

### Answer 1

Assuming the norm is defined on the domain  $\mathbb{D}$

Given  $x, y \in \mathbb{D}$ , and  $\lambda \in \mathbb{R}^1, \lambda \geq 0$

$$f(\lambda x + (1 - \lambda)y) = \|\lambda x + (1 - \lambda)y\|$$

$$f(\lambda x + (1 - \lambda)y) \leq \|\lambda x\| + \|(1 - \lambda)y\|$$

$$f(\lambda x + (1 - \lambda)y) \leq \lambda\|x\| + (1 - \lambda)\|y\|$$

$$f(\lambda x + (1 - \lambda)y) \leq \lambda f(x) + (1 - \lambda)f(y)$$

This means the function is convex.

## Answer2

Consider the probability bound for any of the samples lying in the set  $\mathcal{S}$ .

Focus on the  $j^{th}$  dimension of  $x$  and  $x_0$ ,

If  $x_{0,j} \in [-1, 1]$

$$P(\|x_{i,j} - x_{0,j}\|_\infty \leq 0.99) = \frac{\text{length of the real number line between } [-1,1] \text{ which is within } 0.99 \text{ distance of } x_{0,j}}{\text{length of the real number line between } [-1,1]}$$

$$P(\|x_{i,j} - x_{0,j}\|_\infty \leq 0.99) = \begin{cases} \frac{2*0.99}{2} = 0.99 & \text{if } x_{0,j} < 0.01 \\ \frac{0.99+(1-|x_{0,j}|)}{2} \leq 0.99 & \text{if } x_{0,j} \geq 0.01 \end{cases}$$

If  $x_{0,j} \in (-\infty, -1) \cup (1, \infty)$

$$P(\|x_{i,j} - x_{0,j}\|_\infty \leq 0.99) = \begin{cases} \frac{0}{2} = 0 & \text{if } |x_{0,j}| > 1.99 \\ \frac{|x_{0,j}|-1}{2} < 0.495 & \text{if } |x_{0,j}| \leq 1.99 \end{cases}$$

This means,

$$P(\|x_{i,j} - x_{0,j}\|_\infty \leq 0.99) \leq 0.99$$

Now for the d-dimensional case, for any sample  $x_i \in \mathcal{Q}$  to also  $\in \mathcal{S}$ ,  $\|x_{i,j} - x_{0,j}\|_\infty \leq 0.99$  should hold  $\forall j \in [d]$ . Now each dimension being independent of the other we can write,

$$P(\|x_i - x_0\|_\infty \leq 0.99) \leq (0.99)^d$$

Now for all the  $N$  samples to hold this true, given they are all sampled with uniform probability i.e. the samples are independent of each other

$$P(\|x_i - x_0\|_\infty \leq 0.99 \forall i \in [N]) \leq (0.99)^{dN}$$

Now,

$$P((\|x_i - x_0\|_\infty \leq 0.99 \forall i \in [N])^c) \geq 1.0 - (0.99)^{dN}$$

Also the event

$$(\|x_i - x_0\|_\infty \leq 0.99 \forall i \in [N])^c \subset (\|x_i - x_0\|_\infty > 0.99)$$

This means,

$$P(\|x_i - x_0\|_\infty > 0.99 \forall i \in [N]) \geq P((\|x_i - x_0\|_\infty \leq 0.99 \forall i \in [N])^c) \geq (1.0 - (0.99)^{dN})$$

Now given that  $N = d^{O(1)}$ , that means  $dN = d^{O(1)+1}$

$$P(\|x_i - x_0\|_\infty > 0.99 \forall i \in [N]) \geq (1.0 - (0.99)^{d^{O(1)+1}}) \geq (1.0 - (0.99)^d)$$

Now one can write  $(0.99)^d = \exp(d * \ln 0.99)$

$$P(\|x_i - x_0\|_\infty > 0.99 \forall i \in [N]) \geq (1.0 - \exp(d * \ln 0.99)) \geq (1.0 - \exp(-d))$$

$$P(\|x_i - x_0\|_\infty > 0.99 \forall i \in [N]) \geq (1.0 - \exp(-\Omega(d)))$$

For the optimization algorithm to work we need to sample a point within  $\delta$  neighborhood of  $x^* \in D$ . The above bound can be written for any  $\delta$  as

$$P(\|x_i - x_0\|_\infty > \delta \forall i \in [N]) \geq (1.0 - (\delta/\epsilon)^d)$$

This means the probability of atleast one point lying in the  $\delta$  range is;

$$P((\|x_i - x_0\|_\infty > \delta \forall i \in [N])^c) \leq (\delta/\epsilon)^d$$

Clearly this probability decreases exponentially in  $d$ , so it is not possible to perform the optimization using the given algorithm with high probability by selecting  $O(d)$  points.

**Answer3.1**

$$f(x_t) = f(x_{t-1} - \eta \tilde{\nabla} f(x_{t-1})) = (x_{t-1} - \eta \tilde{\nabla} f(x_{t-1}))^2$$

$$\begin{aligned} f(x_t) &= x_{t-1}^2 + \eta^2 (\tilde{\nabla} f(x_{t-1}))^2 - 2\eta \tilde{\nabla} f(x_{t-1}) x_{t-1} \\ f(x_t) &= x_{t-1}^2 + \eta^2 (\nabla f(x_{t-1}))^2 (1 + (\xi)^2 + 2\xi) - 2\eta \nabla f(x_{t-1}) (1 + \xi) x_{t-1} \end{aligned}$$

Also  $\nabla f(x) = 2 * x$

$$\begin{aligned} f(x_t) &= x_{t-1}^2 (1 + 4\eta^2 (1 + (\xi)^2 + 2\xi) - 4\eta(1 + \xi)) \\ f(x_t) &= f(x_{t-1}) (1 + 4\eta^2 (1 + (\xi)^2 + 2\xi) - 4\eta(1 + \xi)) \end{aligned}$$

