

①

Stochastic Gradient Descent

Gradient Descent \rightarrow very slow if # samples in training data is large.

Eg: #N > 3 million? \rightarrow GD will be very slow.

$$E(b_0, b_1, \dots, b_d) = \frac{1}{2N} \sum_{i=1}^N (\hat{y}_i - y_i)^2$$

Repeat until convergence {

$$b_j := b_j - \alpha \left[\frac{1}{N} \sum_{i=1}^N (\hat{y}_i - y_i) x_{ij} \right]$$

(for $j = 0, 1, 2, \dots, d$)

$$\frac{\partial E(b_0, b_1, \dots, b_d)}{\partial b_j} \downarrow$$

- loop for N to look at all the training samples.

$$N = 3000000000$$



Batch GD.

- looking at all training data/samples.
- read/fill these 3000000000 records into memory to compute the derivative term.
- very slow if N is very large.

Can we do faster?

- Yes: looking only at 1 training sample in each iteration.

Stochastic GD:

$$\hat{y}_i = b_0 x_0 + b_1 x_1 + \dots + b_d x_d$$

$$\text{cost}(B, (x_i, y_i)) = \frac{1}{2} (\hat{y}_i - y_i)^2$$

single data sample.

x	y
x_{11}	y_1
x_{21}	y_2
x_{\dots}	y_3
\vdots	\vdots
x_N	y

$$E(b_0, b_1, \dots, b_d) = \frac{1}{N} \sum_{i=1}^N \text{cost}(B, (x_i, y_i))$$

$$= \frac{1}{N} \sum_{i=1}^N \left[\frac{1}{2} (\hat{y}_i - y_i)^2 \right]$$

\rightarrow Scan through the training example \rightarrow modify param (b_0, b_1, \dots, b_d) one by one

Steps:

- ① Randomly shuffle dataset. } → (safer) Better to shuffle randomly to make the algo converge faster.

② Repeat {

- depends on Training data $\Rightarrow (1 \dots 10)$
- for large N, one iteration is enough.

for $i = 1, 2, \dots, N$ (training samples)

only 1 sample at a time.

{

$$b_j := b_j - \alpha (\hat{y}_i - y_i) \cdot x_{ij}$$

(for $j = 0, 1, 2, \dots, d$)

}

modify $(b_0, b_1, \dots, b_d) \rightarrow \dots \rightarrow \dots$
 $\uparrow \quad \quad \uparrow \quad \quad \uparrow$
 $(x_{11}, y_1), (x_{21}, y_2), (x_{31}, y_3) \dots$
 random samples

$$\frac{\partial}{\partial b_j} \text{cost}[B, (x_i, y_i)]$$

$$\downarrow$$

$$\frac{\partial}{\partial b_j} \left[\frac{1}{2} (\hat{y}_i - y_i)^2 \right]$$

$$= (\hat{y}_i - y_i) x_{ij}$$

- SGD doesn't converge as similar to GD.
- instead the parameters are closer to the global minimum.
good enough for practical applications.

B-GD	SGD
use all N samples in one iteration.	use only 1 sample in each iteration.

Mini Batch GD:

- in between GD and SGD.
- uses param "k" = mini-batch size. (smaller than N)
 $\swarrow \searrow$
 $10 \dots \alpha (2 \text{ to } 100)$

Example: $k=10 \rightarrow$ use 10 samples in updating param in each iteration.

\downarrow
 $(x_1, y_1) \dots (x_{10}, y_{10}) (x_{11}, y_{11}), (x_{12}, y_{12}) \dots$
 used in 1 iteration.

③ Mini Batch GD Algo:

$k = 10$; $N = 1000$ (samples)

Repeat { for $i = 1, 11, 21, 31, \dots, 991$ }

$$b_j := b_j - \alpha \frac{1}{10} \sum_{m=i}^{i+9} (\hat{y}_m - y_m) x_{mj}$$

(for $j = 0, 1, 2, \dots, d$)

}

sees only 10 training data.

