

7) Optimization Algorithms

1) Mini Batch Gradient Descent:

$$X = \underbrace{[x^{(1)}, x^{(2)}, x^{(3)} \dots x^{(1000)}]}_{X^{\{1\}} (n \times 1000)} \mid \underbrace{[x^{(1001)} \dots x^{(2000)}]}_{X^{\{2\}} (n \times 1000)} \mid \dots \mid \underbrace{[x^{(5000)}]}_{X^{\{5000\}} (n \times 1000)}$$

$$Y = \underbrace{[y^{(1)}, y^{(2)} \dots y^{(1000)}]}_{Y^{\{1\}} (1, 1000)} \mid \underbrace{[y^{(1001)} \dots y^{(2000)}]}_{Y^{\{2\}} (1, 1000)} \mid \dots \mid \underbrace{[y^{(5000)}]}_{Y^{\{5000\}} (1, 1000)}$$

Training set = 5.000.000 ise

ör, 5000 tane mini batch oluşturun. Her biri 1000 gözlem içerir.

* Mini batch size = 1000

{ } \leadsto k. batch

Repeat

for $t = 1, 2, \dots, 5000$

Forward prob on $X^{(t)}$ } ^{vectorized imple.}
(1000 example) from $X^{(t)}, Y^{(t)}$

$$\text{Compute cost } J^{(t)} = \frac{1}{1000} \sum_{i=1}^L \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Backprob to compute gradients

Epoch, is a word that means a single pass through the training set.

! With batch gradient descent, a single pass through the training set allows you to take only one gradient descent step.

! With mini-batch gradient descent, a single pass through the training set, allows you to take 5.000 gradient descent steps.

Gözlem sayımız çok fazla ise mini-batch gradient descent kullanalım. Aksi durumda batch gradient descent kullanacağız.

!!

Örneğin; elimizde 2000 training example olsun.

We can divide dataset of 2000 examples into batches of 500 then, it will take 4 iterations to complete one epoch.

↳

Kendi örneğimizle;

5.000.000 training example

1000 batch size

5.000 iterations

for completing 1 epoch.

Ek Bilgi :

* Exponentially Weighted Averages;

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

↓
↓

özell
Bugün

ortalama dün (Dünkü değer)

Son fiyatların (veya neyi ölçüyorsa) önemini artırarak ortalama hesaplar. Basit ortalama da tümünün önemi aynı iken burada sonda olanlara daha çok ağırlık verir.

$V_t \rightarrow$ Exponentially average over $\frac{1}{1-\beta}$ day's temperature.

↳ or: $\beta = 0.9$ ise : \downarrow last 10 days
(yaklaşık)

β çok büyük olursa; Önceki değere çok fazla önem vermiş olur.

β çok küçük olursa; bugünkü değere çok fazla önem vermiş olur.

$$V_t = \beta V_{t-1} + (1-\beta) \theta_t$$

$$V_{100} = 0.9 V_{99} + 0.1 \theta_{100}$$

$$V_{99} = 0.9V_{98} + 0.1\theta_{99}$$

$$V_{g8} = 0.9V_{g7} + 0.1\theta_{g8}$$

$$V_{100} = 0.1 \theta_{100} + 0.9 \theta_{99} \rightarrow (0.1 \theta_{99} + 0.9 \theta_{98}) \rightarrow 0.1 \theta_{98} + 0.9 \theta_{97}$$

$$V_{100} = 0.1 \theta_{100} + 0.1 \times 0.9 \times \theta_{99} + 0.1 \times (0.9)^2 \times \theta_{98} + 0.1 \times (0.9)^3 \times \theta_{97} + 0.1 \times (0.9)^4 \times \theta_{96} + \dots$$

Bias correction in exponentially weighted average =

$$V_0 = 0$$

$$V_1 = 0.98V_0 + 0.02\theta_1$$

$$V_2 = 0.98V_1 + 0.02\theta_2$$

$$= 0.0196\theta_1 + 0.02\theta_2$$

Bu şekilde yaptığımız durumda ilk günlerdeki değerler beklenenden her zaman az olacaktır.

