

# التعلم العميق Deep Learning

## ○ الدرس الأول

### التعلم العميق و الشبكات العصبية

- الأسبوع الأول : مقدمة للتعلم العميق
- الأسبوع الثاني : أساسيات الشبكات العصبية
- الأسبوع الثالث : الشبكات العصبية المجوفة
- الأسبوع الرابع : الشبكات العصبية العميقة

## ○ الدرس الثاني

### تطوير الشبكات العميقة : المعاملات العليا

- الأسبوع الأول : السمات العملية للتعلم العميق
- الأسبوع الثاني : الحصول علي القيم المثالية
- الأسبوع الثالث : ضبط قيم الشبكات العميقة

## ○ الدرس الثالث

### هيكلية مشاريع الـ ML

- الأسبوع الأول : استراتيجيات الـ ML - 1
- الأسبوع الثاني : استراتيجيات الـ ML - 2

## ○ الدرس الرابع

### الشبكات العصبية الملتفة CNN

- الأسبوع الأول : أساسيات الشبكات العصبية الملتفة
- الأسبوع الثاني : حالات عملية من الشبكات العصبية الملتفة
- الأسبوع الثالث : التعرف علي الأشياء
- الأسبوع الرابع : التعرف علي الوجه

## ○ الدرس الخامس

### الشبكات العصبية المتكررة RNN

- الأسبوع الأول : مفهوم الشبكات العصبية المتكررة
- الأسبوع الثاني : المعالجة اللغوية الطبيعية NLP
- الأسبوع الثالث : نماذج التتابع

## الأسبوع الثاني : الحصول علي القيم المثالية Optimization

## الأسبوع الثاني : الحصول علي القيم المثالية Optimization

\_\_\_\_\_

- عقب الانتهاء من هذا الكورس , ستكون قادر علي :
  - التعرف علي عدد من طرق الحصول علي القيم المثلي زي الطريقة العشوائية , او الزخم , او RMS
  - استخدام اسلوب الكميات الكبيرة العشوائي , لتسريع عملية الـ optimization
  - التعرف علي مميزات تضائل معامل التعلم , واستخدامه في الـ optimization

\* \* \* \* \*

هذا الاسبوع نتناول بالتحديد الاساليب والتقنيات الخاصة بما يسمى الـ Optimization و هي الطرق المستخدمة للحصول علي الارقام المناسبة للخوارزم

و من عيوب التعلم العميق و الشبكات العميقة , أنها لا تعمل بفعالية مع البيانات الضخمة Big Data وذلك لان الشبكات العميقة اساسا تستغرق وقتا كبيرا , وإذا ما أضيف عليها مشكلة التعامل مع كم كبير من البيانات , فيعني هذا وقتا مهولا

لذا كان من الضروري الوصول لطرق الحصول علي القيم المثالية في وقت قصير , وهو نقطة بحثنا هنا

وستتناول الان تكنيك : الانحدار الاشتقاقي للمجاميع الصغيرة mini-batch gradient descent

ولفهمها , علينا ان نتناول مثالا مثل هذا :

لو قلنا ان لدينا عدد من عناصر العينة m , وكل عنصر لديه عدد من الـ features بقيمة n , فتكون لدينا مصفوفة X كإيتال هكذا

$$X = \begin{bmatrix} x^{(1)} & x^{(1)} & x^{(1)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(1005)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

حيث أن كل x من هؤلاء هي فيكتور من عمود واحد و الصفوف هي الـ features

إذن يكون الابعاد النهائية للمصفوفة X هي (n x m)

و بالتالي تكون مصفوفة y بالمثل هي :



$$\underline{X} = \begin{bmatrix} x^{(1)} & x^{(1)} & x^{(3)} & \dots & x^{(1000)} & | & x^{(1001)} & \dots & x^{(2000)} & | & \dots & | & \dots & x^{(m)} \end{bmatrix}$$

$(n_x, m)$        $X^{\{1\}} (n_x, 1000)$        $X^{\{2\}} (n_x, 1000)$        $X^{\{5,000\}} (n_x, 1000)$

الجزء الأول من رقم 1 إلى 1000 , الثاني من 1001 إلى 2000 , وهكذا , وهي ما تسمى المجاميع الصغيرة mini-batches

ونقوم بتسمية الأجزاء بطريقة  $X^{\{1\}}$  ,  $X^{\{2\}}$  ,  $X^{\{3\}}$  (الرقم لأعلي) بحيث يشير الرقم لرقم الجزء

لا تنس ان التسمية العليا superscript تختلف حسب نوع القوس . .

