

Gradient descent

At least a local minimum can be found



change $\vec{\theta}$ gradually as
 $\vec{\theta} \rightarrow \vec{\theta} + \Delta\vec{\theta}$

$$\Delta\vec{\theta} = -\xi \frac{\partial J}{\partial \theta}$$

$$(\Delta\theta_j = -\xi \cdot \frac{\partial J}{\partial \theta_j})$$

ξ learning rate

$J(\vec{\theta})$ decreases if ξ is sufficiently small

Decomposing the gradient

(cost function)

$$g(\omega) = \sum_{p=1}^P h(\omega, x_p, y_p)$$

i.e. sum of individual costs over each data points

for e.g. with Least square

$$h(\omega, x_p, y_p) = (x_p^\top \omega - y_p)^2$$

or with two class classification

$$h(\omega, x_p, y_p) = \log(1 + e^{-y_p x_p^\top \omega})$$

This is simple observation, that machine learning cost functions can be decomposed over individual data points, is important because it allows us to write the gradient itself as a summation

of the gradients of each of the P summands

$$\nabla g(\omega) = \nabla \left(\sum_{p=1}^P h(\omega, x_p, y_p) \right) = \sum_{p=1}^P \nabla h(\omega, x_p, y_p)$$

Thus when minimizing only cost function via gradient descent we can think about the k^{th} gradient step in terms of these individual gradients as

$$\begin{aligned}\omega^k &= \omega^{k-1} - \alpha_k \nabla g(\omega^{k-1}) \\ &= \omega^{k-1} - \alpha_k \sum_{p=1}^P \nabla h(\omega^{k-1}, x_p, y_p)\end{aligned}$$

where α_k is always an appropriately chosen length.

for eg

$$\vec{\theta} \rightarrow \vec{\theta} + \Delta \vec{\theta}$$

$$\Delta \vec{\theta} = -\epsilon \frac{\partial J}{\partial \theta}$$

↑
learning rate

Cost function $J(\vec{\theta}) = \frac{1}{2} \sum_{k \text{ all samples}} \left(y^{(k)} - f(\vec{x}^{(k)}, \vec{\theta}) \right)^2$

Vanilla or Batch GD $\frac{\partial J}{\partial \theta_i} = - \sum_{k \text{ all samples}} \left(y^{(k)} - f(\vec{x}^{(k)}; \vec{\theta}) \right) \cdot \frac{\partial f(\vec{x}^{(k)}, \vec{\theta})}{\partial \theta_i}$

The stochastic gradient descent

The idea of taking a sequence of gradient steps in each data point (as opposed to full / batch gradient step over the entire dataset) leads to a procedure called stochastic gradient descent.

The k^{th} iteration of stochastic gradient descent, sometimes called epoch, consists of P sequential pointwise gradient steps written as

$$\omega^{k,p} = \omega^{k,p-1} - \alpha_k \nabla_h (\omega^{k,p-1}, n_p, y_p)$$

$$p = 1, \dots, P.$$

(SGD)

In analogy with the k^{th} batch gradient step in

$$\begin{aligned} \omega^k &= \omega^{k-1} - \alpha_k \nabla g(\omega^{k-1}) \\ &= \omega^{k-1} - \alpha_k \sum_{p=1}^P \nabla_h (\omega^{k-1}, n_p, y_p) \end{aligned}$$

Batch gradient Descent.

here we have used the double superscript $\omega^{k,p}$, which means "the p^{th} individual step of the k^{th} stochastic gradient descent iteration". In this notation point of the k^{th} iteration is written as $\omega^{k,0}$ the corresponding sequence of P individual gradient steps as $(\omega^{k,1}, \omega^{k,2}, \dots, \omega^{k,P})$ and the final output of the k^{th} iteration. (ie the P^{th} stochastic step) as

$\omega^{k,p} = \omega^{k+10}$. After completing the k^{th} iteration we perform the $(k+1)^{\text{th}}$ iteration by cycling through the data again, taking individual gradient steps for $p = 1, \dots, P$.

It can be regarded as a stochastic approximation of gradient descent optimization, since it replaces the actual gradient (calculated from the entire data set) by an estimate thereof (calculated from a randomly selected subset of the data).

Especially in high dimensional optimization problems this reduces the very high computational burden, achieving faster iterations in exchange for lower convergence rate.