Now for a given  $x_{t-1}$ ,

$$\mathbb{E}[f(x_t)] = f(x_{t-1}) \mathbb{E}[(1 + 4\eta^2 (1 + (\xi)^2 + 2\xi) - 4\eta(1 + \xi))]$$

Now  $f(x) > 0 \forall x$ , so for  $\mathbb{E}[f(x_{t+1})] \rightarrow \infty$ ,

$$\mathbb{E}[(1 + 4\eta^2 (1 + (\xi)^2 + 2\xi) - 4\eta(1 + \xi))] \geq 1$$

$$1 + 4\eta^2 (1 + \sigma^2) - 4\eta \geq 1$$

Given  $\eta \geq 0$

$$\eta \geq \frac{1.0}{1.0 + \sigma^2} \geq \frac{2}{\sigma}$$

### 3.2.1

Using given values,

$$\begin{aligned}x_{t+1} &= x_t - 2g_t \\g_t &= 0.9g_{t-1} + 0.2x_t\end{aligned}$$

Using  $g_t$  in the update rule for  $x_{t+1}$

$$x_{t+1} = x_t - 2(0.9g_{t-1} + 0.2x_t) = 0.6x_t - 1.8g_{t-1}$$

Also,

$$x_t = x_{t-1} - 2g_{t-1}$$

$$g_{t-1} = 0.5x_{t-1} - 0.5x_t$$

Using this we can write,

$$x_{t+1} = 1.5x_t - 0.9x_{t-1}$$

Combining with,

$$x_t = x_t + 0x_{t-1}$$

$$\begin{bmatrix} x_{t+1} \\ x_t \end{bmatrix} = \begin{bmatrix} 1.5 & -0.9 \\ 1.0 & -0 \end{bmatrix} \begin{bmatrix} x_t \\ x_{t-1} \end{bmatrix}$$

This means,

$$M = \begin{bmatrix} 1.5 & -0.9 \\ 1.0 & -0 \end{bmatrix}$$

$$\begin{bmatrix} x_{t+1} \\ x_t \end{bmatrix} = U \begin{bmatrix} 0.75000 + 0.58095i & -0.9 \\ 1.0 & -0.75000 - 0.58095i \end{bmatrix} U^T \begin{bmatrix} x_t \\ x_{t-1} \end{bmatrix}$$

$$\begin{bmatrix} x_{t+1} \\ x_t \end{bmatrix} = U \begin{bmatrix} (0.75000 + 0.58095i)^t & 0 \\ 0 & -(0.75000 - 0.58095i)^t \end{bmatrix}^t U^T \begin{bmatrix} x_1 \\ x_0 \end{bmatrix}$$

The eigen values of M are:  $0.75000 + 0.58095i$ ,  $0.75000 - 0.58095i$ , both have a norm of 0.948, and hence the values of  $x_t$  will decay exponentially with  $t$ , i.e.  $x_t = \exp(-\Omega(t))x_0$ .

Now  $f(x_t) = (x_t)^2 = \exp(-2\Omega(t))(x_0)^2 = \exp(-\Omega(t))f(x_0)$

### Answer 3.2.2

$$dh(t) = h(t+1) - h(t) = x_{t+1} - x_t = -\eta g_t$$

$$\frac{d^2h(t)}{d\tau(t)^2} = \frac{dh(t+1) - dh(t)}{(\tau(t+1) - \tau(t))^2} = -\eta(g_{t+1} - g_t)m^2$$

$$g_{t+1} = (1 - \gamma)g_t + \gamma \nabla f(x_{t+1}) = (1 - \gamma)g_t + 2a\gamma x_{t+1}$$

Using this, we can write

$$\frac{d^2h(t)}{d\tau(t)^2} = (\eta_0\gamma_0g_t - \eta_0\gamma_0ax_{t+1})$$

Also,

$$\frac{dh(t)}{d\tau(t)} = -\eta g_t m = -\eta_0 g_t$$

$$\frac{d^2h(t)}{d\tau(t)^2} = -\gamma_0 \frac{dh(t)}{d\tau(t)} - 2\eta_0\gamma_0ax_{t+1}$$

$$x_{t+1} = x_t - \eta g_t = h(t) - \frac{\eta_0}{m} g_t = h(t) - \frac{1}{m} \frac{dh(t)}{d\tau(t)}$$

$$\frac{d^2h(t)}{d\tau(t)^2} = -\gamma_0 \frac{dh(t)}{d\tau(t)} - 2\eta_0\gamma_0ah(t) - 2\eta_0\gamma_0a \frac{dh(t)}{d\tau(t)}$$

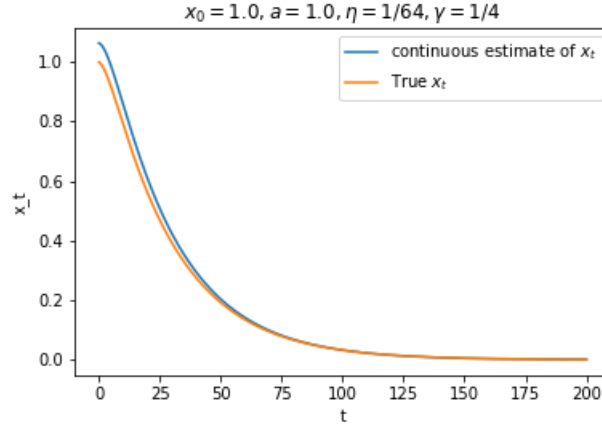
This means,

$$c_1(a, \eta_0, \gamma_0) = -\gamma_0, \quad c_2(a, \eta_0, \gamma_0) = -2\eta_0\gamma_0a, \quad \text{and} \quad \epsilon = -2\eta_0\gamma_0a \frac{dh(t)}{d\tau(t)}$$