فحينما نقوم بعمل تسمية عليا بقوس دائري  $x^{(5)}$  , فنحن نعني العنصر الخامس من العينة

وحينما نقوم بعمل تسمية عليا بقوس مربع  $z[3]$  , فنحن نعني الطبقة الثالثة من الشبكة العصبية

بينما حينما نشير بالقوس المتعرج  $x^{\{15\}}$  فنحن نعني الجزء الخامس عشر من المجاميع المقسمة

و بالتالي فإن أبعاد مصفوفات المجاميع  $x$  ستكون  $(n \times 1000)$  حيث أن صفوفها هي نفس صفوف  $x$  الاصلية و هي عدد الـ features بينما اعمدتها هي عدد 1000 عنصر تم اجتزائه

ونقوم بعمل نفس التجزيئ بنفس العدد المستقطع المستخدم في  $x$  , نقوم به في  $y$  , حتي يقوم الخوارزم بعمل حسابه الدقيق دون اضطراب , وبالتالي ستكون  $y$  :

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \dots & y^{(1000)} \\ y^{(1001)} & \dots & y^{(1500)} \end{bmatrix} \dots \begin{bmatrix} \dots & y^{(m)} \end{bmatrix}$$

$\underbrace{\hspace{10em}}_{Y^{\{1:3\}} \quad (1, 1000)} \quad \underbrace{\hspace{10em}}_{Y^{\{2:3\}} \quad (1, 1000)} \quad \dots \quad \underbrace{\hspace{10em}}_{Y^{\{5,000\}} \quad (1, 1000)}$

علي أن تكون ابعاد كل منها (1 x 1000)

وبالتالي ستقوم الشبكة العصبية بعمل المسار الامامي والخلفي لكل جزء علي حدة , ومن ثم عمل تعديل لقيم  $w$  ,  $b$  ثم تناول الجزء التالي وهكذا , و عقب الانتهاء من العينة بالكامل (5 الاف جزء) يتم تكرار العملية بالكامل عددا من المرات , واقوم بفحص قيمة  $J$  للوصول الي القيمة المناسبة

نبدأ بكتابة الكود الخاص لها

using  $\frac{X^{t+1}}{(\text{as func } 1000)}$

$X, Y$

repeat  $\sum$  for  $t = 1, \dots, 5000$  {

Forward prop on  $X^{\text{test}}$

$$\begin{aligned} z^{(1)} &= W^{(1)} X^{\text{test}} + b^{(1)} \\ A^{(1)} &= \sigma^{(1)}(z^{(1)}) \\ &\vdots \\ A^{(L)} &= \sigma^{(L)}(z^{(L)}) \end{aligned} \quad \left. \begin{array}{l} \\ \\ \end{array} \right\} \text{Vectorized implementation (1000 examples)}$$

for  $X^{(t)}, Y^{(t)}$

Compute cost  $J^{\text{test}} = \frac{1}{1000} \sum_{i=1}^n \ell(y^{(i)}, \hat{y}^{(i)}) + \frac{\lambda}{2 \cdot 1000} \sum_a \|W^{(a)}\|_F^2$

Backprop to compute gradients wrt  $J^{\text{test}}$  (using  $(X^{(t)}, Y^{(t)})$ )

$$W^{(1)} = W^{(1)} - \text{add}(W^{(1)}, \text{d}W^{(1)}), \quad b^{(1)} = b^{(1)} - \text{add}(b^{(1)}, \text{db}^{(1)})$$

"1 epoch"  
pass through training set.

نري أولاً أن هناك for تبدأ من 1 إلى 5000 , وهي التي سنتناول الاجزاء من الجزء الاول للاخير

بعدها اقوم بحساب المسار الامامي لـ  $X\{t\}$  , لا تنس أن  $t$  هو العنصر الذي يحمل الرقم المتغير من 1 إلي 5000

ف نقوم بحساب قيم  $Z$  ثم  $A$  لكل طبقة بالترتيب حتي أصل لقيمة  $y^A$  , ثم نقوم بحساب دالة الخطأ cost function و لا انسي قيمة التنعيم اذا ما كنت ستستخدمها

بعدها اقوم بالمسار الخلفي للوصول لقيم  $dw$  ,  $db$  , ثم اقوم بالخطوة الأهم و هي تعديل قيم  $w$  ,  $b$

كل هذا يسمى 1 epoch اي خطوة واحدة , فيتم الرجوع لدالة for ليتناول الجزء الثاني  $X_{2}$  و يتعامل معه بالقيم المعدلة لـ  $w, b$  , ثم الجزء الثالث و هكذا

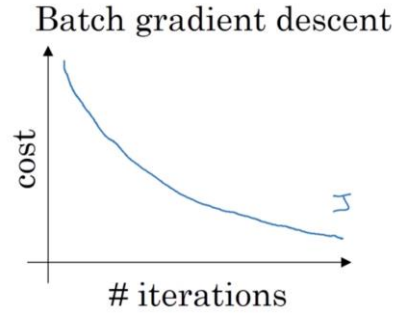
بعدما ينتهي من العينة بالكامل , نقوم بتكرار العملية من بدايتها مرات عديدة , حتي نصل للقيمة المثلي للأوزان  $w, b$