Mini batch GD

randomly choose a batch $\{\vec{x}^{(l)}, y^{(l)}\}_{l \in \Delta}$

$$\Delta \vec{\theta} = -\varepsilon \frac{\partial J_{\Delta}}{\partial \vec{\theta}}$$

$$J_{\Delta}(\vec{\theta}) = \frac{1}{2} \sum_{l \in \Delta} (\vec{x}^{(l)} - f(\vec{x}^{(l)}; \vec{\theta}))^2$$

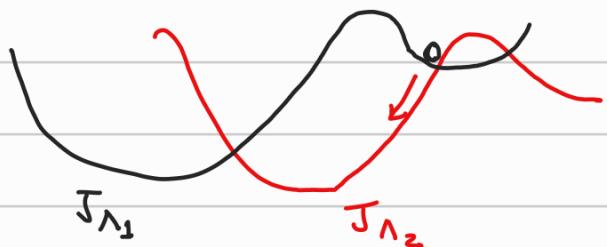
Advantages of Mini Batch GD

- ① If the batch is sufficiently large, J and J_{Δ} should not be different too much.
(as long as few outliers are dominant)

By using J_{λ} instead of J , we can reduce the cost at each step.

(eg total 10K samples in training data batch size
8, 16, or 32)

The location of minima depends on batch

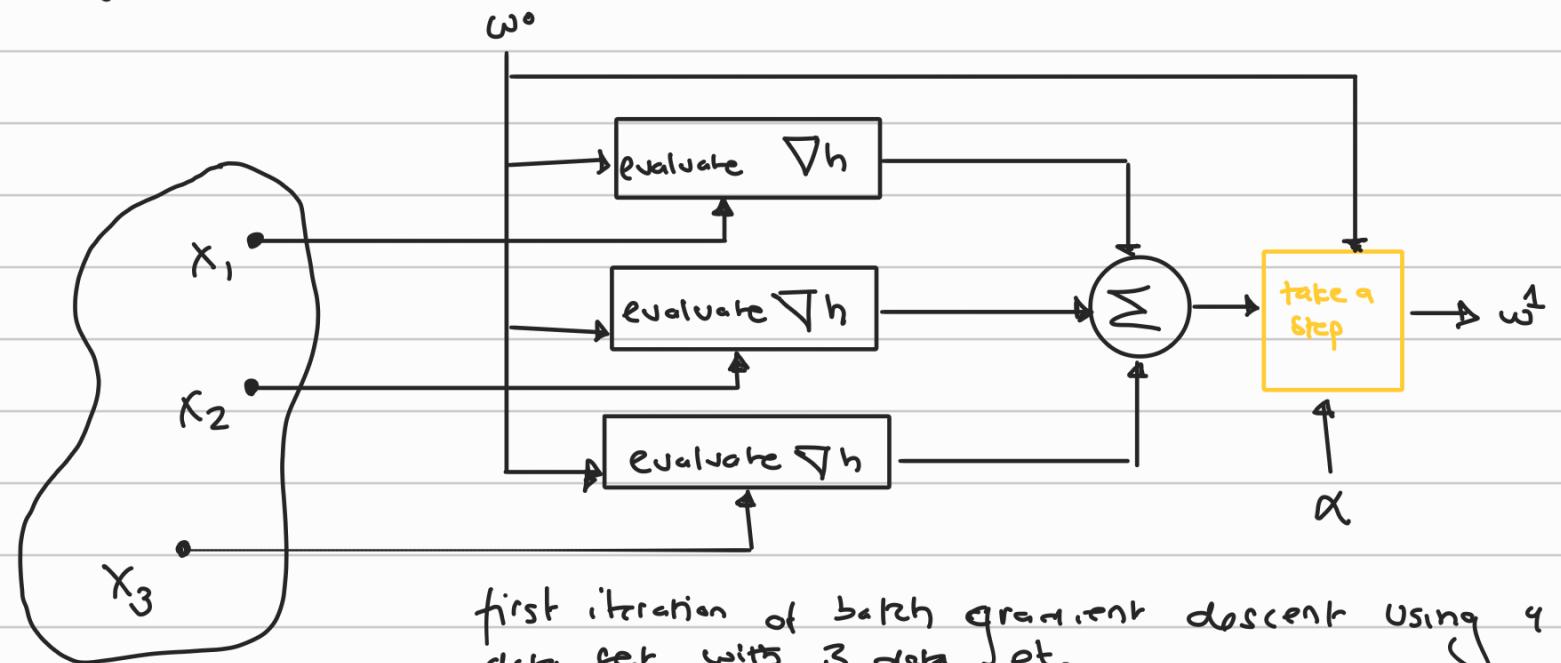


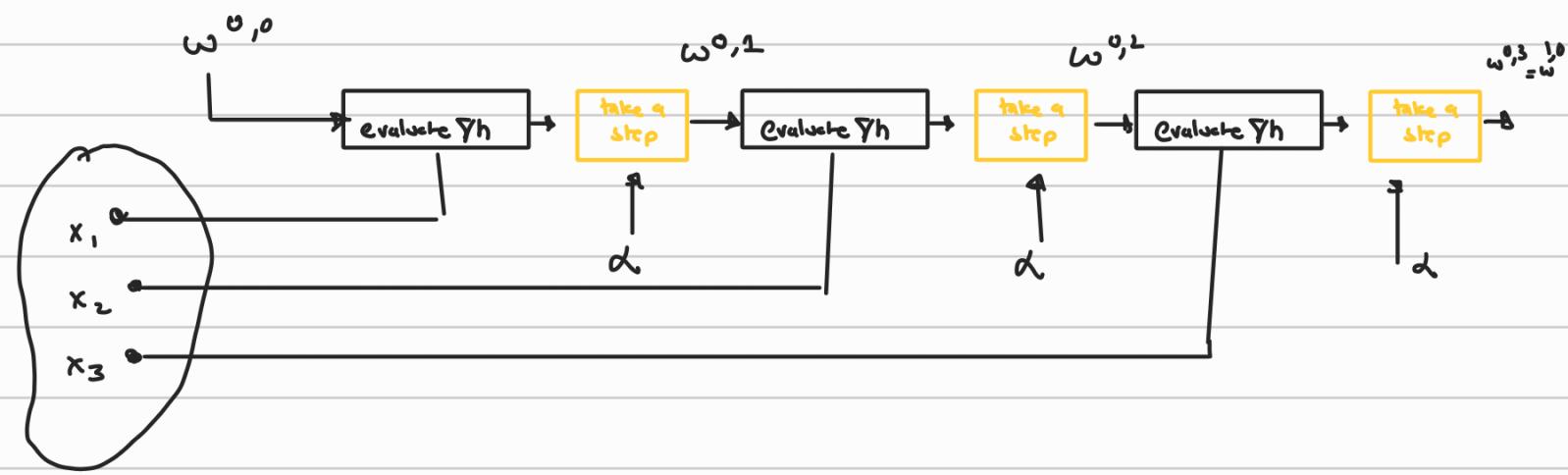
⇒ It is less likely to be trapped in a local minimum.

At least, we have a better chance to find a local minimum which is robust under a batch choice, which is likely to be a good one.

③ (similar to 2)

Implicit regularization effect preventing overfitting and improve generalization performance.





first iteration of stochastic gradient descent, using a data set

The standard or batch gradient descent we often use "step" and "iteration" interchangeably i.e. each iteration consists of one full gradient step in all P data points simultaneously

Conversely with stochastic method we refer to single "iteration" or epoch" as consisting of all P individual gradient steps one in each data point, executed sequentially for $p = 1, \dots, P$.

