\psi(r) = r^2$$

$$\psi'(r) = 2rdr$$

$$\frac{\partial L}{\partial F(x_i, w)} = \frac{(F(x^i, W) - y_i)^2}{\partial F(x_i, w)} = \frac{\partial(\psi(r))}{\partial r_i}$$

$$\frac{\partial L}{\partial F(x_i, w)} = 2rdr = 2(F(x^i, W) - y_i)$$

1.3.14 Combining it All Together [Code and Report] (Updated: 27/05/2021)

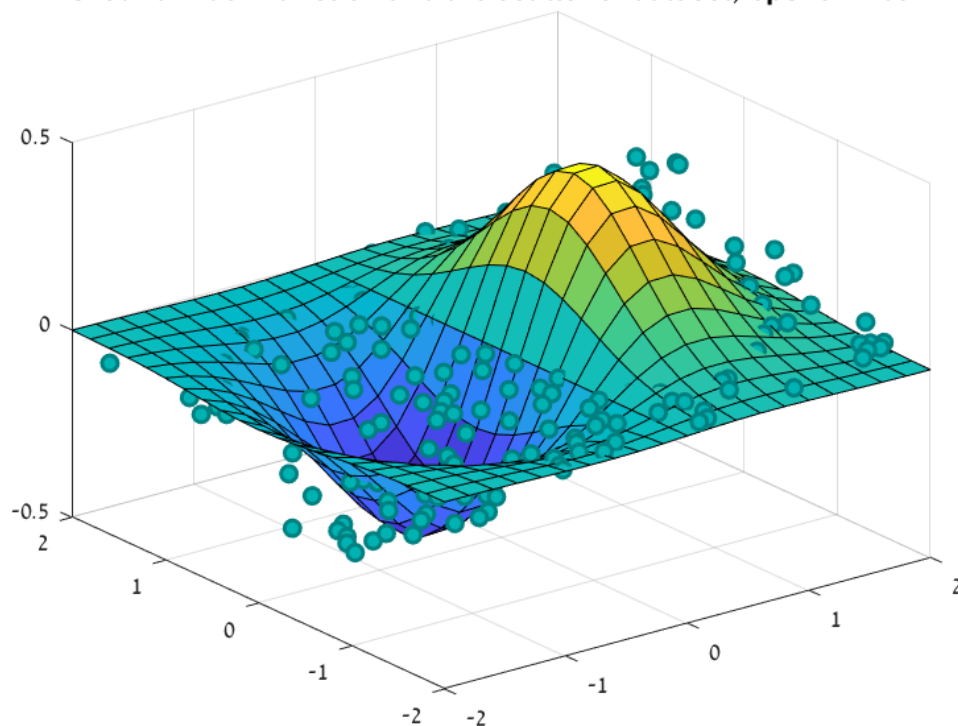
In order to train the neural network and plot the results, follow the steps below:

In the plots we see the ground truth function and upon it, a scatter of the reconstruction of the test set that was estimated using our trained network, both of which are limited to the subspace $[-2,2] \times [-2,2]$. Each plot shows the scattering and the ground function for a different epsilon which is the stopping condition on the gradient in the BFGS algorithm.