\* \* \* \* \*

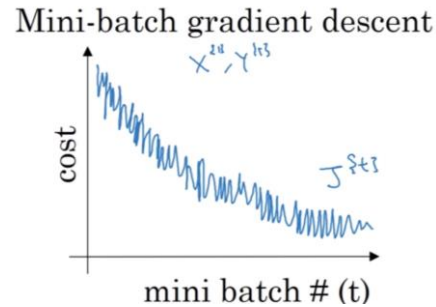
دعنا نتناول عدد من خصائص الـ mini-batch gradient descent

أولا إذا تكلمنا عن معادلة الخطأ cost function و الذي نرمز لها بـ  $J$

ففي الأسلوب التقليدي Batch Gradient Descent نري أن قيمة  $J$  دائما ما تنخفض مع كل محاولة جديدة iteration هكذا



بينما في تكنيك المجاميع الصغيرة , فقد تري ان الشكل يختلف قليلا :





والسبب في هذا التذبذب , أنه في خلال نفس المحاولة قد يصادف الخوارزم جزء  $X\{1\}$  مثلا يكون قيمة الخطا فيه قليلة , بينما بعدها قد يصادف  $X\{6\}$  تكون قيمة الخطأ كبيرة , لأنها عينة كبيرة

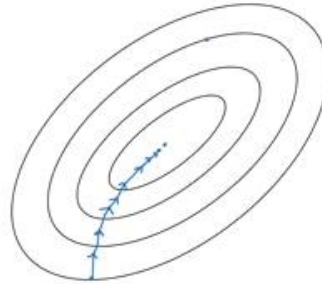
لكن ستلاحظ أن الاتجاه العام للمنحني يكون لاسفل , كما أنه اسرع كثيرا

نصل لسؤال مهم , علي أي أساس اقوم بتحديد عدد عناصر الجزء الواحد , هنا قلناه أنه 1000 , لكن هل احدد رقم اكبر او اصغر ؟ ؟

لنجيب علي السؤال , علينا التعرف علي القيمتين المتطرفتين في هذا الأمر :

نتعامل بداية مع القيمة القصوي , وهي أن نفرض أن عدد عناصر الجزء هو نفسه العدد الكلي للعناصر وهو  $m$  , وهو ما يجعل عدد الأجزاء هو واحد فقط , وبهذا سنصل لمفهوم Batch Gradient Descent مرة أخرى

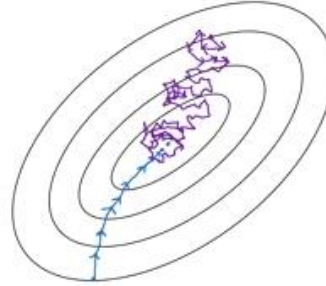
وإذا اردنا تتبع مسار قيمة  $J$  هنا سيكون كالتالي :



نري أنه مسار مباشر دون التفاف , لكن يعيبه الوقت الطويل اللازم لتنفيذه , فكل خطوة فيه يستلزم المرور بالملايين الخمسة للعناصر

وإذا تناولنا الحد الأدنى لعدد عناصر الجزء , سيكون 1 , لانه يستحيل ان يتم رقم اقل من هذا , وهذا معناه أنه سيكون لدينا خمسة ملايين جزء , كل جزء من عنصر واحد , وعقب معالجة كل جزء فيهم , يقوم الخوارزم بعمل update لقيم  $w$  ,  $b$

و هذا النظام يسمى Stochastic gradient descent اي الانحدار الاشتقاقي العشوائي , وستري أن خطوات تعديل قيمة الاوزان و بالتالي تغيير معادلة الخطأ  $J$  سريعة جدا , لكن مشكلتها الأساسية انها ستكون متذبذبة المسار , حيث أن كل عنصر من عناصر العينة يأخذها في اتجاه , وهو ما يجعل وصولها للقيم المثلي شديد الصعوبة و الوقت الطويل , وقد لا يصل في النهاية



المنطقة الوسطي بينهما هي ما نتناوله , وهو المجاميع الصغيرة , حيث نستفيد من ميزة كلا منهما , التعديل السريع , والاتجاه شبه المباشر



نتكلم الان عن عدد من الخوارزميات الأخرى , والتي قد تكون أفضل في معالجة الشبكات العصبية من الانحدار الاشتقاقي . .