### Answer 3.2.3

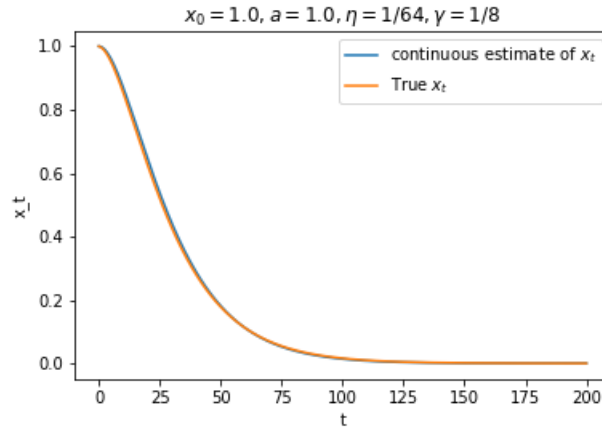
- $x_0 = 1, a = 1, \eta = \frac{1}{64}, \gamma = \frac{1}{4}$

Solution:  $h(t) = 1.270711 \exp(-0.0366117t) - 0.207107 \exp(-0.213388t)$



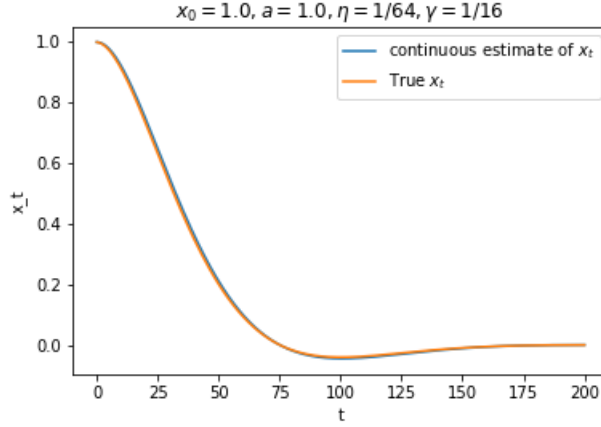
- $x_0 = 1, a = 1, \eta = \frac{1}{64}, \gamma = \frac{1}{8}$

Solution:  $h(t) = \frac{1}{16} \exp(-\frac{t}{16}) + (t + 16)$

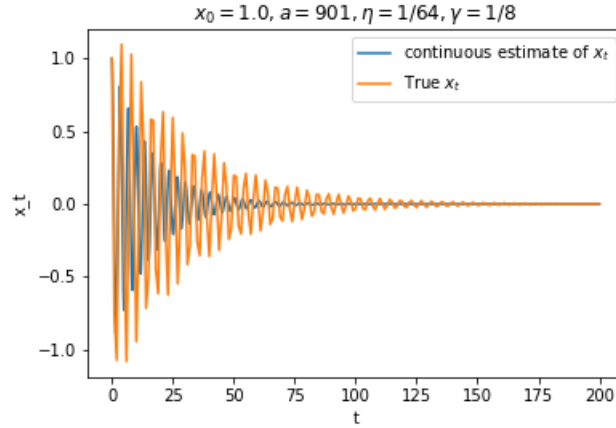


- $x_0 = 1, a = 1, \eta = \frac{1}{64}, \gamma = \frac{1}{16}$

Solution:  $h(t) = \exp(\frac{-t}{32})(\sin(\frac{t}{32}) + \cos(\frac{t}{32}))$



- $x_0 = 1, a = 901, \eta = \frac{1}{64}, \gamma = \frac{1}{8}$   
 Solution:  $h(t) = \frac{1}{30} \exp\left(\frac{-t}{16}\right) \left(\sin\left(\frac{15t}{8}\right) + 30 \cos\left(\frac{15t}{8}\right)\right)$



The estimate breaks down when the value of  $a$  is high, as this leads to a high value of  $\nabla f(x_{t+1})$ , and this the effect of past gradients also gets scaled which leads to very bad steps being taken, and hence there is a wild oscillation that can be seen in the last plot.

For a given  $a, \eta$ , for smaller  $\gamma$  the algorithm follows the past gradients more than the current one, due to this after reaching minima where current gradient is close to 0, the algorithm still takes bigger steps due to the past gradients, and this stabilizes slowly over there. If the value of  $\gamma$  is kept high, the algorithm will stabilize quicker when it reaches the minima.



#### Answer 4

Gradient of  $f(x)$  can be written as:

$$\nabla f(x) = [2ax_1, 2x_2]^T$$

Given the partial derivative of  $f(x)$  w.r.t  $x_i$  is only a function of  $x_i$ , the update rule for each coordinate can be written separately without considering the update rule for the other.

$$[x_t]_1 = [x_{t-1}]_1 - \frac{2a[x_{t-1}]_1}{\sqrt{\sum_{s \leq (t-1)} (2a[x_s]_1)}}$$

$$[x_t]_1 = [x_{t-1}]_1 - \frac{[x_{t-1}]_1}{\sqrt{\sum_{s \leq (t-1)} ([x_s]_1)}}$$

$$[x_t]_2 = [x_{t-1}]_2 - \frac{[x_{t-1}]_2}{\sqrt{\sum_{s \leq (t-1)} ([x_s]_2)}}$$

Clearly both the coordinates have the same functional form of the update rule (independent of a), and given that they both start at the same starting point  $[x_0]_i = c$ ,  $i \in [1, 2]$ , this means both the coordinates will stay the same  $\forall t$ , i.e.

$$[x_t]_1 = [x_t]_2$$

Now we can write  $\gamma$  as:

$$\begin{aligned}\gamma &= \sqrt{\frac{\sum_{s \leq t} [\nabla f(x_s)]_1^2}{\sum_{s \leq t} [\nabla f(x_s)]_2^2}} \\ \gamma &= \sqrt{\frac{(2a)^2 \sum_{s \leq t} [x_s]_1^2}{2^2 \sum_{s \leq t} [x_s]_2^2}}\end{aligned}$$

