

Problem Set 1:

1. (a) (i): $\theta^{(t)} = \theta^{(t-1)} - \alpha \nabla J(\theta^{(t-1)})$ $\nabla J(\theta^{(t-1)}) = \frac{\partial (\frac{1}{2} \beta \theta^2)}{\partial \theta} = \beta \theta$

$$= \theta^{(t-1)} - \alpha \beta \theta^{(t-1)}$$

$$= (1 - \alpha \beta) \cdot \theta^{(t-1)}$$

$$= (1 - \alpha \beta)^t \theta^{(0)}$$

assume $\lim_{t \rightarrow \infty} |\theta^{(t)} - \theta^+| = 0$

$$\lim_{t \rightarrow \infty} |(1 - \alpha \beta)^t \theta^{(0)} - \theta^+| = 0$$

$$\lim_{t \rightarrow \infty} (1 - \alpha \beta)^t \theta^{(0)} = \theta^+$$

$$\because \theta^{(0)} \neq 0 \quad \therefore \lim_{t \rightarrow \infty} (1 - \alpha \beta)^t = \frac{\theta^+}{\theta^{(0)}}$$

~~$\theta^{(0)}$~~ $\neq 0$ \neq let $\underline{\theta^+ = 0}$

$$\lim_{t \rightarrow \infty} (1 - \alpha \beta)^t = 0$$

$$-1 < (1 - \alpha \beta) < 1$$

$$-2 < -\alpha \beta < 0$$

$$2 > \alpha \beta > 0$$

$$\because \beta > 0 \quad \therefore \alpha < 0$$

$$\therefore \underline{\frac{2}{\beta} > \alpha > 0}$$

close form:

$$\nabla J(\theta) = \beta \theta = 0$$

$$\therefore \beta > 0$$

$$\therefore \underline{\theta^* = 0}$$

$$\therefore \text{global minimum } \underline{\theta^* = 0} = \arg \min_{\theta} J(\theta)$$

1.(a) (i)

$$|\theta^{(T)} - \theta^*| \leq \epsilon$$

$$\because \theta^* = 0$$

$$\therefore |\theta^{(T)}| \leq \epsilon$$

$$|(1-\alpha\beta)^T \theta^{(0)}| \leq \epsilon$$

$$|(1-\alpha\beta)^T| \cdot |\theta^{(0)}| \leq \epsilon$$

$$\therefore \theta^{(0)} \neq 0$$

$$\therefore |(1-\alpha\beta)^T| \leq \frac{\epsilon}{|\theta^{(0)}|} \Rightarrow |(1-\alpha\beta)|^T \leq \underbrace{\frac{\epsilon}{|\theta^{(0)}|}}_{\text{constant}}$$

$$\because |(1-\alpha\beta)| \leq 1$$

$$\therefore 0 \leq |1-\alpha\beta| < 1$$

when $|1-\alpha\beta|=0$, $T=1$

when $0 < |1-\alpha\beta| < 1$, then when $|1-\alpha\beta| \nearrow$, $T \nearrow$

$\therefore \begin{cases} \text{when } 1-\alpha\beta=0, \text{ i.e. } \alpha=\frac{1}{\beta}, T=1 \\ \text{when } 1-\alpha\beta<0, \text{ i.e. } \alpha>\frac{1}{\beta}, \text{ when } |1-\alpha\beta|\nearrow, 1-\alpha\beta\nearrow, \alpha\nearrow, T\nearrow \end{cases}$

$\begin{cases} \text{when } 1-\alpha\beta>0, \text{ i.e. } \alpha<\frac{1}{\beta}, \text{ when } |1-\alpha\beta|\nearrow, 1-\alpha\beta\nearrow, \alpha\nearrow, T\nearrow \end{cases}$

\therefore When choose an inappropriate α (e.g. ~~10^{-100}~~) can significantly increase T

1.(b)

$$\nabla J(\theta) = \begin{bmatrix} \beta_1 \theta_1 \\ \beta_2 \theta_2 \\ \vdots \\ \beta_d \theta_d \end{bmatrix}$$

let \vec{B} be diagonal ($\beta_1, \beta_2, \dots, \beta_d$)

$$\begin{aligned}\theta^{(t)} &= \theta^{(t-1)} - \alpha \nabla J(\theta^{(t-1)}) \\ &= \theta^{(t-1)} - \alpha \vec{B} \cdot \theta^{(t-1)} \\ &= (\vec{I} - \alpha \vec{B}) \cdot \theta^{(t-1)} \\ &= (\vec{I} - \alpha \vec{B})^T \cdot \theta^{(0)}\end{aligned}$$

$$\lim_{t \rightarrow \infty} \|\theta^{(t)} - \theta^+\|_2 = 0$$

$$\rightarrow \lim_{t \rightarrow \infty} \|(\vec{I} - \alpha \vec{B})^T \theta^{(0)} - \theta^+\|_2 = 0$$

$$\lim_{t \rightarrow \infty} (\vec{I} - \alpha \vec{B})^T \theta^{(0)} = \theta^+$$

when $\theta^+ = \vec{0}$ and $|\vec{I} - \alpha \vec{B}| < 1$, it converges.

$$\vec{I} - \alpha \vec{B}$$

$\therefore 1 < |\vec{I} - \alpha \vec{B}| < 1$ for all i in $1, 2, \dots, d$

$\therefore \beta_i > 0$

$$\therefore \frac{2}{\beta_i} > \alpha > 0$$

$$\therefore 0 < \alpha < \min\left(\frac{2}{\beta_1}, \frac{2}{\beta_2}, \dots, \frac{2}{\beta_d}\right)$$

(12)(C)(iii)

Observation:

- ① for the same learning rate, both A and rotated A result in the same iterations to converge.
- ② once $\text{lr} \geq 0.5$, it starts to overshoot during the training. It bounces between the global maximum.
- ③ the iterations number decreases then increases when lr increases

$$\text{d) } J(\theta^t) = J(\theta^{t-1}) + \nabla J(\theta^{t-1})^T (\theta^t - \theta^{t-1}) + \frac{1}{2} (\theta^t - \theta^{t-1})^T \nabla^2 J[\theta^{t-1} + c(\theta^t - \theta^{t-1})] (\theta^t - \theta^{t-1})$$

$$\theta^t = \theta^{t-1} - \alpha \nabla J(\theta^{t-1})$$

$$\therefore \theta^t - \theta^{t-1} = -\alpha \nabla J(\theta^{t-1})$$

$$\begin{aligned} J(\theta^t) &= J(\theta^{t-1}) - \alpha \nabla J(\theta^{t-1})^T \nabla J(\theta^{t-1}) + \frac{1}{2} (-2 \nabla J(\theta^{t-1}))^T \nabla^2 J[\theta^{t-1} + c(-\alpha \nabla J(\theta^{t-1}))] \\ &\quad - \alpha \nabla J(\theta^{t-1}) \end{aligned}$$

$$= J(\theta^{t-1}) - \alpha \nabla J(\theta^{t-1})^T \left[I - \frac{1}{2} \alpha \nabla^2 J[\theta^{t-1} - c \nabla J(\theta^{t-1})] \right] \nabla J(\theta^{t-1})$$

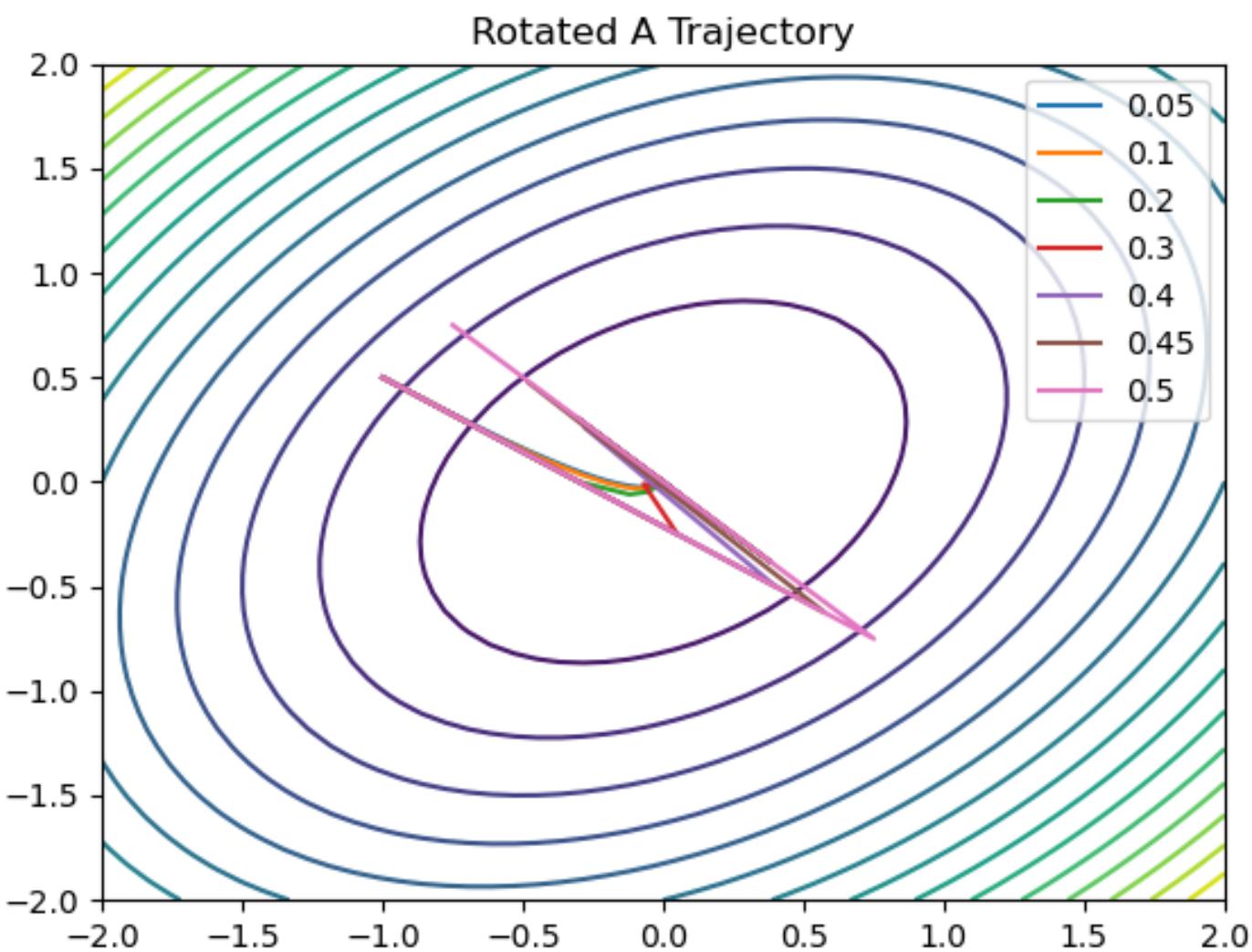
~~$\because \frac{1}{2} \alpha \nabla^2 J(\dots) \leq \frac{1}{2} \alpha \beta_{\max}$ for all points θ .~~