منها ما يسمى (Exponentially weighted averages المتوسطات الأسية)

ولنفهم ماهيتها , علينا أن نتناول عددا من الأرقام مثل هذه :

لو أخذنا درجات الحرارة في لندن علي مدي أيام السنة , علي اعتبار ان نيتا 1 تمثل يوم 1 يناير , ونيتا 180 تمثل يوم 1 يونيو , ونيتا 365 تمثل 31 ديسمبر , فستكون الدرجات كالتالي :

$$\theta_1 = 40^{\circ}\text{F}$$

$$\theta_2 = 49^{\circ}\text{F}$$

$$\theta_3 = 45^{\circ}\text{F}$$

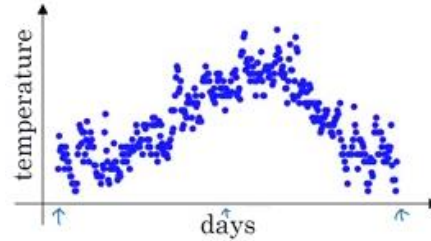
$\vdots$

$$\theta_{180} = 60^{\circ}\text{F}$$

$$\theta_{181} = 56^{\circ}\text{F}$$

$\vdots$

إذا ما اردنا أن نرسمها علي جراف فستكون كالتالي :



ومشكلة هذا الرسم انه شديد التذبذب , وذلك بسبب كمية كبيرة من المعلومات لا احتاجها , فانا اريد رسم أكثر نعومة و إيضاحا من هذا , فسنقوم بعمل بحساب المتوسط الأسّي

وهذا عبر تحديد معامل معين بيتا  $\beta$  , الذي سنقوم بضربه في قيمة اليوم السابق (ليس درجة الحرارة لكن قيمة  $V$  ) , ونجمعه علي قيمة 1 ناقص بيتا , مضروبة في درجة حرارة اليوم نفسه , ليكون اسمه  $V_t$  , كالتالي :

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

وهو ما سيجعل قيمة  $V_t$  لا تمثل فقط درجة اليوم , لكنها مرتبطة بقوة بالدرجات السابقة لها , وبالتالي فأني ارتفاع او انخفاض في درجات الحرارة في اليوم , لا تجعل القيمة تتغير بشكل كبير , لكن بدرجة قليلة

وإذا اخترنا لبيتا قيمة 0.9 , تكون المعادلات كالتالي :

$$\begin{aligned} V_0 &= 0 \\ V_1 &= 0.9 V_0 + 0.1 \theta_1 \\ V_2 &= 0.9 V_1 + 0.1 \theta_2 \\ V_3 &= 0.9 V_2 + 0.1 \theta_3 \\ &\vdots \\ V_t &= 0.9 V_{t-1} + 0.1 \theta_t \end{aligned}$$

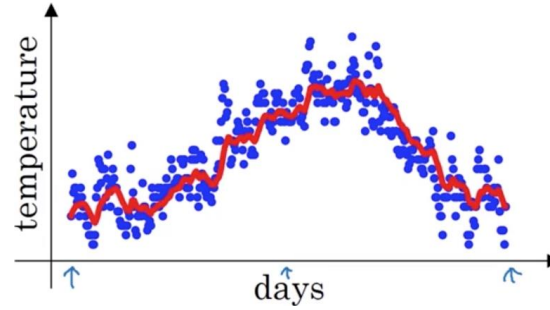
وبالتالي :

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

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

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

ويكون الرسم كالخط الأحمر هنا :



فالخط هنا أكثر هدوء و استقرار , و اقل تذبذب , ويجعل التغيرات ابطئ قليلا

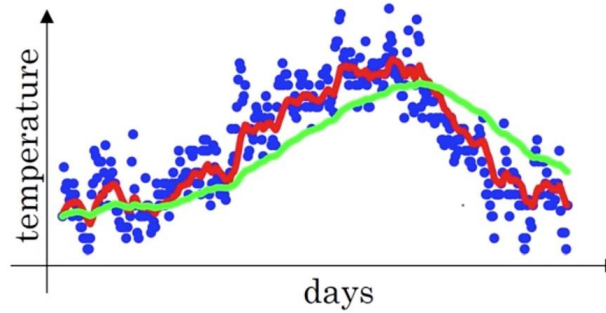
وكلما زاد مقدار بيتا , زاد البطئ في التغير حتي يتحول الرسم لما يشبه المعادلات الخطية , وكلما قلت قيمة بيتا , كلما زاد التذبذب

وذلك لأن قيمة بيتا يتم ضربها في قيمة الايام السابقة , فكلما زادت قيمة بيتا , قلما قل تأثير قيمة اليوم و زاد ارتباطها بالايام السابقة