Now given

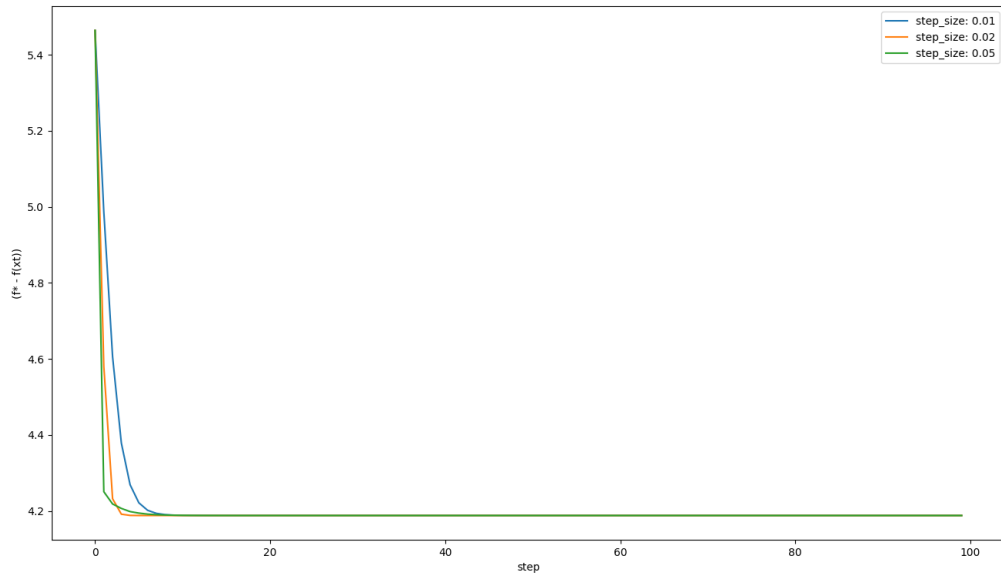
$$\begin{aligned}[x_s]_1 &= [x_s]_2 \quad \forall s \\ \sum_{s \leq t} [x_s]_1^2 &= \sum_{s \leq t} [x_s]_2^2\end{aligned}$$

Thus we can write  $\gamma$  as:

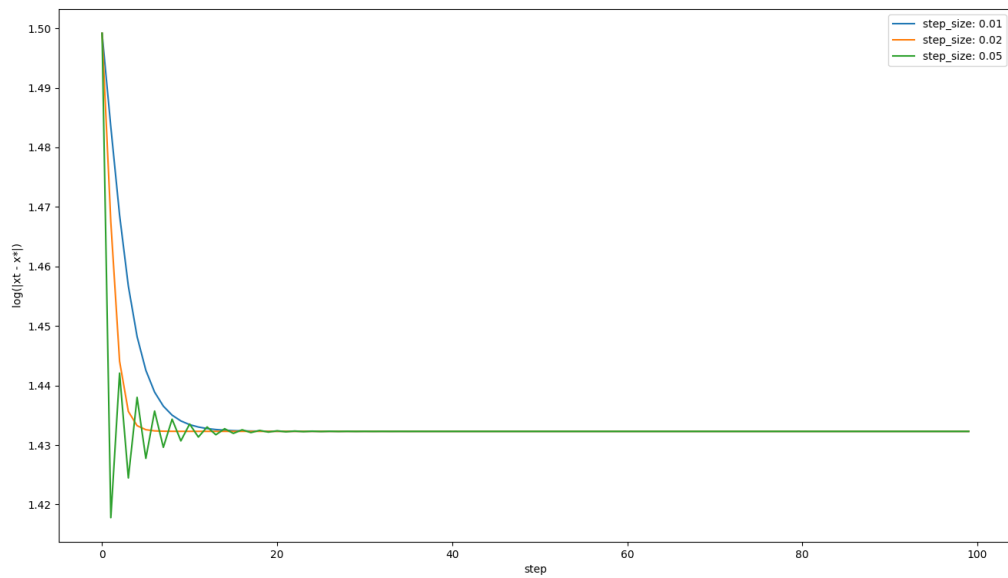
$$\gamma = a$$

## Answer 5.1

$$f^* - f(x_t)$$



$$|x_t - x^*|$$

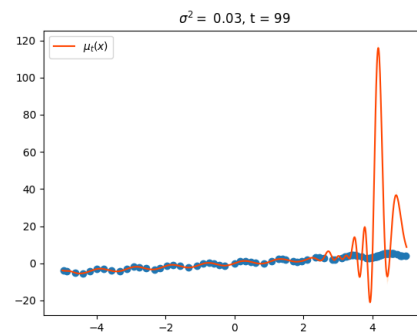
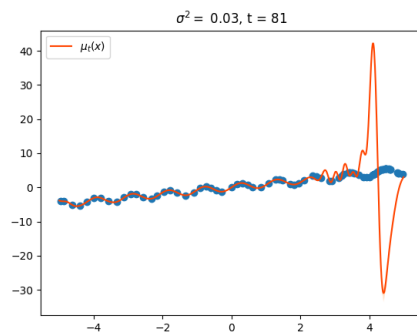
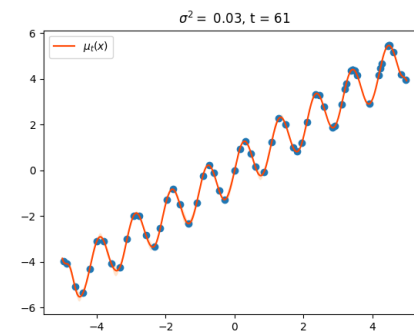
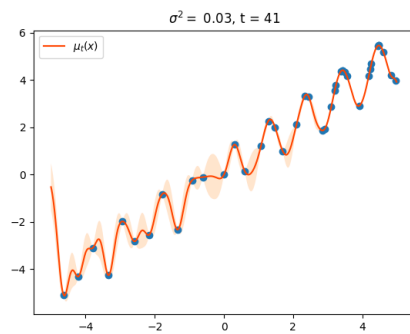
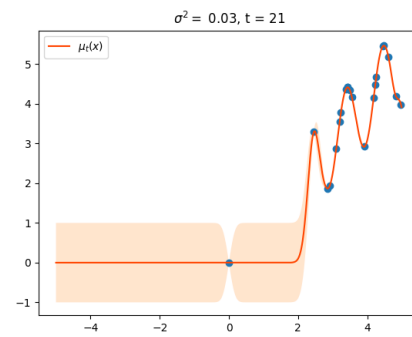
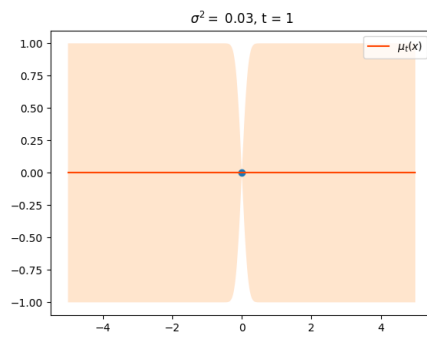