↳ Bias Correction burada devreye girer;

V_t yerine; $\frac{V_t}{1-\beta^t}$ → örneğin: $t=2 : 1-\beta^t = 1-(0.98)^2 = 0.0396$

$$\frac{V_2}{0.0396} = \frac{0.0196\theta_1 + 0.02\theta_2}{0.0396}$$

Yani, Instead of taking V_t , take V_t divided by $1-\beta^t$. (t : current data)

2) Momentum (or Gradient Descent with momentum)

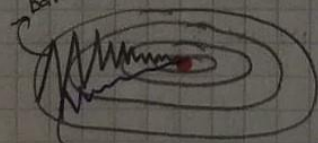
Basic idea; compute an Exponentially weighted Average of your gradients, and then use that gradient to update your weights.

$$V_{dw} = \beta V_{dw} + (1-\beta) dw$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W = W - \alpha V_{dw} ; b = b - \alpha V_{db}$$

↳ Batch Gradient Descent



↳ Momentum (less oscillation, straight line, faster)
(Balgabanna)

(↑ slower learning, ↔ faster learning } Bizim istediğimiz.)

Initialize; $V_{dw}=0$, $V_{db}=0$

On iteration t :

Compute dW , db on the current mini-batch.

$$V_{dw} = \beta V_{dw} + (1-\beta) dW$$

$$V_{db} = \beta V_{db} + (1-\beta) db$$

$$W = W - \alpha V_{dw}, \quad b = b - \alpha V_{db}$$

Hyperparameters: α , β

ör $\beta = 0.9 \Rightarrow$ average of last ≈ 10 gradients
(yaklaşık)

Bu β iyi çalışır genelde.

3) RMSprop (Root mean square prop)

On iteration t :

Compute dW , db on current Mini-batch.

$$S_{dw} = \beta S_{dw} + (1-\beta) dW^2 \quad \text{Element-wise}$$

$$S_{db} = \beta S_{db} + (1-\beta) db^2$$

$$W = W - \alpha \frac{dW}{\sqrt{S_{dw}}}, \quad b = b - \alpha \frac{db}{\sqrt{S_{db}}}$$

* Larger "Learning Rate" kullanabiliriz; Faster.

4) Adam Optimization Algorithm

Taking momentum and RMSprop, then put them together.

Initialize: $V_{dw}=0$, $S_{dw}=0$, $V_{db}=0$, $S_{db}=0$

On iteration t :

Compute dW , db using current mini-batch.

$$\text{Momentum} \rightarrow V_{dw} = \beta_1 V_{dw} + (1-\beta_1) dW, \quad V_{db} = \beta_1 V_{db} + (1-\beta_1) db \quad \text{"}\beta_1\text{"}$$

$$\text{RMSprop} \rightarrow S_{dw} = \beta_2 S_{dw} + (1-\beta_2) dW^2, \quad S_{db} = \beta_2 S_{db} + (1-\beta_2) db^2 \quad \text{"}\beta_2\text{"}$$

(Bias correction)

$$V_{dw}^{\text{corrected}} = V_{dw} / (1 - \beta_1^t), \quad V_{db}^{\text{corrected}} = V_{db} / (1 - \beta_1^t)$$

$$S_{dw}^{\text{corrected}} = S_{dw} / (1 - \beta_2^t), \quad S_{db}^{\text{corrected}} = S_{db} / (1 - \beta_2^t)$$

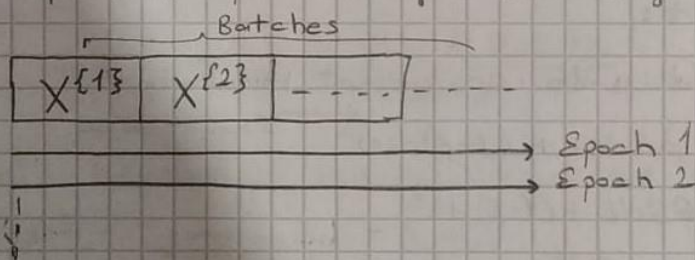
$$W := W - \alpha \frac{V_{dw}^{\text{corrected}}}{\sqrt{S_{dw}^{\text{corrected}} + \epsilon}}, \quad b := b - \alpha \frac{V_{db}^{\text{corrected}}}{\sqrt{S_{db}^{\text{corrected}} + \epsilon}}$$

(zero olursa diye)

Learning Rate Decay:

- * Help speed up your learning algorithm,
- * Slowly reduce your learning rate over time.

1 epoch \Rightarrow 1 pass through all training data



$$\alpha = \frac{1}{1 + \text{decay_rate} * \text{epoch_num}} \quad \alpha_0 \rightarrow (\text{initial } \alpha)$$

! Learning rate 'in her iterasyonlarda azalması' istenir, ki minima of cost function bulunabilsin.

Not: Our goal is to reach a global minimum.

But at times when our cost function has an irregular curve (mostly in deep networks), maybe we reach a local minimum, which is not as good as global minimum. One way to overcome this is by using the concept of momentum.