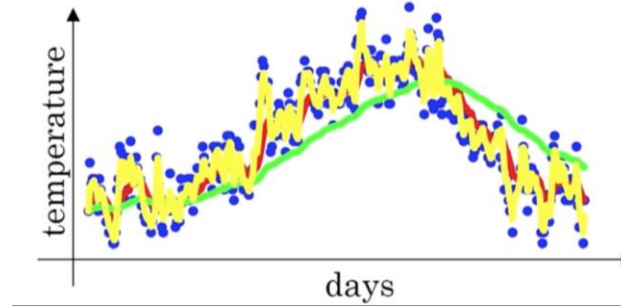
و يقال أن الارقام الحالية لا تمثل العلاقة بين درجات الحرارة لكل يوم , لكن لمتوسط عدد من الايام , يرتبط بقيمة بيتا , بالمعادلة

$$1 / ( 1 - \beta )$$

وهذا معناه أن , في حالة تحديد بيتا برقم معقول , مثلا 0.8 , فيكون المقام 0.2 , بالقسمة عليه يكون 5 , فوقتها يكون الرسم يمثل متوسط لكل خمسة ايام  
وإذا زاد الرقم ووصل مثلا لـ 0.98 , فالمقام سيكون 0.02 , وبالقسمة عليه سيكون خمسين , وهذا معناه متوسط خمسين يوما , وهو ما سيتم رسمه بالخط  
الأخضر :



اما اذا ما اخترنا رقم صغير نسبيا لبيتا , مثلا نصف , وهو ما يجعل المقام بـ 0.5 , و يكون المتوسط يومين , أي سيكون الرسم لمتوسط يومين فقط , وهو ما  
يجعل قيمة المتوسطات غير مرتبطة كثيرا بالقيم السابقة , و سيكون قيمة تنذبها كبيرا , مثل الخط الاصفر :



وإذا قمنا بمزيد من البحث في نفس النظرية . .

فالمعادلة المستخدمة هي :

$$v_t = \beta v_{t-1} + (1 - \beta)\theta_t$$

وبالتالي

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

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

$$v_{98} = 0.9v_{97} + 0.1\theta_{98}$$

وإذا قمنا باستبدال  $v_{99}$  بمعادلتها , ثم قيمة  $v_{98}$  بمعادلتها الداخلية و هكذا , تصير المعادلة :

$$V_{100} = 0.1\theta_{100} + 0.1 \cdot 0.9 \cdot \theta_{99} + 0.1 \cdot 0.9^2 \cdot \theta_{98} + 0.1 \cdot 0.9^3 \cdot \theta_{97} + \dots$$

أي أن قيمة أي  $v$  لأي يوم , هو قيمة  $(\beta-1)$  في ثبنا لنفس اليوم , مجموعة علي  $(\beta-1)$  في  $(\beta)$  في ثبنا اليوم السابق , مجموعة علي  $(\beta-1)$  في  $(\beta)$  تربيع , في اليوم السابق , مجموعة علي  $(\beta-1)$  في  $(\beta)$  تكعيب في اليوم السابق وهكذا

ولاحظ ان الصيغة الرياضية الخاصة ببينا هنا :

$$(1-x)^{-1/x}$$



علي اعتبار ان  $x$  هي  $(\beta-1)$  , الصيغة ديه مع كبر ارقام البيت , تؤدي للرقم  $e$

وبالتالي المعادلة العامة هتكون :

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

يعني الخوارزم الخاص بيها :

$$V_\theta := 0$$

$$V_\theta := \beta v + (1-\beta) \theta_1$$

$$V_\theta := \beta v + (1-\beta) \theta_2$$

ولو هنتعامل ككود ,

$\rightarrow V_0 = 0$   
 Repeat {  
     Get next  $O_c$   
      $V_0 := \beta V_0 + (1-\beta) O_c \leftarrow$   
 }

واللي نكون فيها قيمة لـ  $v$  و نقوم بعمل update لها بضرب بيتا في  $v$  الاصلية , مجموعة علي 1 ناقص بيتا , مضروبة في ثيتا لدرجة اليوم نفسه و ميزة هذا الكود انه سريع و سهل التعامل مع البروسييسور , ولا يستهلك ذاكرة كثيرة

\* \* \* \* \*

ولدي تطبيق هذا الخوارزم , ستلحظ سلوكا غريبا , خاصة في حسابات الايام الاولى من درجات الحرارة . .

فنحن نفترض دائما  $v_0$  علي انها صفر . .

وبما أننا نتعامل مع هذه المعادلة :

$$v_0 = 0$$

$$v_1 = \beta v_0 + (1 - \beta) \theta_1$$

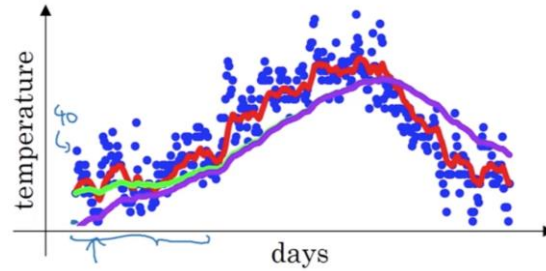
$$v_2 = \beta v_1 + (1 - \beta) \theta_2$$

$$v_3 = \beta v_2 + (1 - \beta) \theta_3$$

فتجد أنه لو كانت درجة حرارة اول يوم مثلا تساوي 40 , فستكون قيمة  $v_1$  والتي يفترض ان تمثل متوسط اليوم الأول , ستكون بيتا في  $v_0$  والتي تساوي صفر , مضروبة في 0.1 في 40 , والتي ستساوي فقط 4 درجات