As the step size is increased, the algorithm reaches the optimum value faster, but for the highest

step size of 0.05, the algorithm starts to oscillate around the minima, but eventually decays to the same minima as the others.

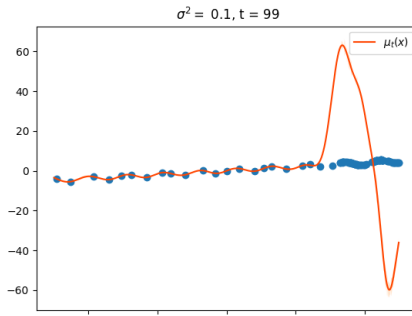
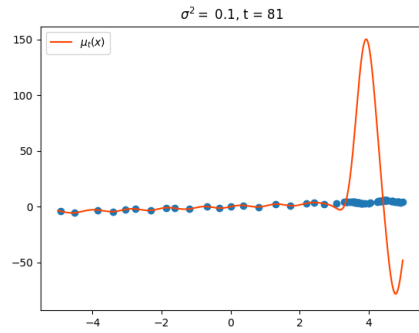
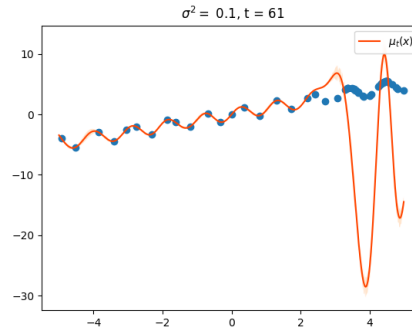
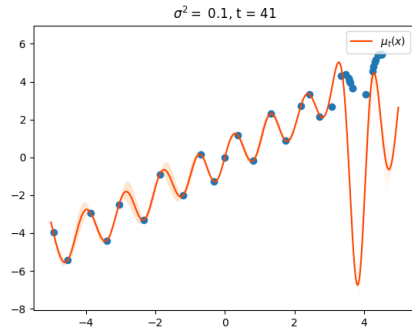
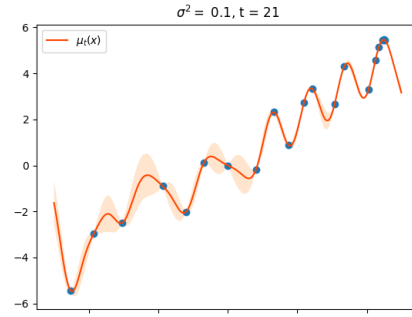
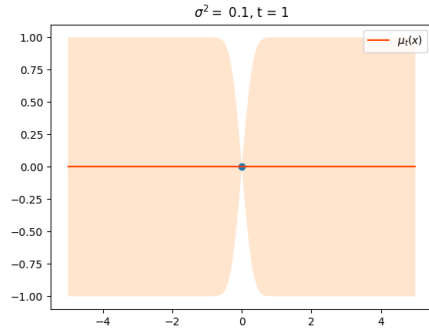
## Answer 5.2

- $\sigma^2 = 0.03$



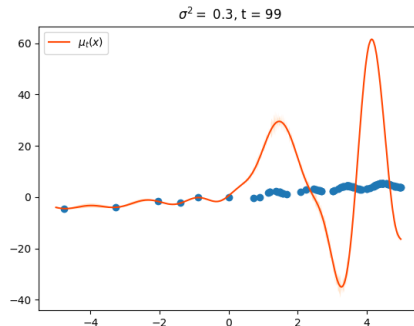
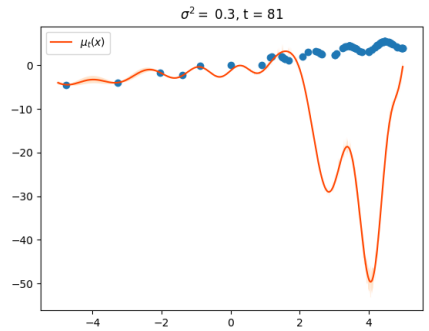
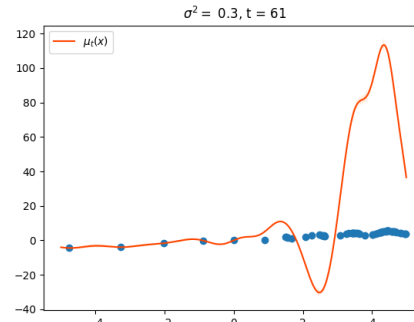
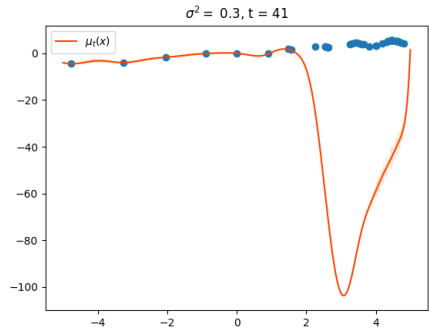
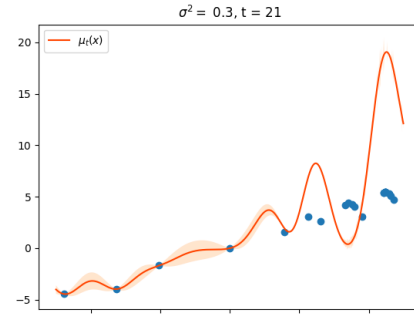
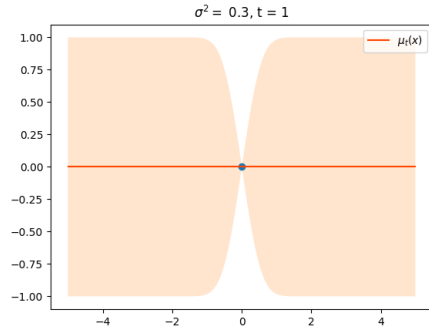
Final value:  $(x, f) = (4.4782464, 5.46450978)$