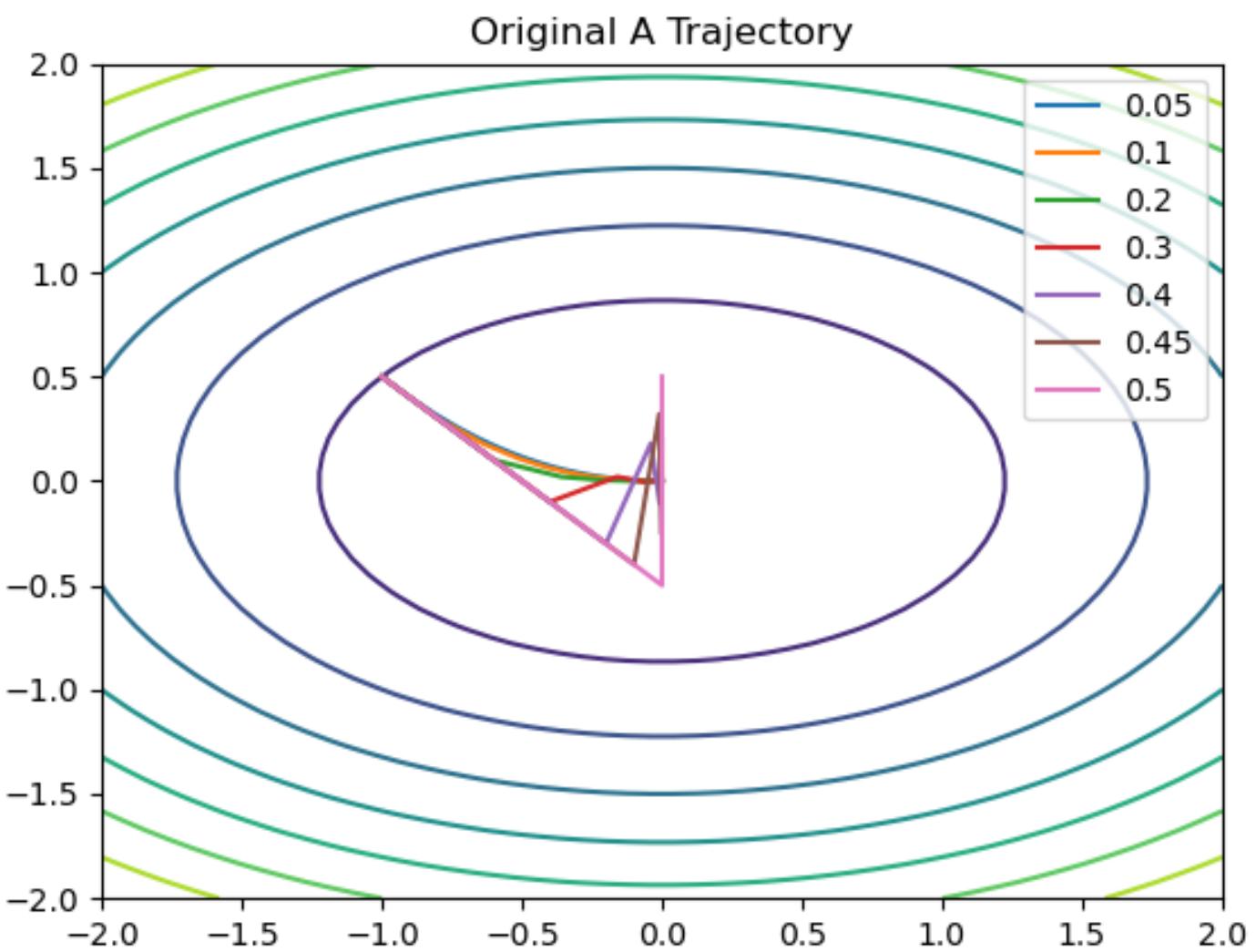
$$\therefore J(\theta^t) - J(\theta^{t-1}) -$$

~~When $I - \frac{1}{2} \alpha \nabla^2 J(\theta^{t-1} - c \nabla J(\theta^{t-1}))$ is positive definite.~~

$$J(\theta^t) < J(\theta^{t-1})$$

$$\therefore 1 - \frac{1}{2} \alpha \cdot \beta \leq$$





for $J\theta^*$ to converge

$$\alpha \nabla J\theta^{k+1}^\top [I - \frac{1}{2}\alpha \nabla^2 J(\theta^{k+1} - \alpha \nabla J\theta^{k+1})] \nabla J\theta^{k+1} > 0$$

$$\because H = \nabla^2 J(\theta) \leq \beta_{\max} \text{ for all } \theta$$

$$\therefore I - \frac{1}{2}\alpha \beta_{\max} > 0$$

$$\therefore \frac{2}{\beta_{\max}} > \alpha > 0$$

$$\therefore \text{if } 0 < \alpha < \frac{1}{\beta_{\max}}$$

$$J\theta^k < J\theta^{k+1}$$

(e)

$$2.\text{ i)} J(\theta) = (\mathbf{x}\theta - \mathbf{y})^T \cdot \mathbf{W} (\mathbf{x}\theta - \mathbf{y})$$

$$= \sum_{i=1}^n \sum_{j=1}^n w_{ij} (\mathbf{x}\theta - \mathbf{y})_i (\mathbf{x}\theta - \mathbf{y})_j$$

when $\mathbf{W} = \frac{1}{2} \text{diag}(w_1, w_2, w_3, \dots, w_n)$

$$J(\theta) = \frac{1}{2} \sum_{i=1}^n w_{ii} (\mathbf{x}\theta - \mathbf{y})_i^2$$

$$= \sum_{i=1}^n w_{ii} = \frac{1}{2} \sum_{i=1}^n w^{(i)} (\theta^T \mathbf{x}^{(i)} - y^{(i)})^2$$

$(\mathbf{x}\theta)_i = \theta^T \mathbf{x}^{(i)}$ is the i th data point

y_i is the i th observation.

$$2.\text{ ii)}, \quad \nabla J(\theta) = \nabla_{\theta} (\mathbf{x}\theta - \mathbf{y})^T \mathbf{W} (\mathbf{x}\theta - \mathbf{y})$$

$$= \nabla_{\theta} [(\mathbf{x}\theta)^T \mathbf{W} \mathbf{x}\theta - (\mathbf{x}\theta)^T \mathbf{W} \mathbf{y} - \mathbf{y}^T \mathbf{W} \mathbf{x}\theta + \mathbf{y}^T \mathbf{W} \mathbf{y}]$$

$$= \nabla_{\theta} [\theta^T \mathbf{x}^T \mathbf{W} \mathbf{x}\theta - \theta^T \mathbf{x}^T \mathbf{W} \mathbf{y} - \mathbf{y}^T \mathbf{W} \mathbf{x}\theta]$$

$$= \nabla_{\theta} [\theta^T \mathbf{x}^T \mathbf{W} \mathbf{x}\theta - 2 \theta^T \mathbf{x}^T \mathbf{W} \mathbf{y}]$$

$$= 2 \mathbf{x}^T \mathbf{W} \mathbf{x}\theta - 2 \mathbf{x}^T \mathbf{W} \mathbf{y}$$

$$= 0$$

$$\mathbf{x}^T \mathbf{W} \mathbf{x}\theta = \mathbf{x}^T \mathbf{W} \mathbf{y}$$

$$\theta = \underline{(\mathbf{x}^T \mathbf{W} \mathbf{x})^{-1} \mathbf{y}^T \mathbf{W} \mathbf{x}} \quad (\mathbf{x}^T \mathbf{W} \mathbf{x})^{-1} \mathbf{x}^T \mathbf{W} \mathbf{y}$$

2.(a) iii:

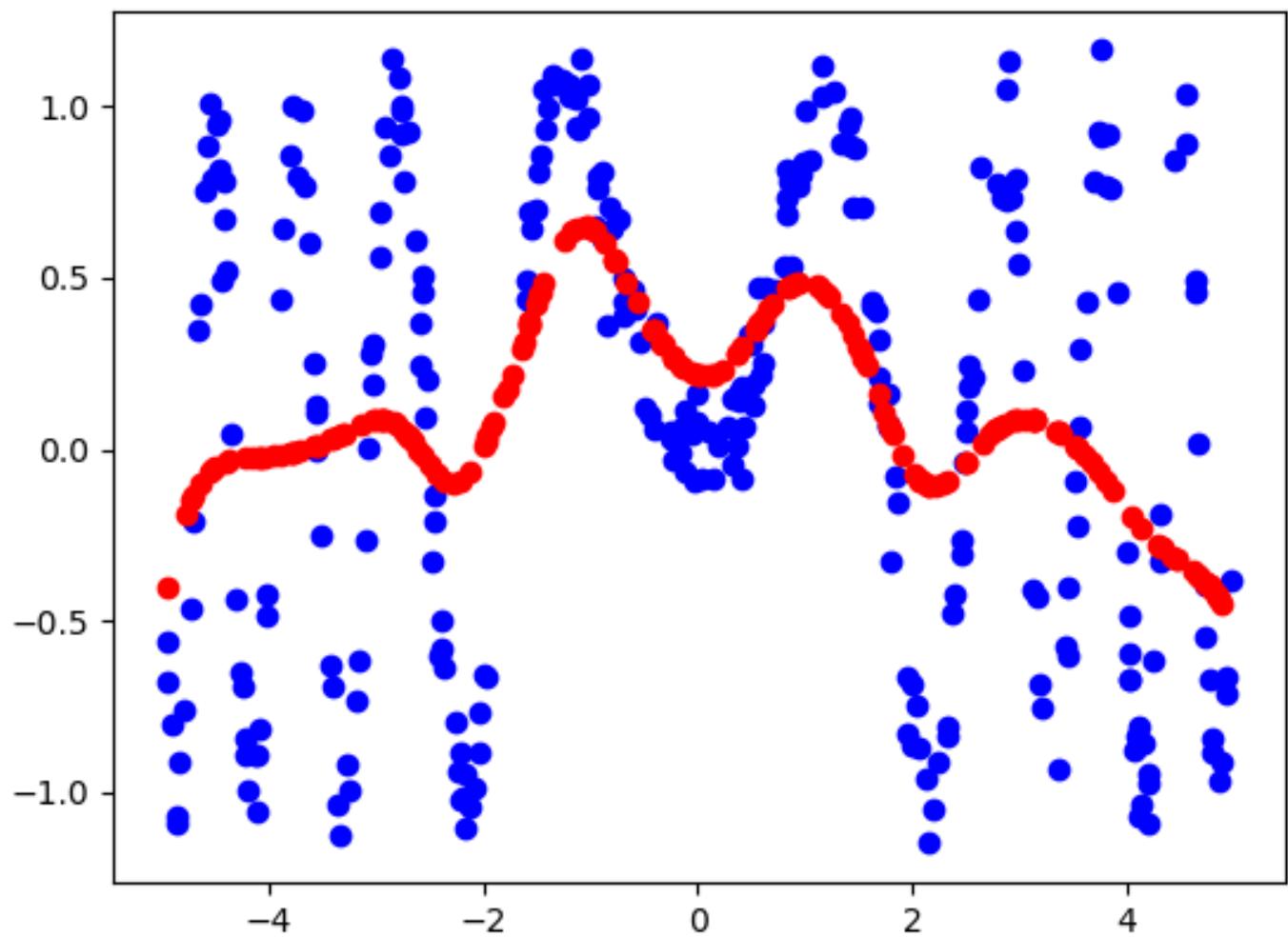
$$\underset{w}{\operatorname{argmax}} \cdot p(y_i | x_i, w) = \underset{w}{\operatorname{argmax}} \prod_{i=1}^N p(y_i | x_i, w)$$

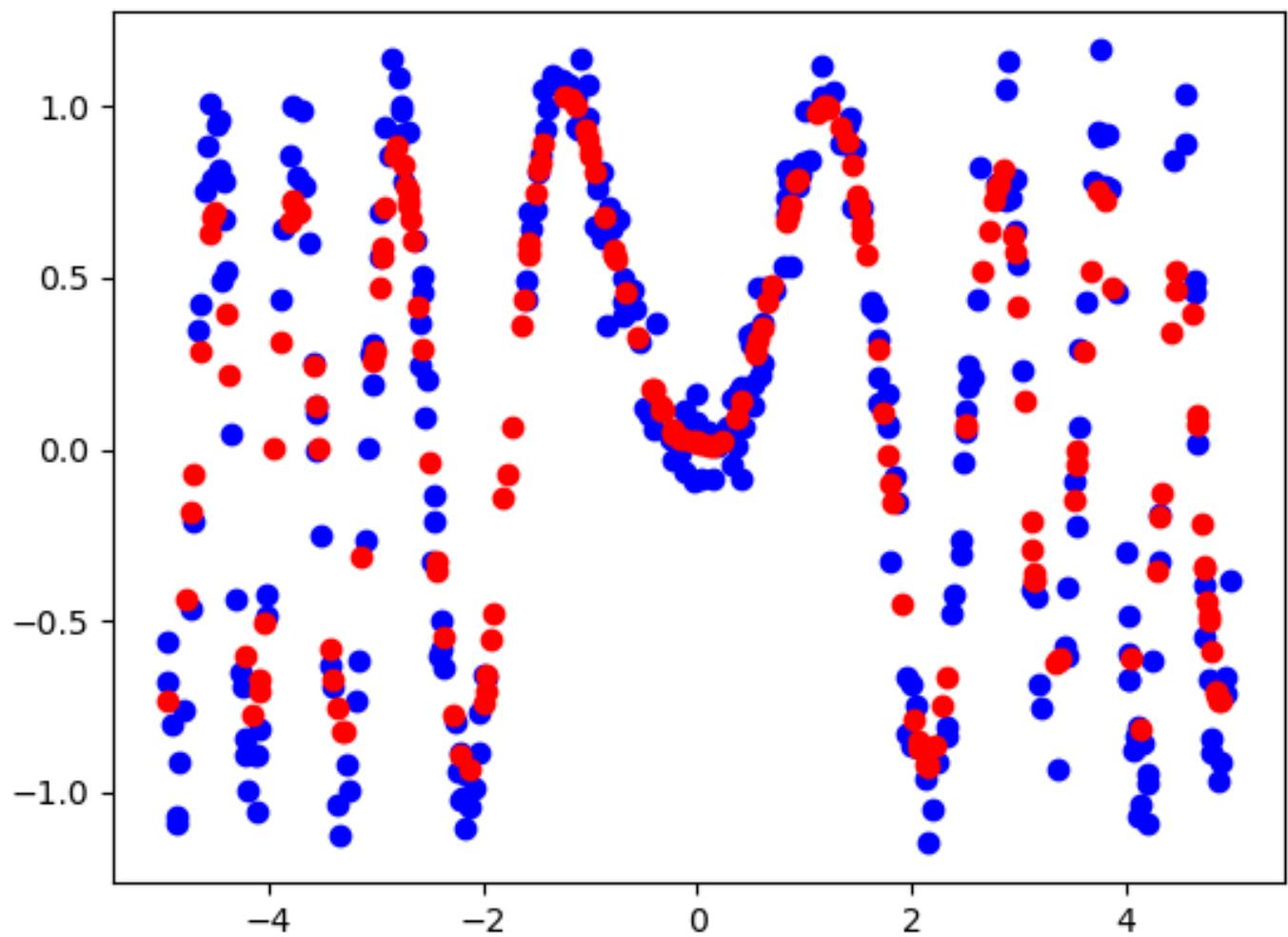
take log:

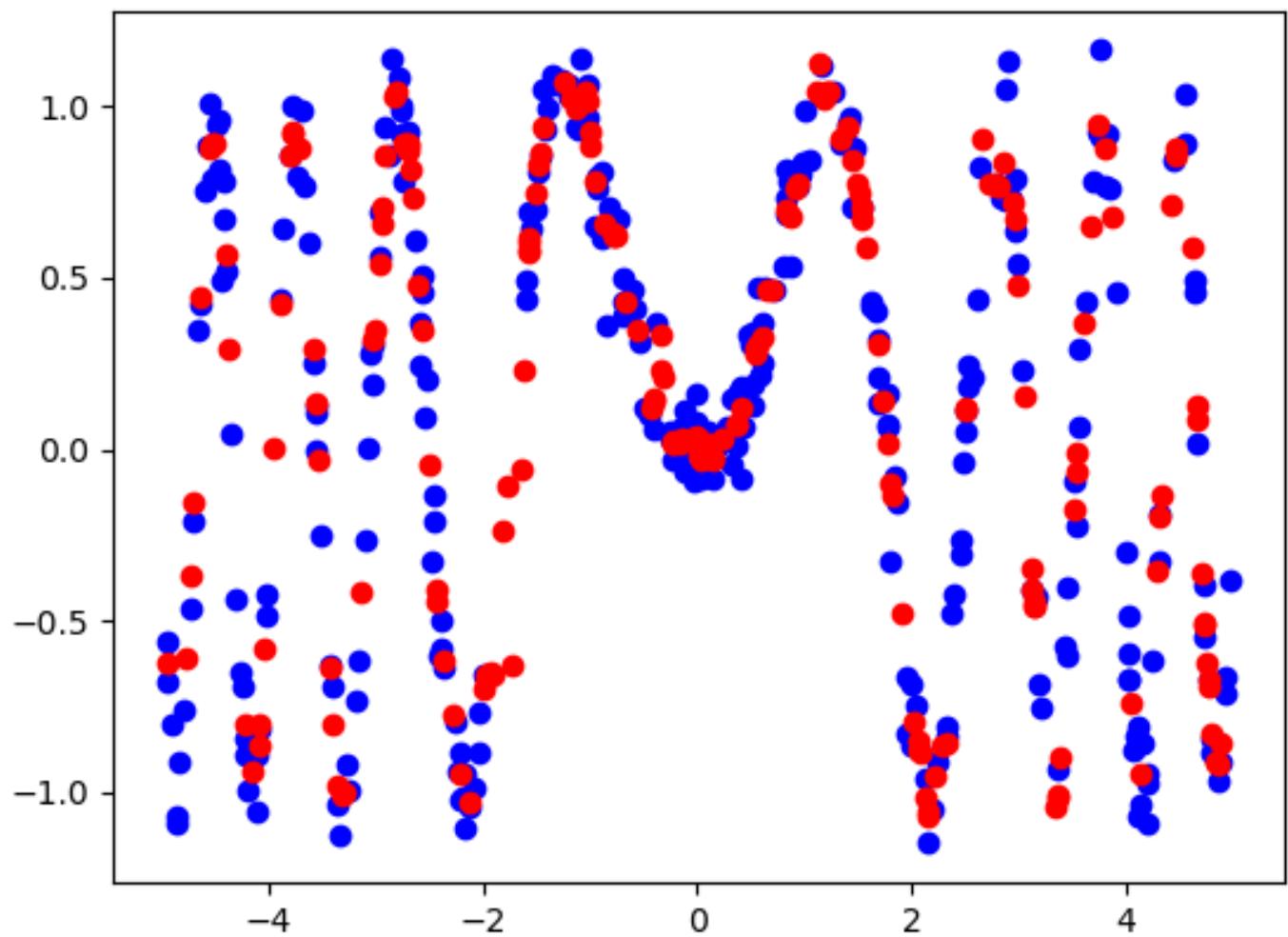
$$\underset{w}{\operatorname{argmax}} \ln \left(\prod_{i=1}^N p(y_i | x_i, w) \right)$$
$$= \underset{w}{\operatorname{argmax}} \ln \left\{ \frac{1}{\sqrt{2\pi}} \right\}^N \cdot \prod_{i=1}^N \frac{1}{\sigma^{x_i}} e^{-\frac{(y_i - \theta^T x_i)^2}{2\sigma^{x_i 2}}} \}$$

