

Q2-b

Stochastic gradient decent is little different.
gradient decent iterates over all points and for large data that is very time consuming.

In SGD, we shuffle the data; Pick Random K ^{Points} data of them and run the gradient steps. In SGD, we also have the learning rate λ .
The gradient step is similar to Q2-a ($F_{\text{data}}(\theta)$).

Procedure SGD ($F_{\text{data}}, \theta(0)$)

$\theta \leftarrow \theta(0)$

while not converged do

for shuffle data $\{1, 2, \dots, N\}$ pick K data points

for $k \in \{1, 2, \dots, K\}$ do

$\theta_k \leftarrow \theta_k + \lambda F_{\text{data}}(\theta_k)$

return θ_k