- $\sigma^2 = 0.1$



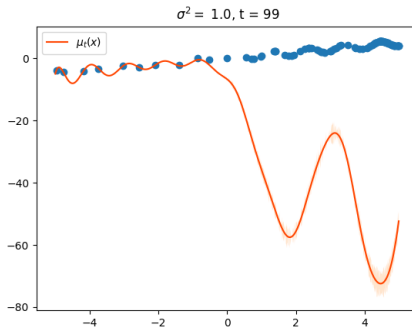
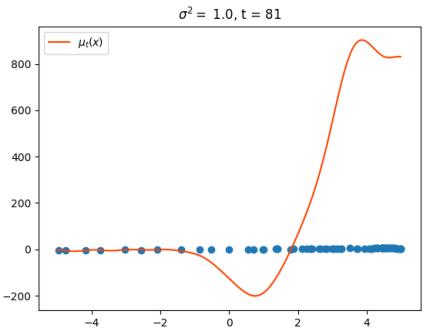
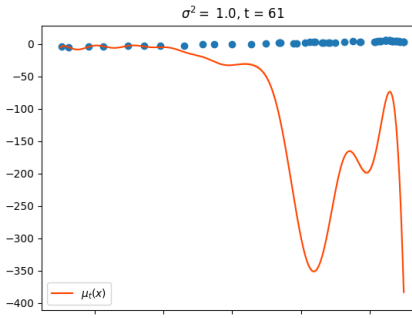
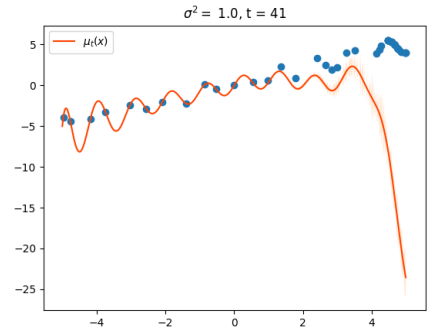
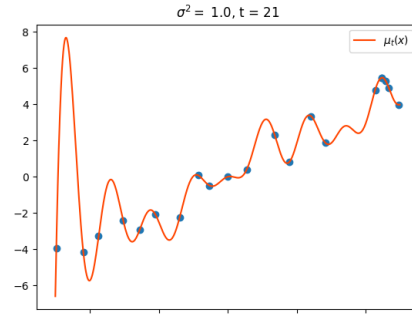
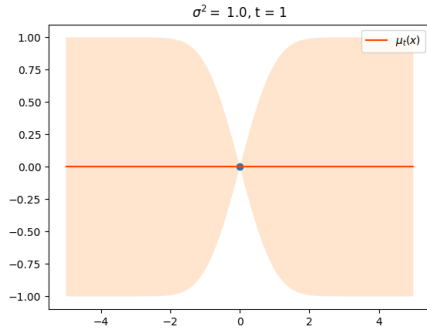
Final value:  $(x, f) = (4.48194939, 5.46429969)$

- $\sigma^2 = 0.3$



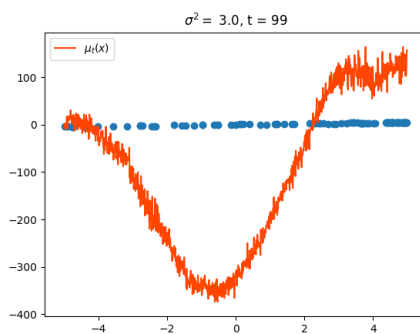
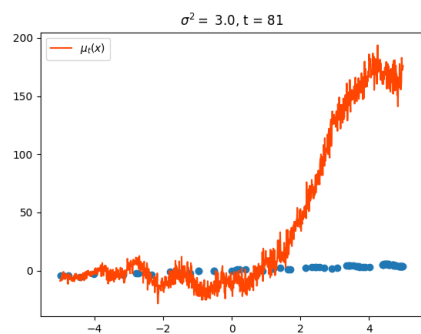
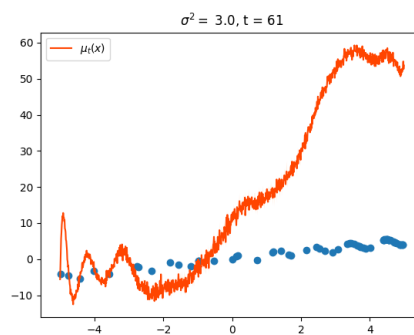
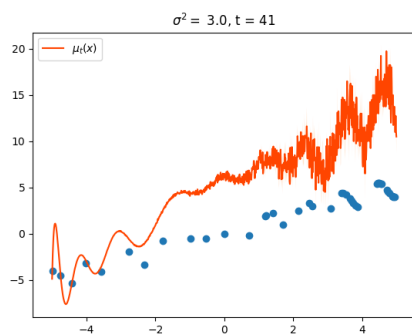
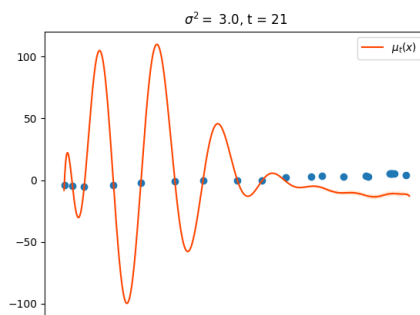
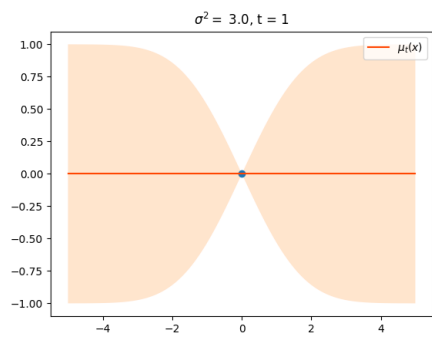
Final value:  $(x, f) = (4.47029106, 5.46331253)$