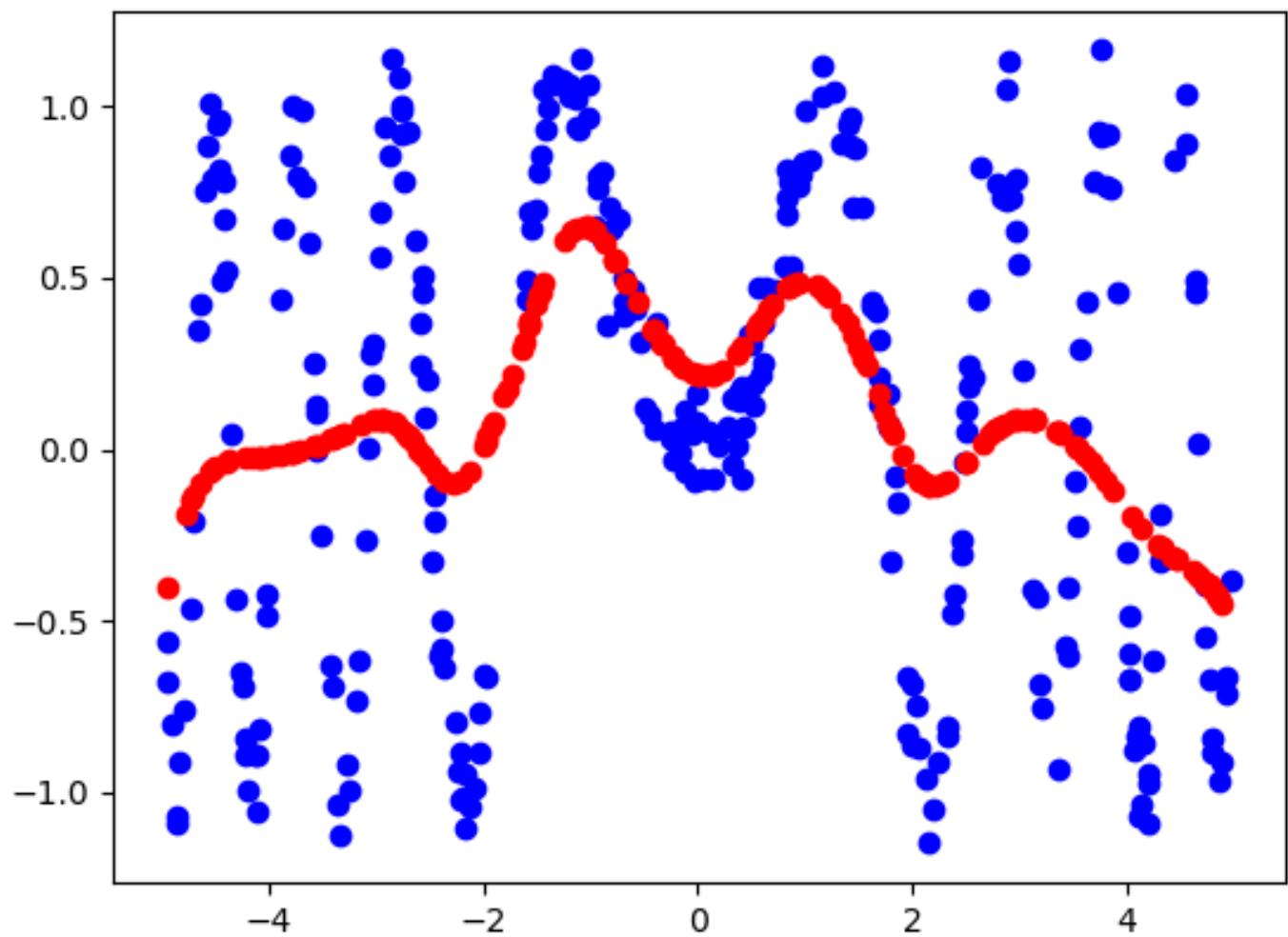
$$= N \ln \frac{1}{\sqrt{2\pi}} + \sum_{i=1}^N \left\{ \ln \frac{1}{\sigma^{x_i}} - \frac{(y_i - \theta^T x_i)^2}{2\sigma^{x_i 2}} \right\}$$

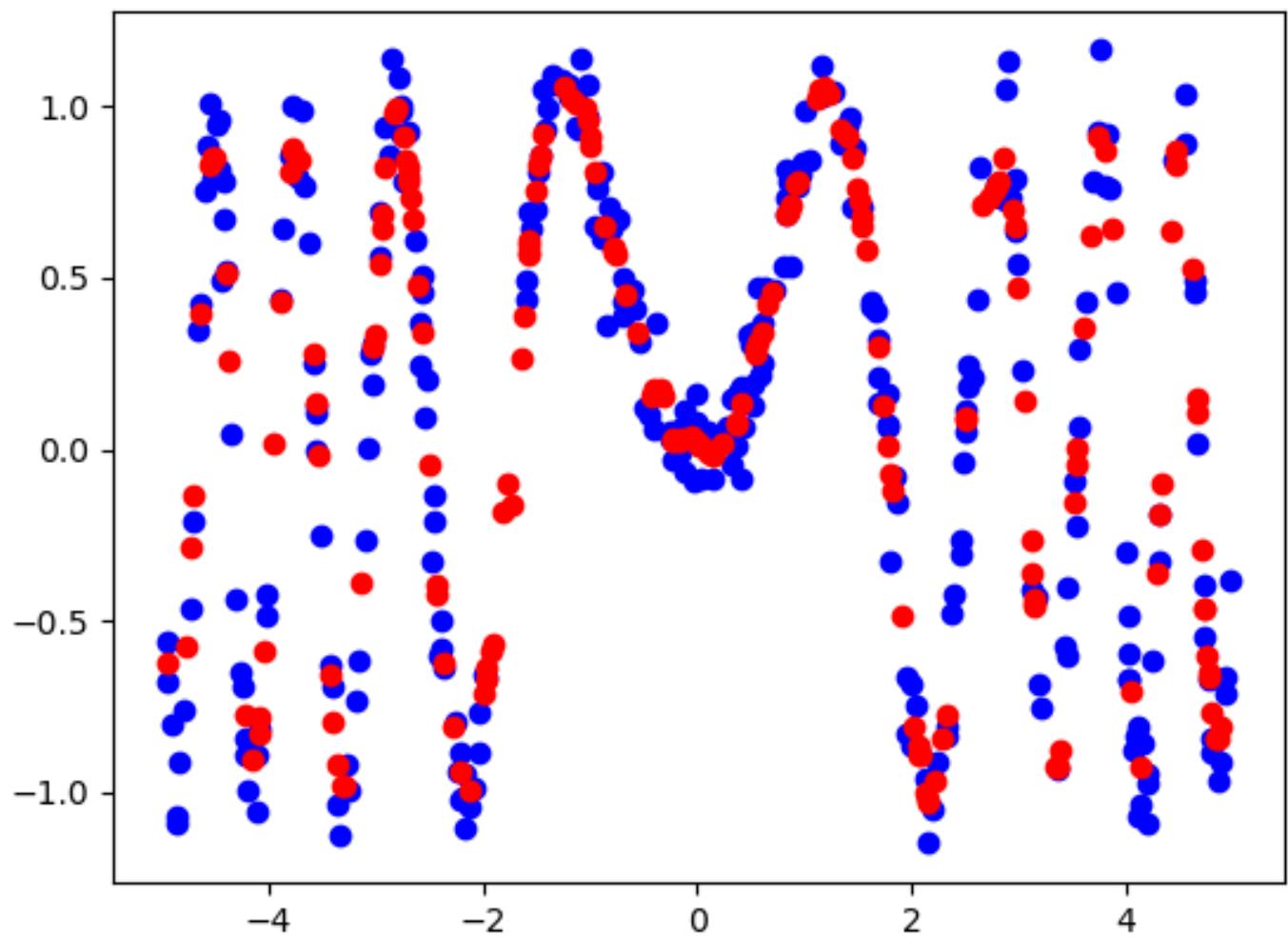
$$= N \ln \frac{1}{\sqrt{2\pi}} - \sum_{i=1}^N \left(\ln \sigma^{x_i} + \frac{(y_i - \theta^T x_i)^2}{2\sigma^{x_i 2}} \right)$$