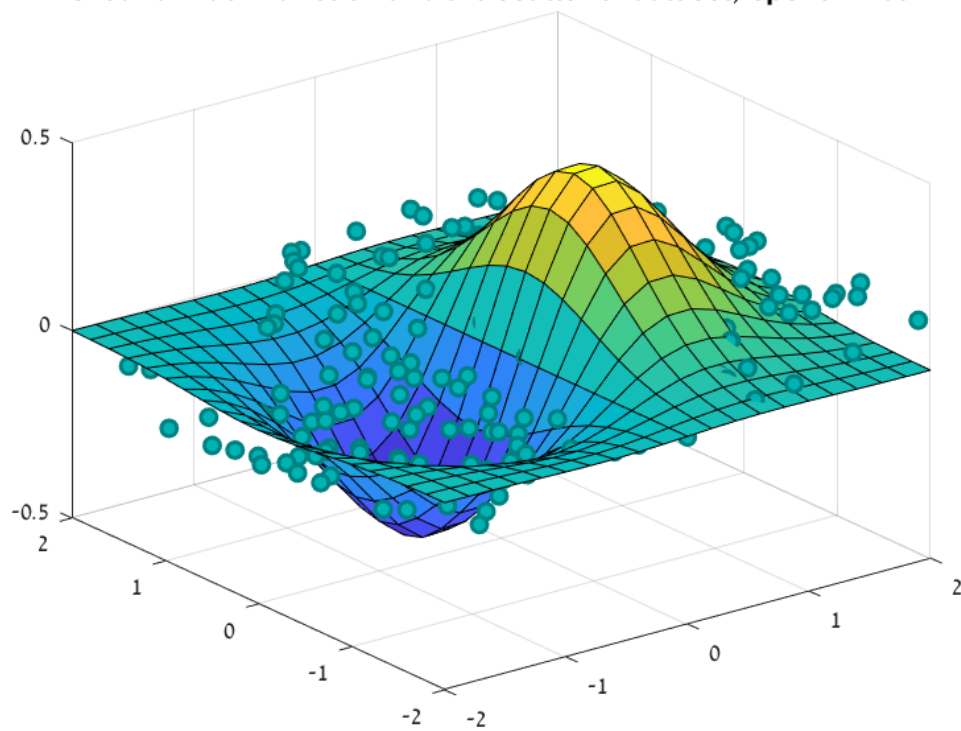
Conclusions:

We can see from the graphs that the smaller the epsilon the more accurate the scatter matches the ground truth function, meaning that the reconstruction of the test set using the trained network works best when the stopping condition is as close to zero as possible.

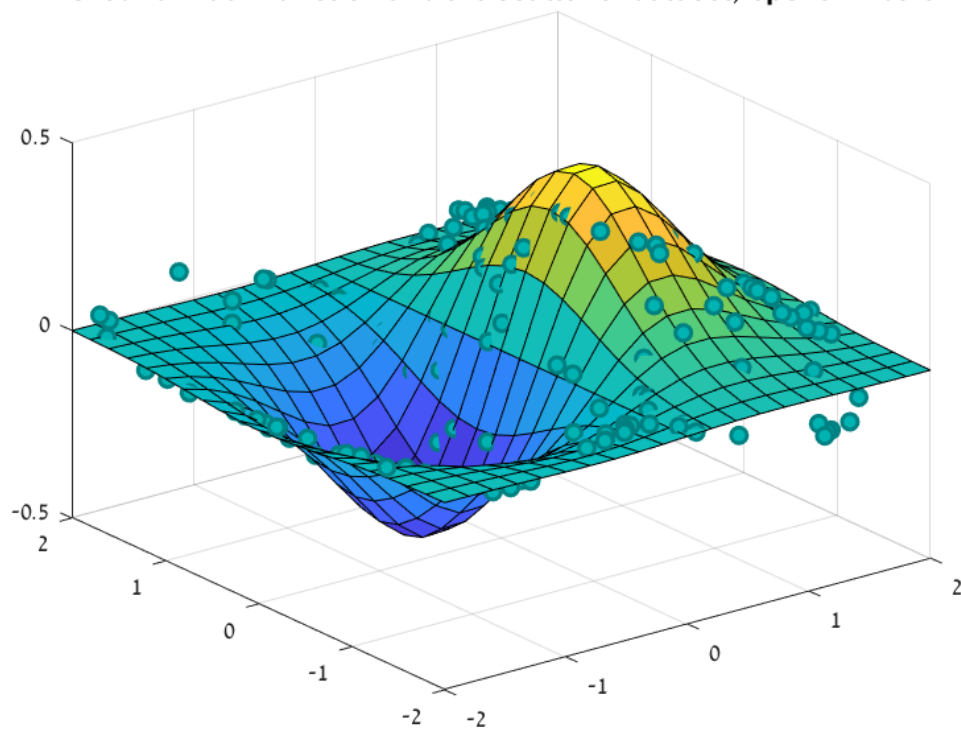
Ground Truth Function and the scatter of dataset; epsilon=10e-1



Ground Truth Function and the scatter of dataset; epsilon=10e-2



Ground Truth Function and the scatter of dataset; epsilon=10e-3



Ground Truth Function and the scatter of dataset; epsilon=10e-4