- $\sigma^2 = 1.0$



Final value:  $(x, f) = (4.4790136051787535, 5.4645061782361495)$

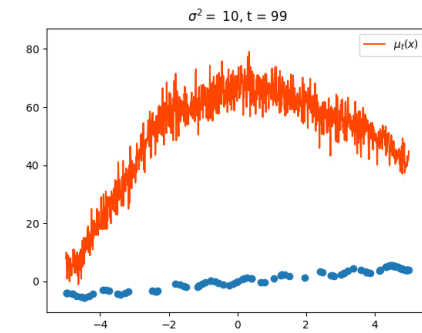
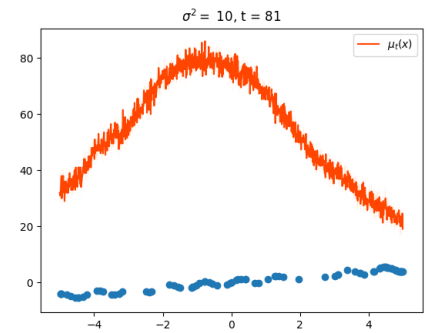
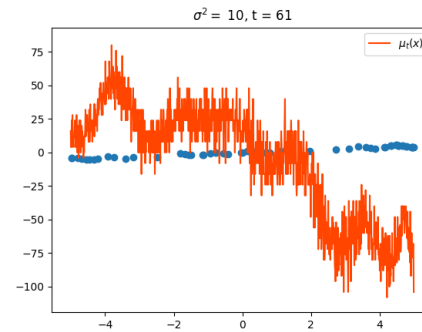
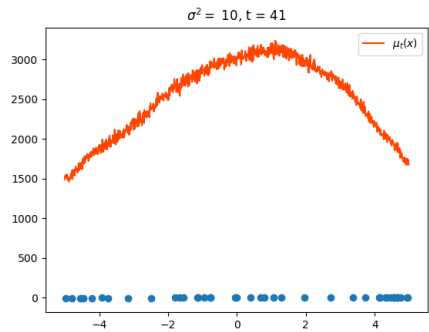
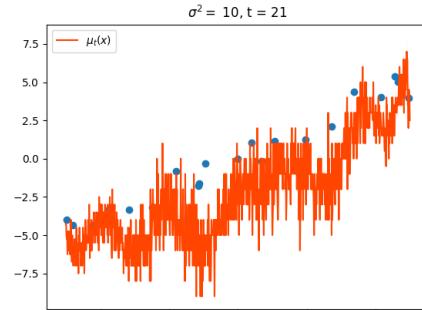
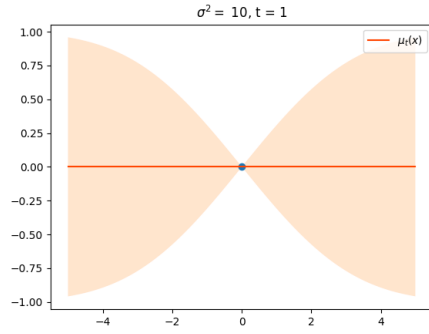
- $\sigma^2 = 3$