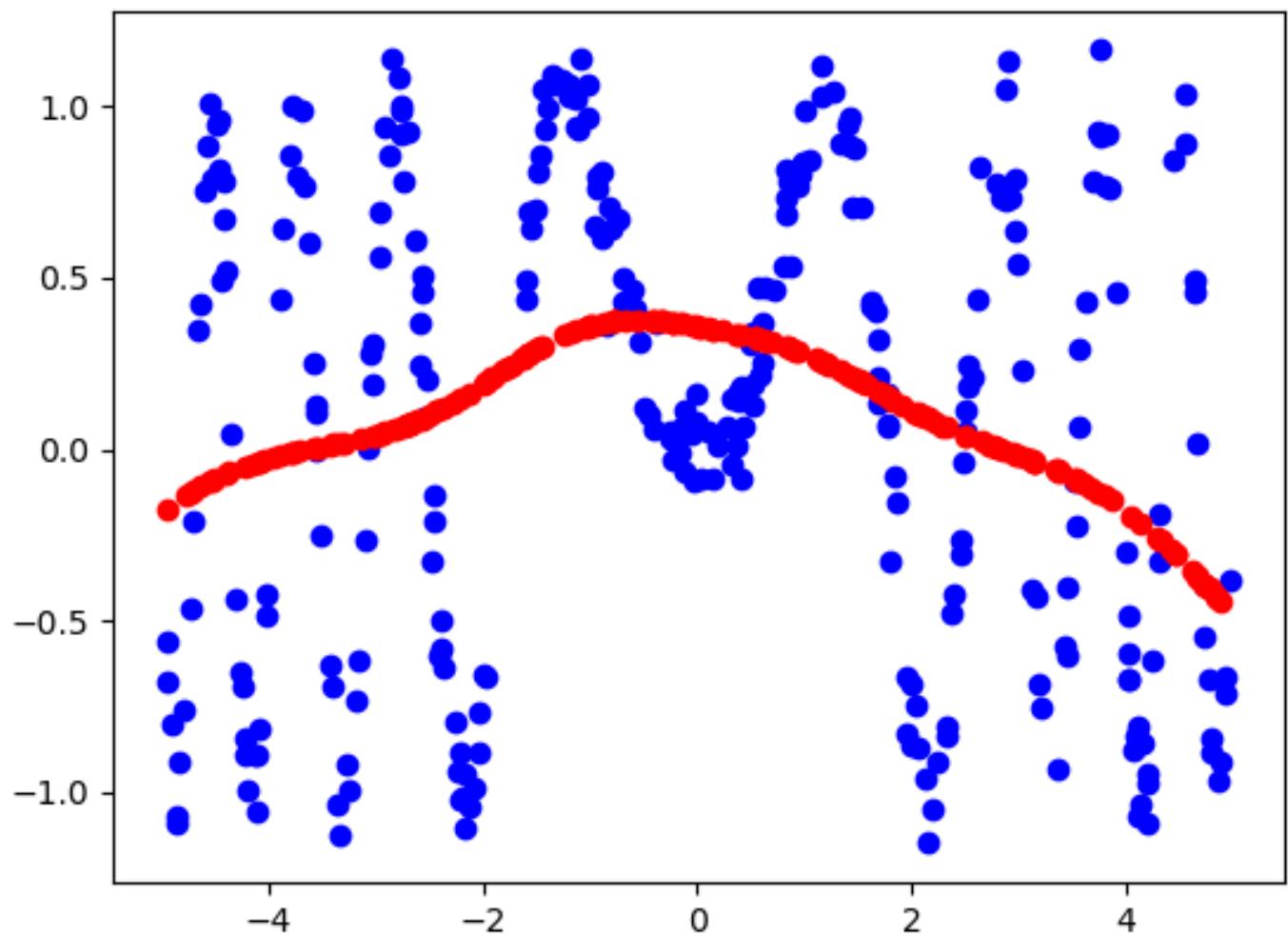


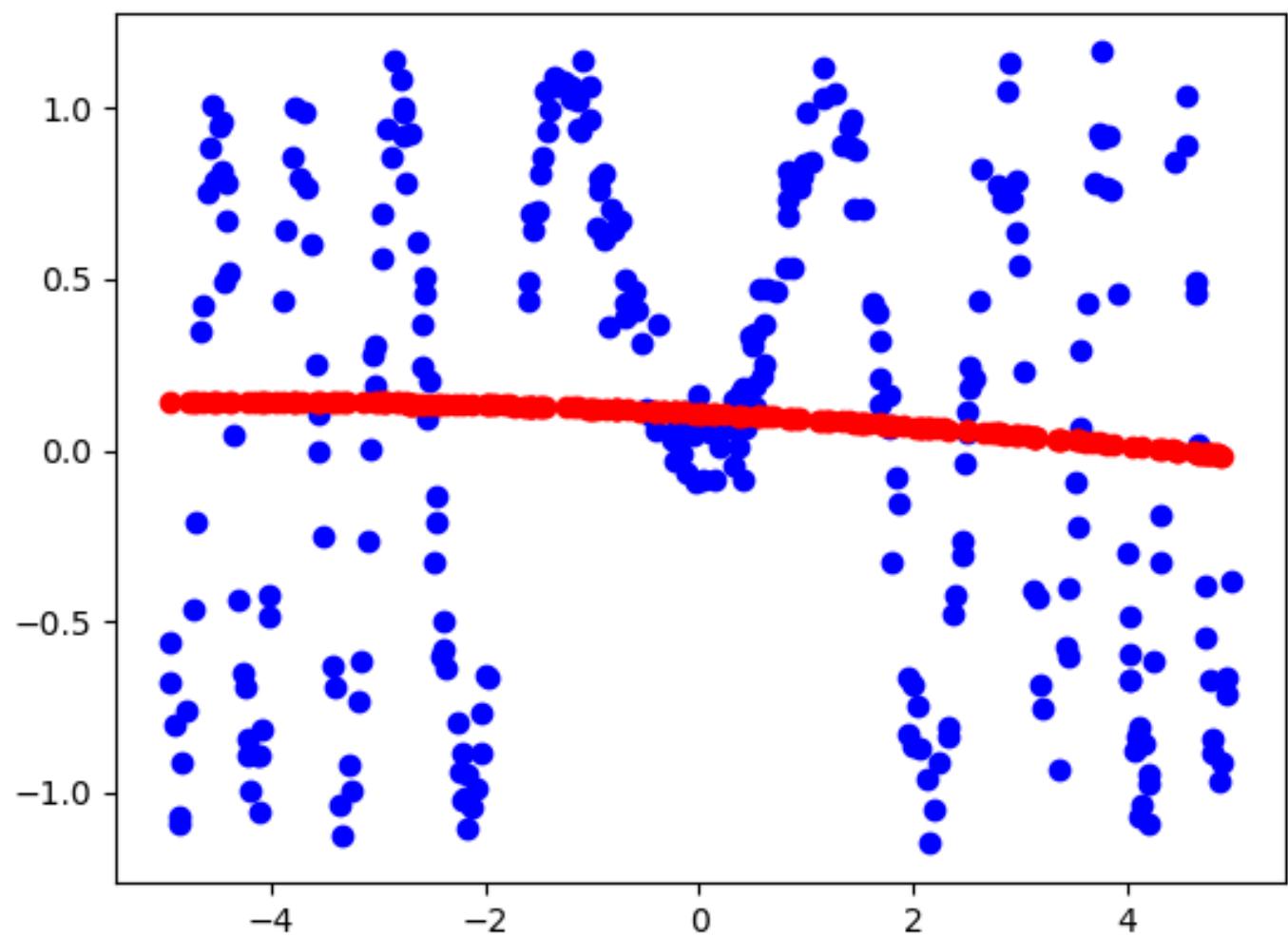












3) (a)

