

به نام خدا

تمرین دوم درس یادگیری عمیق

دکتر محمدی

مجتبی نافذ 96431335

سوال یک:

تفاوت $batch$, $mini-batch$, $stochastic$:

در حقیقت الگوریتم پایه ی بهینه سازی $lost function$ الگوریتم $Gradient Descent$ است که همان $batch gradient descent$ نامیده میشود

و $stochastic$, $min-batch$ به نوعی $variation$ هایی از $batch$ میباشد.

تفاوت این الگوریتم ها در مقدار دیتایی ایست که در هر $epoch$ میگیرند.

الگوریتم $batch$: هر $epoch$ روی کل دیتاست می باشد.

الگوریتم $mini-batch$: هر $epoch$ روی یک $batch$ و در واقع یک بخشی از دیتا است (معمولا 16 و 32 و 64 و ...)

الگوریتم $stochastic$: هر $epoch$ روی یک داده از دیتاست می باشد.

مقایسه SGD با GD :

الگوریتم $batch$:

رو دیتاست بزرگ میتواند خیلی کند باشد چون در هر $epoch$ فقط یک آپدیت دارد.

در دیتاست بزرگ دیتا قابلیت fit شدن در حافظه را ندارد. چون حجم دیتا بیشتر از حافظه است.

در سطوح $non-convex$ ممکن است در در مینیمم محلی گیر کند.

الگوریتم SGD:

در کل SGD شانس بیشتری از GD برای فرار مینیمم محلی را دارد.

و قطعا سریع تر به دقت های بالا خواهد رسید و همگرا خواهد شد.

به دلیل عملکرد $greedy$ خود تخمینی از گرادیان میزند.

به دلیل فرکانس نوسانات در داده ها ، نزدیک نقطه ی مینیمم $overshooting$ خواهد داشت.

مشکلات SGD :

به دلیل عملکرد $greedy$ خود تخمینی از گرادیان میزند.

به دلیل فرکانس نوسانات در داده ها ، نزدیک نقطه ی مینیمم $overshooting$ خواهد داشت.

در مناطقی که شیب کم است سرعت همگرایی به شدت کند خواهد بود و در شیب صفر حرکت نخواهد کرد. و انگونه بیشتر در

مینیمم محلی هم گیر هود کرد.

رفع مشکلات SGD با استفاده از Momentum :

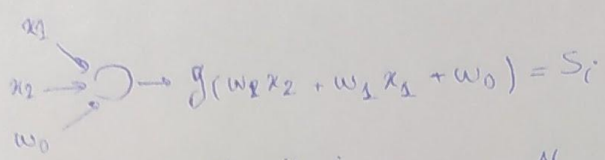
رفع مشکل سرعت همگرایی در SGD : در روش $momentum$ در واقع ما علاوه بر گرادیان در این نقطه به سرعت حرکت در

نقاط قبلی هم باید توجه کنیم یعنی رفته رفته با توجه به سرعت (گرادیان) مکان های قبلی شتاب بگیریم و $step size$ مان را بیشتر

کنیم و با گام های بلندتری حرکت کنیم. که این گونه ما در جاهایی که شیب کم است سرعت همگرایی بالاتری خواهیم داشت و گاهی

میتوانیم از مینیمم های محلی فرار کنیم.

سوال ۲
الف



use MSE

$$L = \frac{1}{N} \sum_{i=1}^N L^{(i)}(s^i, y^i) = \frac{1}{N} \left(\sum_{i=1}^N (y_{(i)} - s_{(i)})^2 \right) =$$

$$L = \frac{1}{N} \sum_{i=1}^N \left[y_{(i)} - \max(0, (w_2 x_2^{(i)} + w_1 x_1^{(i)} + w_0)) \right]^2 \Rightarrow$$

$$\frac{\partial L}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial w_1} \left[y_{(i)} - \max(0, (w_2 x_2^{(i)} + w_1 x_1^{(i)} + w_0)) \right]^2 =$$

$$\frac{\partial L}{\partial w_1} = \begin{cases} \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial w_1} [y_{(i)} - (w_2 x_2^{(i)} + w_1 x_1^{(i)} + w_0)]^2 & : w_2 x_2 + w_1 x_1 + w_0 > 0 \\ \frac{1}{N} \sum_{i=1}^N \frac{\partial}{\partial w_1} [y_{(i)}]^2 & : w_2 x_2 + w_1 x_1 + w_0 \leq 0 \end{cases}$$

$$\frac{\partial L}{\partial w_1} = \begin{cases} \frac{1}{N} \sum_{i=1}^N -2 [y_i - w_2 x_2^{(i)} - w_1 x_1^{(i)} - w_0] x_1 & : w_2 x_2 + w_1 x_1 + w_0 > 0 \\ 0 & : w_2 x_2 + w_1 x_1 + w_0 \leq 0 \end{cases}$$

$$w_1^{(t+1)} = w_1^{(t)} - \eta \frac{\partial L}{\partial w_1}$$

\$N \rightarrow\$ تعداد داده ها
batch

در اینجا \$w_0\$ و \$w_2\$ را هم در نظر بگیریم

$$\frac{\partial L}{\partial w_0} = \begin{cases} \frac{1}{N} \sum_{i=1}^N -2 [y_i - w_2 x_2^{(i)} - w_1 x_1^{(i)} - w_0] & : w \cdot x > 0 \\ 0 & : w \cdot x \leq 0 \end{cases}$$

$$\frac{\partial L}{\partial w_2} = \begin{cases} \frac{1}{N} \sum_{i=1}^N -2 [y_i - w_2 x_2^{(i)} - w_1 x_1^{(i)} - w_0] x_2 & : w \cdot x > 0 \\ 0 & : w \cdot x \leq 0 \end{cases}$$

در اینجا \$w_0\$ و \$w_2\$ را هم در نظر بگیریم
و \$w_1\$ را هم در نظر بگیریم

$$w_i^{(t+1)} = w_i^{(t)} + \eta \frac{\partial L}{\partial w_i}$$

sigmoid



$$L = \frac{1}{N} \sum_{i=1}^N \left[y_i - \frac{1}{1 + e^{-z}} \right]^2 \rightarrow$$

$$z = w_2 x_2 + w_1 x_1 + w_0$$

$$\frac{\partial L}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N -2 [y_i - \sigma(z_i)] \sigma(z_i) (1 - \sigma(z_i)) x_1$$

مغزی درون Z

$$\frac{\partial L}{\partial w_1} = \frac{1}{N} \sum_{i=1}^N [G(z_i) - y_i] G(z_i) (1 - G(z_i)) x_1^{(i)}$$

$$z_i = w_2 x_2^{(i)} + w_1 x_1^{(i)} + w_0$$

دفع شود $\frac{\partial L}{\partial w_2}$

$$\frac{\partial L}{\partial w_2} = \frac{1}{N} \sum_{i=1}^N [G(z_i) - y_i] G(z_i) (1 - G(z_i)) x_2^{(i)}$$

$$\frac{\partial L}{\partial w_0} = \frac{1}{N} \sum_{i=1}^N [G(z_i) - y_i] G(z_i) (1 - G(z_i))$$

و محاسبه

← در مورد آپدیت هر

$$w_{i(t+1)} = w_{i(t)} + \eta \frac{\partial L}{\partial w_i}$$

سوال ۳:

نتیجه این که همگرایی momentum از GD بیشتر است.