Final value:  $(x, f) = (4.48194939, 5.46429969)$

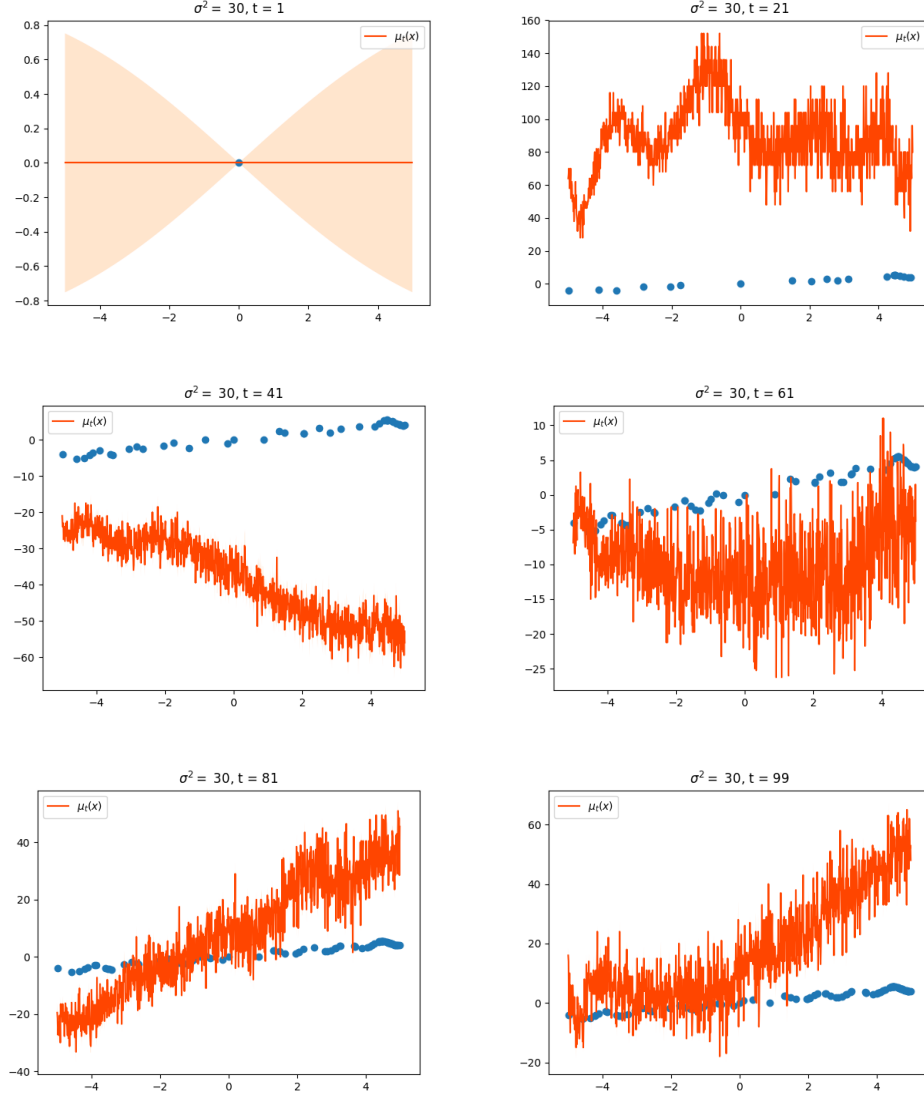
- $\sigma^2 = 10$





Final value:  $(x, f) = (4.470956012680244, 5.463499059630465)$

- $\sigma^2 = 30$

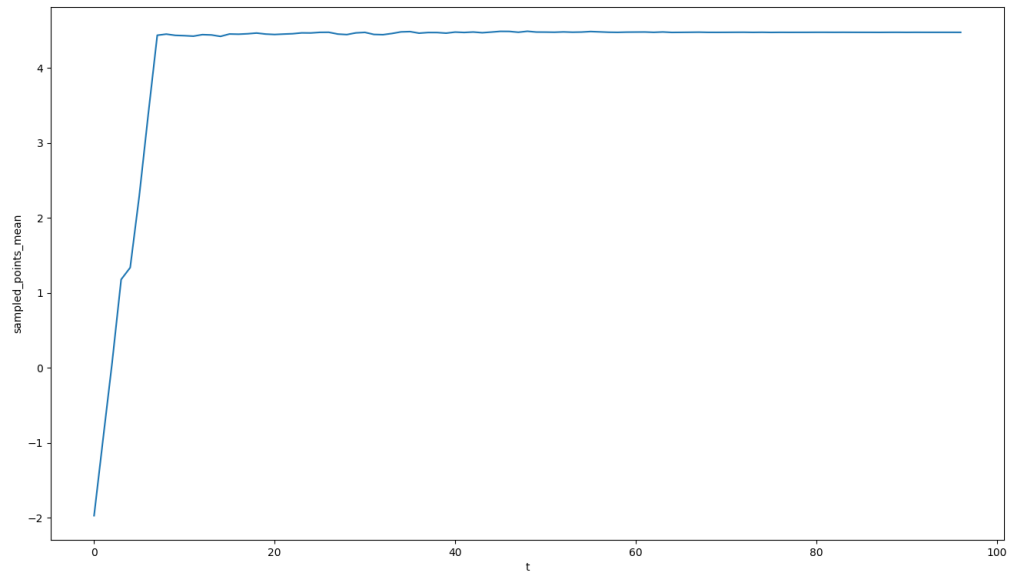


Final value:  $(x, f) = (4.481363086251285, 5.464365317322305)$

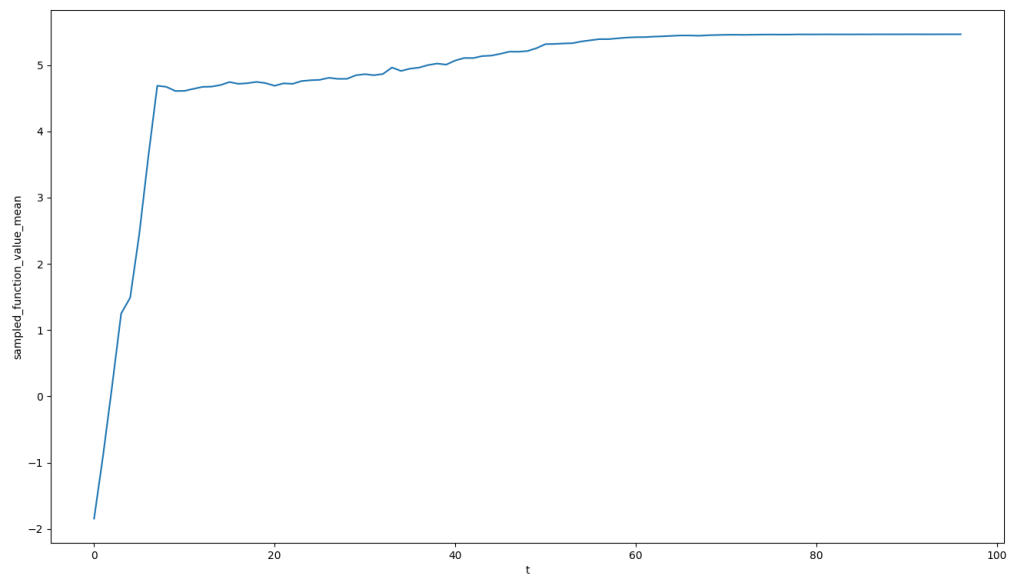
For all the sigma values, the optimization algorithm is able to find the optima value, but it seems to fit the function, only for the smaller sigma values. It could be because for higher sigma, the kernel gives the same inverse distance for all the points, which might be making it difficult to find the perfect fit.

### Answer 5.3

$$\frac{1}{T} \sum_{s \in [T]} x_s^{(t)}$$



$$\frac{1}{T} \sum_{s \in [T]} f(x_s^{(t)})$$



I used  $T = 500$  for the Metropolis-Hasting algorithm. The plots clearly show a very fast convergence.

to the optima value. Though one thing I noticed was that if the value of  $T$  is kept low, the samples produced by the Metropolis-Hasting algorithm are not very good, and lead to very noisy updates initially.