$$1) J(\theta) = \frac{1}{2} \sum_{i=1}^n (\theta_i \hat{x}_i - y_i)^2$$

$$2) \theta^{t+1} = \theta^{t-1} - \eta \nabla J(\theta^t)$$

(b) under fitting

$$mse = 0.33053126821375245$$

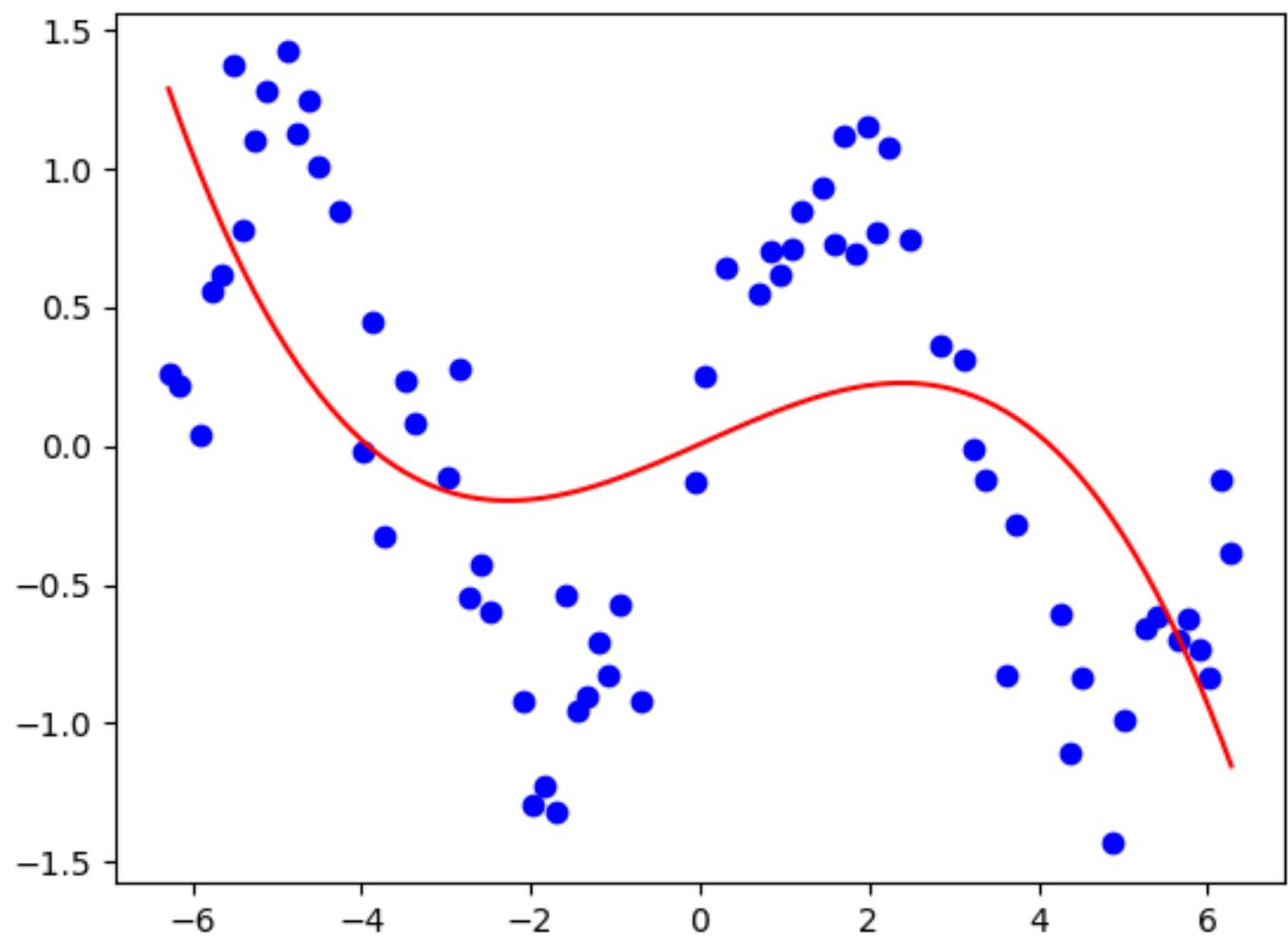
(C) best tau: 0.05

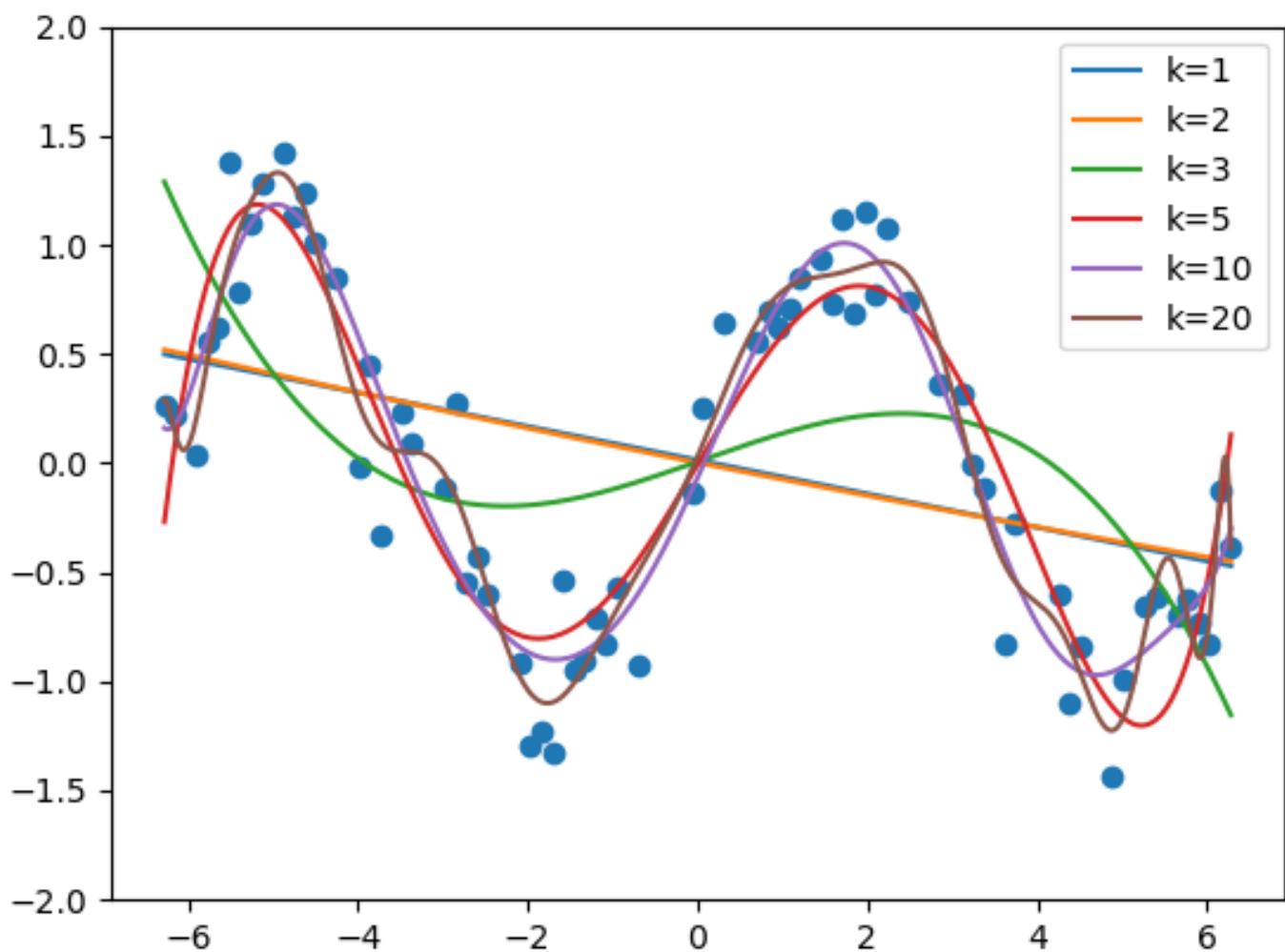
$$mse = 0.01690617458$$

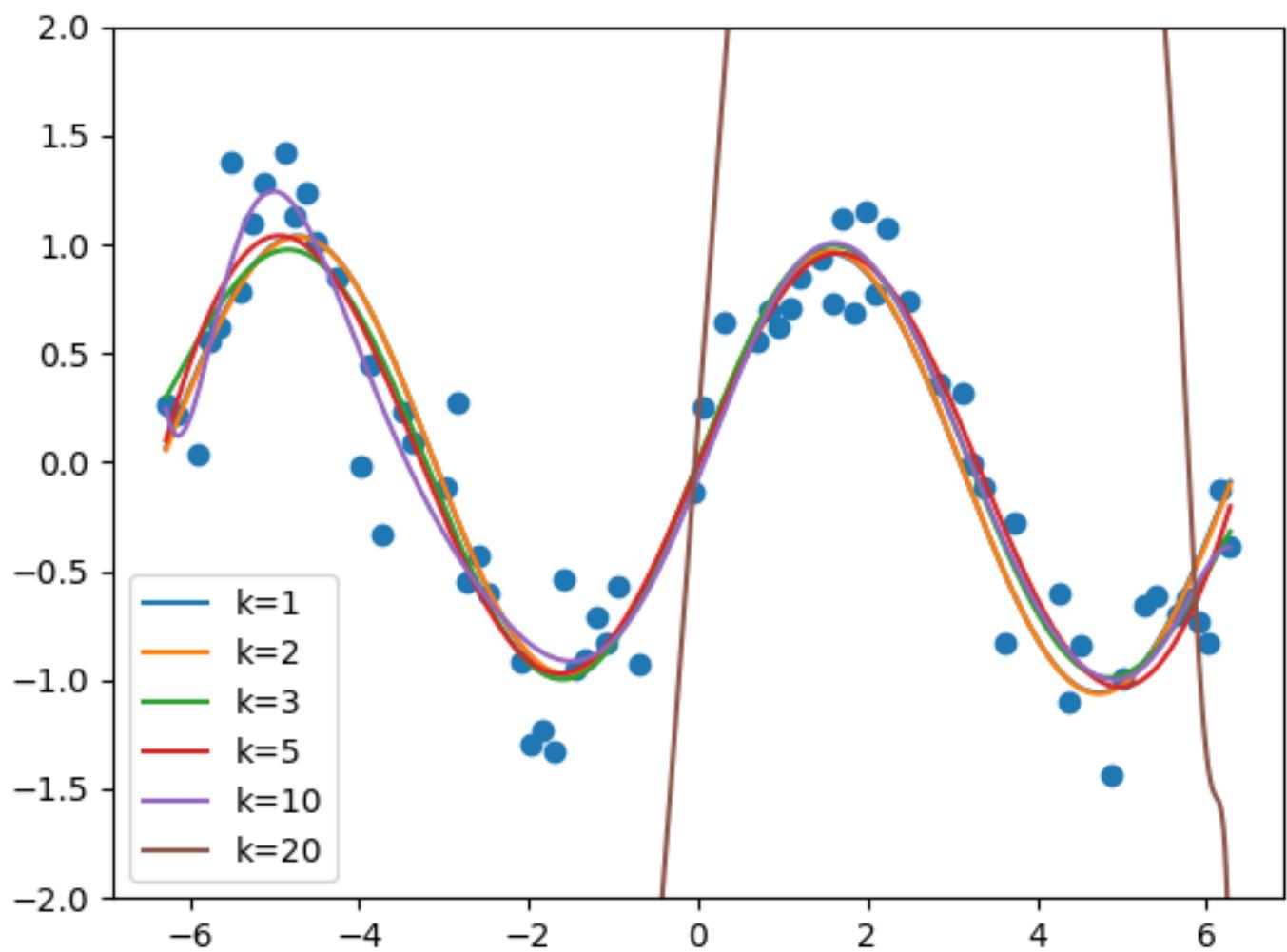
(G) for small dataset

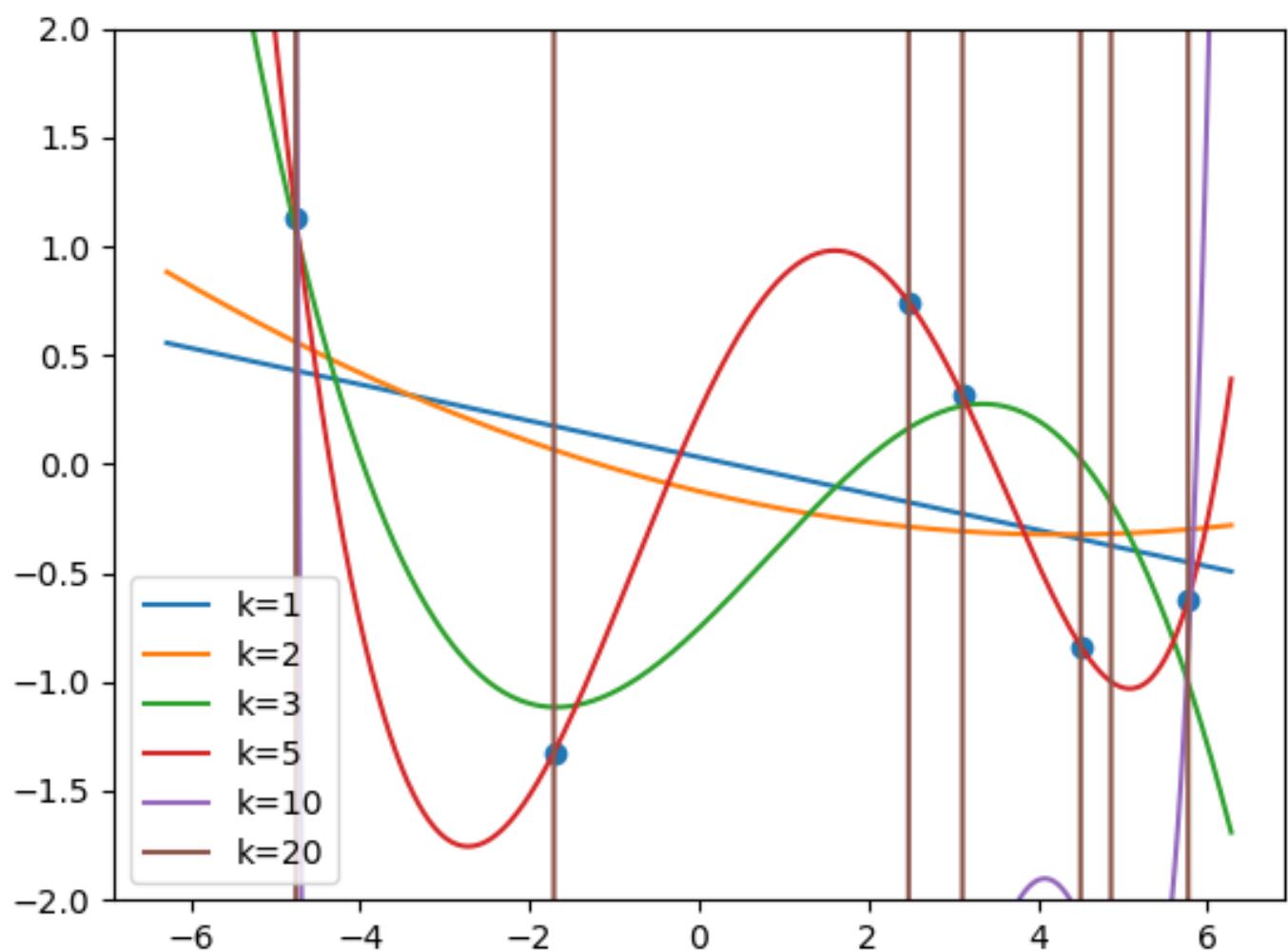
fitting change drastically when k changes

it goes from underfit to overfit as k \nearrow .









$$4. a) \Rightarrow J(\beta) = 0$$

$$\rightarrow x\beta - y = \vec{0}$$

$$x\beta = y.$$