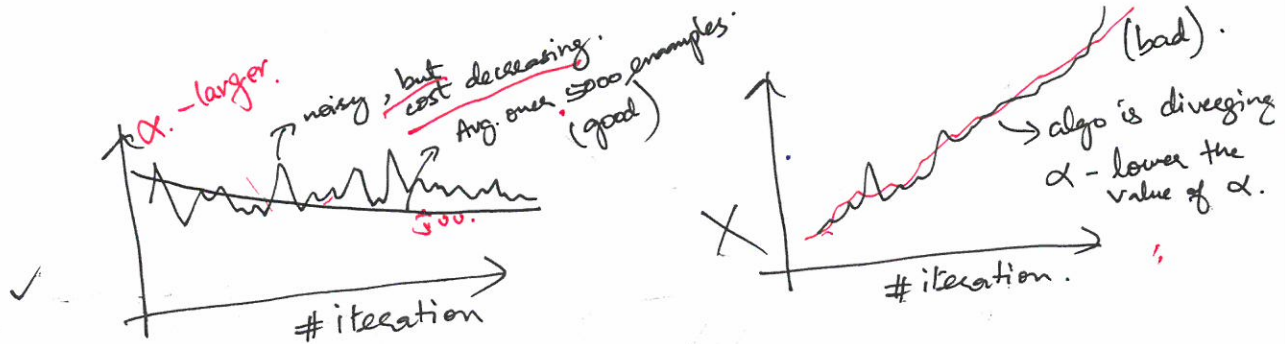
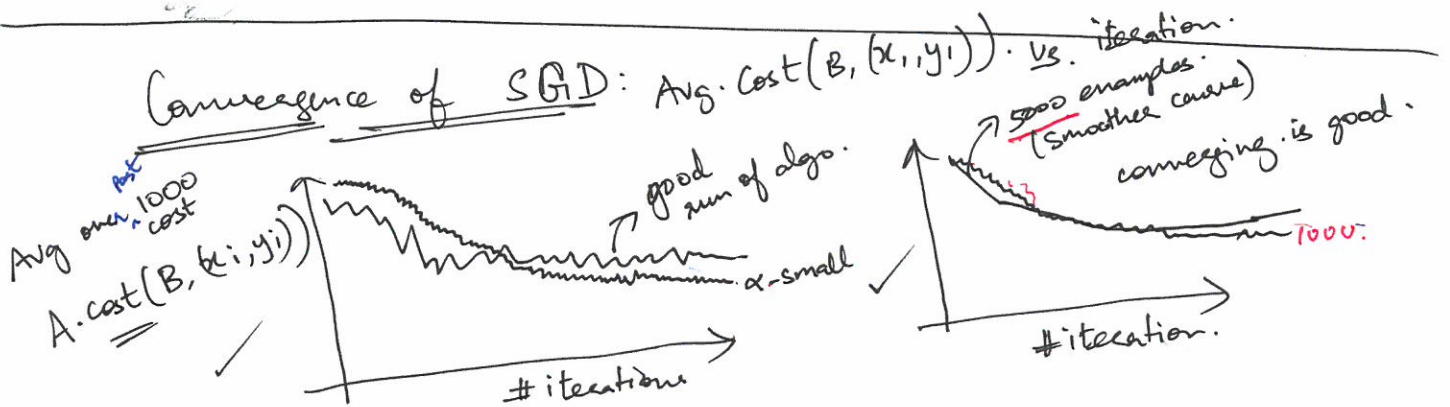
Why k -examples than 1-example?

Vectorization

- parallelize your gradient computation over $\frac{k}{10}$ samples.
- use linear algebra api's that uses vectorizations.

Disadv. of Mini-B GD:

- Additional parameters to tune 'k' - minibatch size.
- choose $k = 2$ to 100.



Note:- slowly decrease α over time \rightarrow in order for B to converge.

fix vs. $\alpha = \left(\frac{\text{constant 1}}{\text{iteration.No.} + \text{constant 2}} \right) \Rightarrow \text{const 1, const 2}$

extra parameters.

o.l.