في اليوم التالي إذا كانت درجة الحرارة 45 مثلا , ستكون 0.9 في 4 , زائد 0.1 في 45 , اي 8 درجات

وهكذا ستجد أنه في اول ايام فإن درجات الحرارة المحسوبة بعيدة جدا عن الدرجات الحقيقية , مما يرسم هذا الشكل

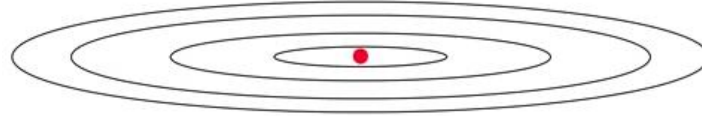




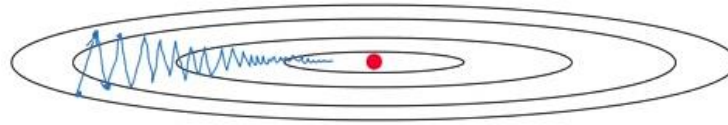
ولاستخدام هذا التكنيك في الشبكات العميقة , سنتعرف علي مفهوم هام يسمى الانحدار الاشتقاقي بالزخم Gradient descent with momentum

وهو يعتبر تكنيك سريع و قوي للوصول للقيم المثالية , مقارنة بالانحدار الاشتقاقي التقليدي

والفكرة تبدأ من هذا الشكل :



إذا كان لدينا القيمة المثلي هي النقطة الحمراء , و سنقوم بالبدا من أحد النقاط , فيمكن أن يكون المسار التقليدي للانحدار الاشتقاقي هو :



ومشكلة هذا المسار ان خطواته طويلة و كثيرة , و أن التحرك الرأسى يعمل ازعاج و نوع من تضییع الوقت و الجهد , لان المسار الرأسى لا يقوم بحل اي مشكلة او التقدم للقيم المثلي , فكأننا نريد الحفاظ علي التقدم الأفقي و ليس الرأسى

وهنا نقوم باستخدام فكرة الـ momentum والذي تم استقاؤه من فكرة المتوسطات الأسية

فإذا قمنا بإنشاء قيمة تسمى  $V_{dw}$  كذلك قيمة تسمى  $V_{db}$  بحيث نطبق عليها مبدئ المتوسطات الأسية . .

بحيث يكون :

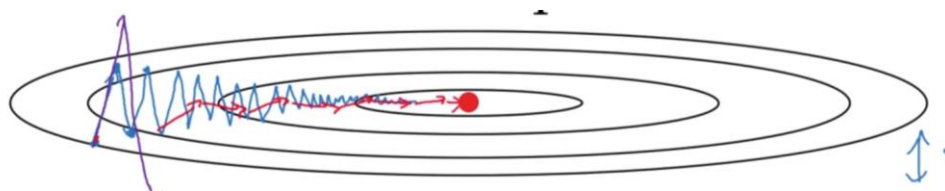
$$v_{dW} = \beta v_{dW} + (1 - \beta) dW$$

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

وبالتالي نعمل update لقيمة  $b$  ,  $w$  بدل من استخدام  $dw$  ,  $db$  نستخدم :

$$W = W - \alpha v_{dW}, b = b - \alpha v_{db}$$

و لأن تكنيك المتوسطات الأسية يجعل الخطوة اهدى و اقل , ويجعل التغييرات بتمشي بشكل ابطئ شوية , فهذا هو السبب الرئيسي , الذي يجعل خطوات الانحدار ابطئ رأسيا , لكنها تتجه افقيا بشكل جيد , كما في الخط الأحمر



وبالتالي ببساطة , يتم استخدام القوانين

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

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

وأحيانا ويتم ازالة قيمة  $(\beta-1)$  , بحيث تكون المعادلة :

$$v_{dW} = \beta v_{dW} + dW$$

و ايضا هذه الصيغة فعالة , لكن الأولي أفضل

وغالبا ما تكون قيمة بيتا بـ 0.9 , و غالبا ايضا لا يتم استخدام فكرة معامل الخطأ  $\frac{V_t}{1 - \beta^t}$ . وذلك لأن القيم سرعان ما يتم تصحيحها سريعا

\* \* \* \* \*

أحد الطرق الأخرى المستخدمة , عوضا عن الانحدار الاشتقاقي , او ذو الزخم , هو ما يسمى الانحدار الاشتقاقي باستخدام الجذر التربيعي RMS Prop

والـ RMS اختصار root mean square

وتقوم فكرتها علي تعديل المعادلات , بإنشاء رمز جديد يسمى Sdw و Sdb والذان هما :

$$S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db$$

بحيث بيتا 2 , هي كمان قيمة لبيتا , لكن تم عمل 2 لمنع التداخل مع بيتا الأولي

ويتم عمل update لقيم  $w, b$  هكذا :

$$W := W - \alpha \frac{V_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}} + \epsilon}}$$