$$x^T(x x^T)^{-1} x \beta = x^T(x x^T)^{-1} y$$

~~$$x^T(x^T)^{-1} x \beta = x^T(x x^T)^{-1} y$$~~

~~$$\beta = x^T(x x^T)^{-1} y$$~~

~~$$\text{further } x(\beta + N(x)) = x\beta + x \cdot \cancel{N(x)}$$~~

when ~~C^T~~

let C to be the null space of X .

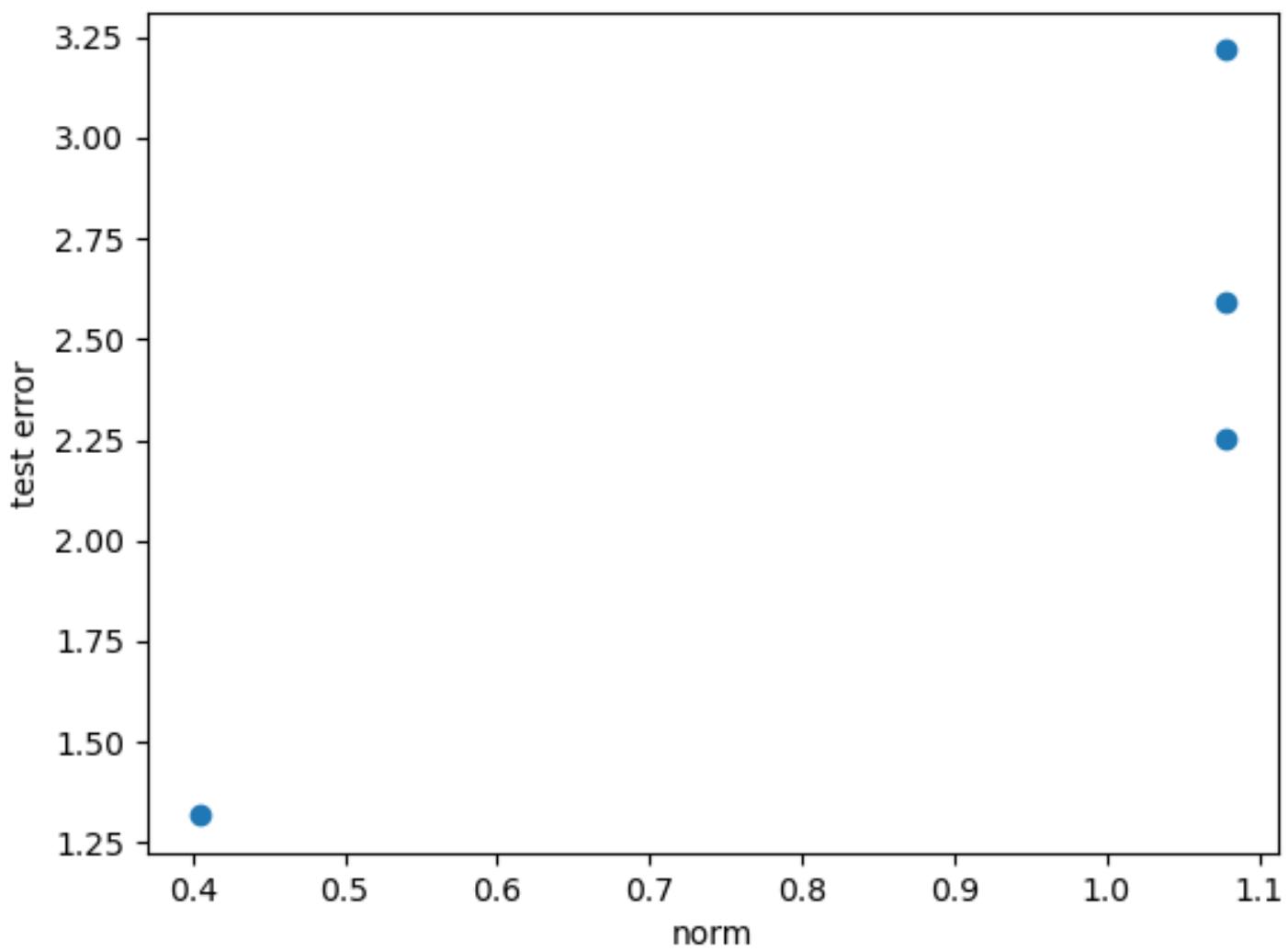
$$X \cdot (\beta + C) = x\beta + X \cdot C = x\beta.$$

$$\therefore \beta = x^T(x x^T)^{-1} y + N(x)$$

b)

$$\|\beta\|_2 = \|x^T(x x^T)^{-1} y\|_2 + \|C\|_2$$

$$= \|R\|_2 + \|C\|_2$$



$$4.(e) J(\theta, \phi) = \frac{1}{4n} \sum_{i=1}^n (x^{(i)\top} (\theta^{(i)} - \phi^{(i)}) - y^{(i)})^2$$