The mini batch gradient descent.

The cost function

$$g(\omega) = \sum_{p=1}^P h(\omega, x_p, y_p)$$

can be written more generally as

$$g(\omega) = \sum_{j=1}^J \sum_{p \in R_j} h(\omega, x_p, y_p)$$

where $R_1, R_2, R_3, \dots, R_J$ denote a certain partitioning of set $\{1, 2, \dots, P\}$ into J index sets.

These index sets divide the dataset into J subsets, each called a minibatch, where every data point belongs to one and only one mini-batch.

Analogous to

$$\begin{aligned}\nabla g(\omega) &= \nabla \left(\sum_{p=1}^P h(\omega, x_p, y_p) \right) \\ &= \sum_{p=1}^P \nabla h(\omega, x_p, y_p)\end{aligned}$$

we have

$$\begin{aligned}\nabla g(\omega) &= \nabla \left(\sum_{j=1}^J \sum_{p \in R_j} h(\omega, x_p, y_p) \right) \\ &= \sum_{j=1}^J \nabla \left(\sum_{p \in R_j} h(\omega, x_p, y_p) \right)\end{aligned}$$

where the gradient is now decomposed over each mini-batch (as opposed to each data point) and mini-batch gradient descent is then the algorithm wherein we take gradient steps sequentially using each mini-batch. Ideally we want all mini-batches to have the same size. — a parameter we call the batch size. or be as equally sized as possible when J does not divide P .

Notice, a batch size 1 turns mini-batch gradient descent into stochastic gradient descent, whereas a batch size of P turns it into the standard or batch gradient descent.