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

مع ملاحظة أنه تم وضع قيمة ايسلون مضافة إلي  $S$  في المعادلتين , لتجنب القسمة علي صفر , في حالة كانت قيمة  $S$  شديدة الضالة

\* \* \* \* \*



تحدثنا في الفقرتين السابقتين , عن انحدار الزخم , و انحدار الجذر التربيعي , ونقوم الان بتناول اسلوب (آدم للقيم المثلي) Adam optimization algorithm ) والذي يقوم بدمج الاسلوبين معا

هنتعامل بكلا من الـ  $V$  ,  $S$  لكلا من  $w$  ,  $b$  , وهيكون قيمهم الأولي بصفر , كدة :

Initialize  $V_{dW} = 0, S_{dW} = 0, V_{db} = 0, S_{db} = 0$ .

بعدها يتم استخدام نفس المعادلات السالف شرحها

$$V_{dW} = \beta_1 V_{dW} + (1 - \beta_1) dW$$

$$V_{db} = \beta_1 V_{db} + (1 - \beta_1) db$$

$$S_{dW} = \beta_2 S_{dW} + (1 - \beta_2) dW^2$$

$$S_{db} = \beta_2 S_{db} + (1 - \beta_2) db$$

بعدها يتم حساب ما يسمى  $V$  و  $S$  قيمة correct

$$V_{dW}^{\text{corrected}} = V_{dW} / (1 - \beta_1^t)$$

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

$$S_{dW}^{\text{corrected}} = S_{dW} / (1 - \beta_2^t)$$

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

ومنها يتم عمل update لقيم w , b

$$W := W - \alpha \frac{V_{dW}^{\text{corrected}}}{\sqrt{S_{dW}^{\text{corrected}} + \epsilon}}$$

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

وستري أن هذه الصيغة , فعالة بقوة لعدد كبير من الخوارزميات المرتبطة بالشبكات العصبية

و نري هنا عددا من الـ hyper parameters والواجب ضبطها جيدا قبل العمل

و غالبا ما نقوم بضبطها بهذه القيم :

$\alpha$  : needs to be tuned

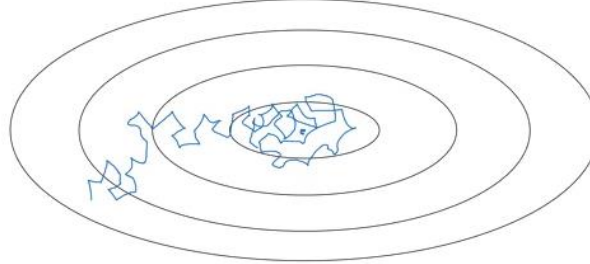
$\beta_1 : 0.9 \leftarrow (dW)$

$\beta_2 : 0.999 \leftarrow (dW^2)$

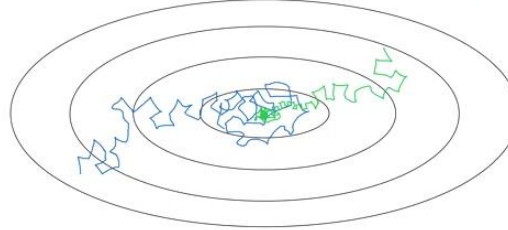
$\epsilon : 10^{-8}$

\*\_\*\*

ننتقل لنقطة أخرى في ذات الموضوع , وهي المتعلقة بمعامل التدريب **learning rate** المسمي ألفا  
لو تم تحديد ألفا برقم كبير فستقفز عملية الـ **optimization** قفزات واسعة و لن يصل للحل أبدا  
وإذا اخترنا قيمة متوسطة , فسنري أنه لن يصل بسهولة , وسيظل الخوارزم يدور حول نفسه حينما يقترب



لكن إذا قمنا بعمل تقليل تدريجي لقيمة ألفا , سنري أن الخوارزم يصل الي ارقام قريبة جدا من المطلوب , مثل الخط الأخضر:



فعلينا أن نقوم بعمل تقليل تدريجي لقيمة ألفا , ويكون التقليل مرتبط بعدد المحاولات , ومن الصيغ المستخدمة لهذا هي :

$$\alpha = \frac{1}{1 + \text{decay-rate} * \text{epoch-num}} * \alpha_0$$

حيث الفا 0 , هي القيمة الاولى لالفا , و الـ decay rate رقم نقوم بتحديدده و ليكون 1 , و الـ epoch number هو عدد المحاولات و بالتالي كلما زادت المحاولات , زاد المقام , فقلت الفا النهائية , فمثلا لو قلنا أن الفا 0 تساوي 0.2 , فيكون الجدول كالتالي :