by definition

$$\hookrightarrow = \frac{1}{4n} \|X(\theta^{(i)} - \phi^{(i)}) - Y\|_2^2.$$

$$\therefore = \frac{1}{2} J(\beta) \quad \text{where } \beta = \theta^{(i)} - \phi^{(i)}.$$

$$= \frac{1}{2} \cdot \frac{1}{2n} \|X\beta - Y\|_2^2$$

base on (a) $J(\beta)$ has many infinite many optimal solution

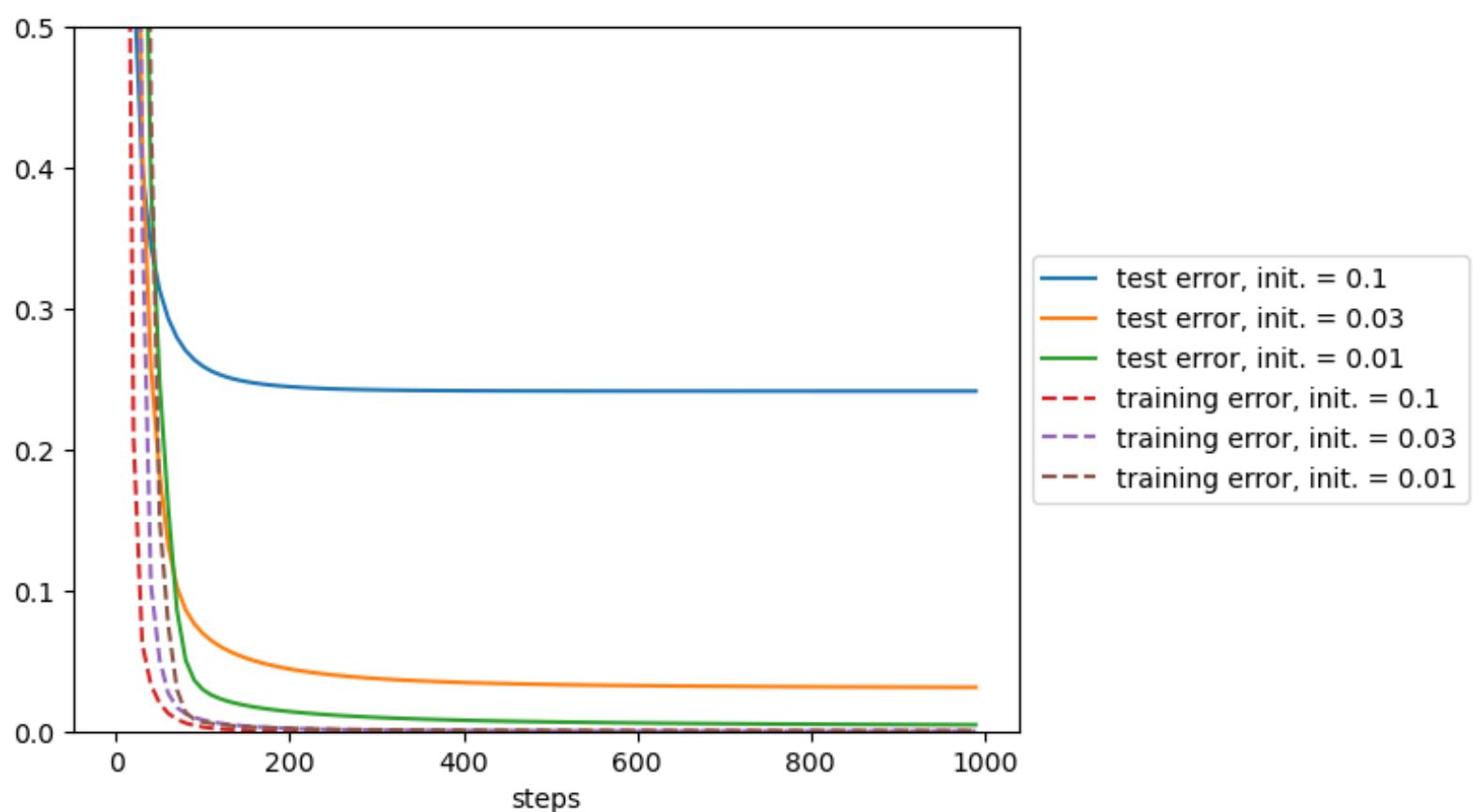
$$\therefore \theta^{(i)} - \phi^{(i)} \text{ has } \forall \quad \forall \quad \forall \quad \forall$$

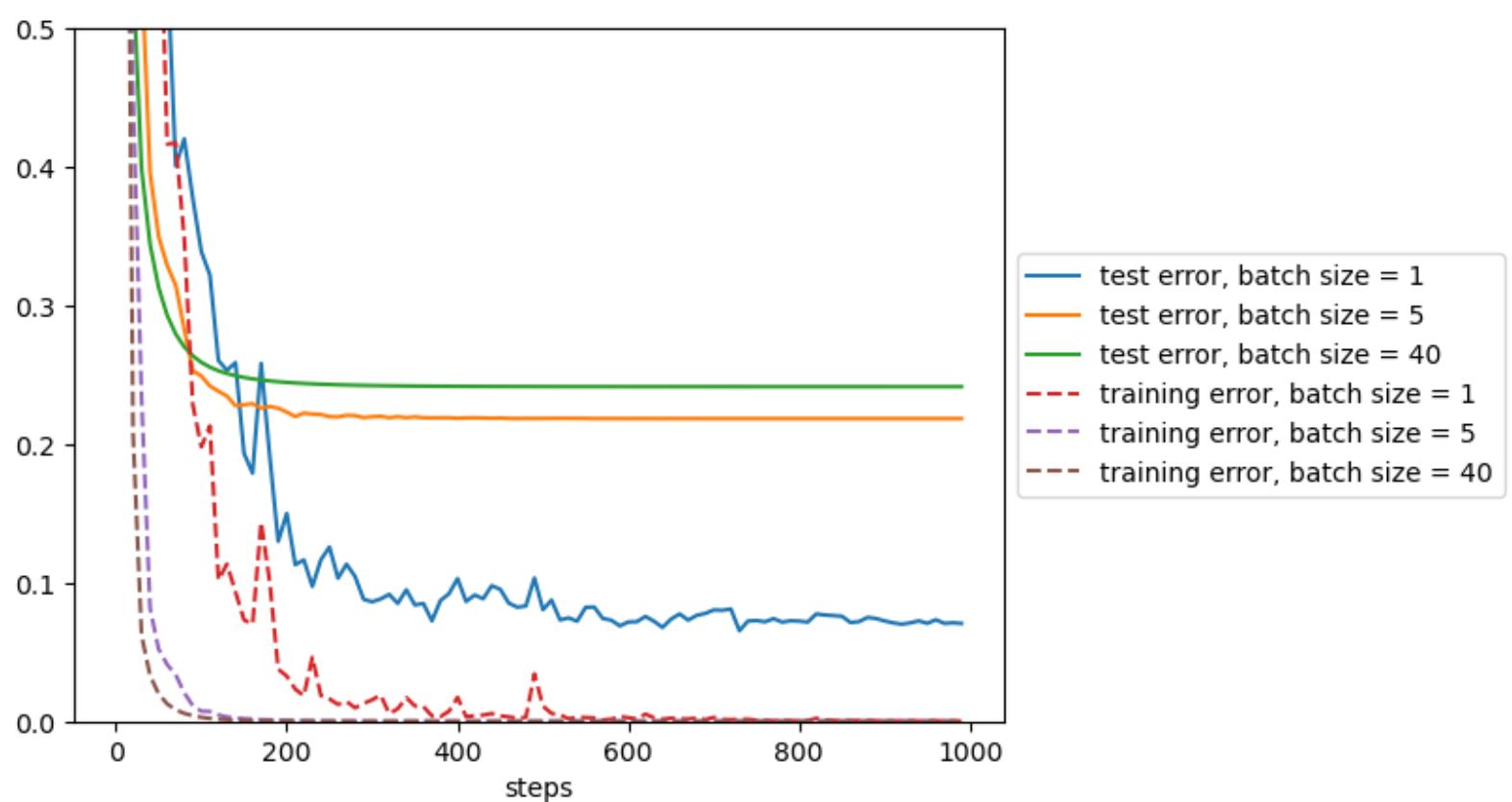
$$\therefore \theta^{(i)}, \phi^{(i)} \text{ has } \forall \quad \forall \quad \forall \quad \forall$$

4.(f) all can fit training set very well but
 $\text{init} = 0.01$ has the best test error

(4)(g) for the same init = 0.1

SGD has a lower test error
 thus generalize better.





$$\begin{aligned}
 5. \textcircled{a}) \quad J_{\lambda}(\beta) &= \frac{1}{2} \|x\beta - y\|_2^2 + \frac{\lambda}{2} \|\beta\|_2^2 \\
 &= \frac{1}{2} ((x\beta - y)^T (x\beta - y) + \lambda \beta^T \beta) \\
 &= \frac{1}{2} (\beta^T x^T x \beta - \beta^T x^T y - y^T x \beta + y^T y + \lambda \beta^T \beta) \\
 &= \frac{1}{2} (\beta^T x^T x \beta - 2y^T x \beta + \lambda \beta^T \beta).
 \end{aligned}$$

$$\nabla J_{\lambda}(\beta) = \frac{1}{2} (2(x^T x)\beta - 2(y^T x)^T + 2\lambda I\beta) = 0$$

$$(x^T x + \lambda I)\beta = x^T y$$

$$\beta = (x^T x + \lambda I)^{-1} x^T y$$