Epoch	Alpha
0	0.2
1	0.1
2	0.06
3	0.05
4	0.04
5	0.03

ويمكن تغيير هذه الارقام عبر التحكم في و الـ decay rate

كما أن هناك صيغ أخرى لتقليل الفا , مثل الصيغة الأسية :

$$\alpha = 0.95^{\text{epoch-num}} \cdot \alpha_0$$

ففي حين تزيد عدد المحاولات , يزداد أس الرق المكتوب و ليكن 0.95 , فيقل قيمته , فتقل الفا

### أو الصيغة الجذرية

$$\alpha = \frac{k}{\sqrt{\rho \pi h \omega}} \cdot \alpha_0$$

### أو الصيغة السلمية



والتي فيها نحدد رقم معين لالفا لرينج عدد محاولات و يقل بالتتابع

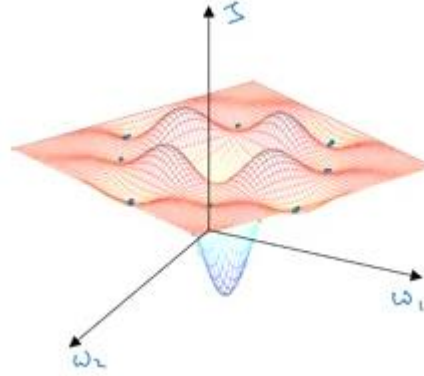
أو الطريقة اليدوية , وهي التي نغير فيها الفا بأنفسنا , وهي تصلح غالبا حينما يكون عدد البيانات هائل , فتستغرق المرة الواحدة عدد من الساعات , ففي خلال هذه الساعات و مع مراقبة المعالجة , اتخذ قرار بتقليل او زيادة الفا

\* \* \* \* \*

أخيرا , سنتكلم عن مشكلة لن تحدث , ومشكلة أخرى قد تحدث

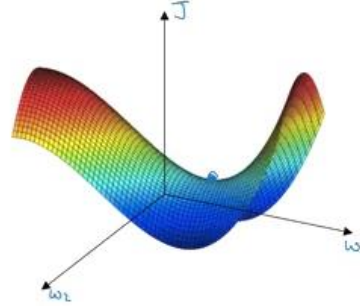
المشكلة التي لن تحدث هي الوقوع في فخ local optima

والمقصود بها , هو القيم الدنيا المحلية , في حين هناك قيم دنيا عظمي مثل هذا :



فكما هو واضح , هناك عددا من النقاط الزرقاء وهي دنيا محلية , بينما القاع به الدنيا العظمي , وهذه المشكلة تحدث فقط في المسائل البسيطة ذات البعدين , وهي شبه مستحيل حدوثها مع الـ NN

والواقع ان التصميمات الخاصة بالشبكات العميقة , تكون اقرب لهذا . .



والتي تحمل عددا كبيرا من الأبعاد , قد تصل لعشرات الآلاف , وبالتالي لن يتواجد قيمة محلية , ولكن ما يسمى الـ saddle point او نقطة السرج فكما نري في الشكل , لا يوجد امكانية لوجود قيمة محلية في الاساس , والاشتقاقات دائما توصلنا للقيمة الدنيا العظمي . .

و هي تسمى نقطة السرج لان الشكل نفسه يشبه سرج الحصان

لكن المشكلة التي قد تحدث , هي ما يسمى مشكلة التسطح plateau problem

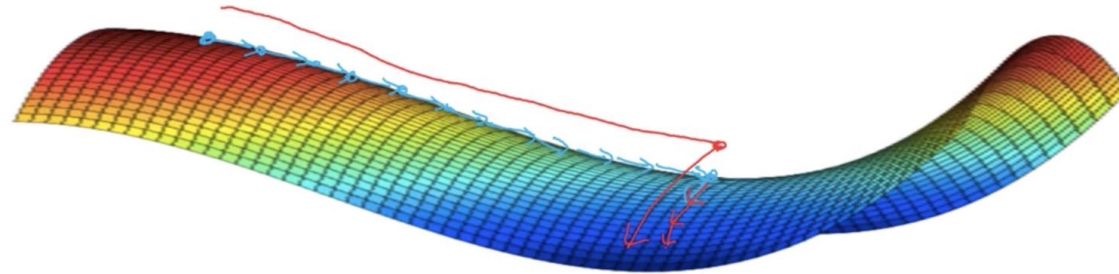
والتي يكون خلالها السطح اساسا شبه افقي بالكامل , مما يعني ان الاشتقاقات اقرب للصفر , وبالتالي معدل التغيير يكون بطيء للغاية

فالانتقال من النقطة اليسري لليمني سيستغرق وقتا رهيبا , لتسطح المستوي

وهذه المشكلة قد تحل بالتكنيكات السالف ذكرها في ايجاد الانحدارات الاشتقاقية



## Problem of plateaus



\* \* \* \* \*

## نهاية الأسبوع الثاني