

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 على 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 على الأشياء
- الأسبوع الرابع : التعرف على الوجه

○ الدرس الخامس : الشبكات العصبية المتكررة RNN

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

درس 5: الشبكات العصبية المتكررة RNN

الأسبوع الأول : مفهوم الشبكات المتكررة



نتعرف الان في الكورس الخامس , علي ما يسمى : نموذج التتابع sequence model , ومنه الشبكات المتكررة RNN , وهو المستخدم بكثير في العديد من تطبيقات الـ DL مثل التعرف علي الصوت و غيرها

و يسمى هذا النوع : نموذج التتابع , لأن البيانات التي يتم التعامل معها سواء كـ inputs or outputs هي دالة في الزمن , مثل الملفات الصوتية , او الكلام مكتوب .

و من أمثلة التطبيقات المعتمدة علي نموذج التتابع :

- تحويل الصوت الي كلام :
 - يتلقي ملف صوتي و يترجمه إلي text
- إنتاج الموسيقى :
 - وهي الخوارزميات التي تقوم بانتاج قطع موسيقية كاملة بالشبكات العصبية

- استنتاج التقييم : sentiment classification :
 - وهي إمكانية استنتاج تقييم معين للمنتج بناء علي التعليقات , فلو كتب شخص في تعليق كلمة "زفت" فيتم استنتاج تقييم اوتوماتيك لها
 - تحليل الـ DNA
 - وهي تكون عبر تحليل الرموز المعقدة لشفرة الـ DNA الي مصطلحات بسيطة
 - الترجمة الآلية
 - أي الترجمة من لغة إلي لغة عبر خوارزم محدد
 - تحليل فيديو
 - لمتابعة كل اطار من الفيديو و تحليله و التعرف علي الاشخاص به
 - التعرف علي الاسماء
 - مثل ان يتم تحليل نصوص كبيرة , وتحديد اي اسماء اشخاص فيه و الاشارة لها
- فتطبيقات نموذج التسلسل sequence model كثيرة , وتختلف في شكلها و نوعها , فممكن ان تجد ان المدخل هو قيمة تتابعية و المخرج ليس كذلك او العكس , كما ان طول المدخل و المخرج قد يتساوي و قد يختلف

و للتعرف علي مفهوم : نموذج التتابع , علينا ان نشرح بالتفصيل أحد تطبيقاته , وليكن التطبيق الاخير المذكور في الفيديو السابق , وهو المسمي : مدي تواجد الاسماء named-entity , وهو الذي يقوم بالبحث في نصوص عديدة , عن اسم اي شخص و الاشارة اليه , لاحظ ان هذا الأمر ليس سهل فهو لا يبحث عن ام شخص معين بالية search عادية (مثلا يبحث عن اسم دونالد ترامب بالتحديد) , ولكن يبحث عن اي اسم حتي لو لم يعرفه الخوارزم , ليكون الخوارزم جاهز لاي اسماء .

و غالبا ما يستخدم في محركات البحث , حتي يقوم اوتوماتيك بتحديد من هم الساسة و الفنانين و رجال الأعمال الذين تم ذكرهم في الاخبار , وسردها بطريقة معينة .

و بالطبع نفس الأمر ينطبق علي اسماء شركات , اسماء دول , اسماء مدن , عملات , كتب , وهكذا .

و من الطرق التقليدية (والقديمة) لتفعيل هذا الأمر , ان يكون المدخل هو نص كلامي , ويكون المخرج هي عبارة عن سلسلة من الارقام 0 او 1 , لكل كلمة , حيث 0 هو ليس اسم , و 1 هو اسم

فمثلا لو كان النص :

Yesterday Mr Donald Trump visited the French president Marconi

ستكون مصفوفة المخرج

0 0 1 1 0 0 0 0 1

و من الممكن ان يكون التصنيف لأكثر من متغير , حيث 1 اسم شخص , 2 اسم مدينة , 3 اسم عملة , 4 اسم كتاب , 0 شئ مختلف
فلو كانت الجملة

yesterday peter went to NYC library to get ML book of Andrew with 5 dollar

فيكون المخرج

0 10 0 2 0 0 0 4 0 0 1 0 0 3

إلا أن هذه الطريقة قديمة نوعا , وحاليا هناط طرق أحدث في تحديد اسم الشخص او المدينة او البلد , ومكانها في النص و تكرارها و هكذا
و لبدء التعامل مع هذا النظام (القديم) , فلو قلنا ان هذه هي جملة المدخل :

Harry Potter and Hermione Granger invented a new spell.

فنقوم بتحديد features لكل كلمة فيهم , و ان يكون رمزها $x < t >$ مع اختيار الاقواس الزاوية هذه لتفريقها عن باقي الاقواس

Harry Potter and Hermione Granger invented a new spell.

$x^{<1>}$ $x^{<2>}$ $x^{<3>}$... $x^{<9>}$

كلّك في المخرج , سنقوم بتحديد قيم y بنفس الارقام و نفس الاقواس الزاوية للقيم المخرجة

$y^{(1)}$ $y^{(2)}$ $y^{(3)}$... $y^{(9)}$

كذلك سنقوم بتحديد T_x للدلالة علي عدد كلمات المدخل (هنا 9) , T_y للدلالة علي عدد المخارج (هنا كذلك 9 لكن في تطبيقات اخري قد لا تساوي عدد المدخل) ولا تنس أن الاقواس الدائرية تستخدم للدلالة علي عنصر محدد من عناصر العينة , لذا فإن :

$x^{(i)}_{(t)}$

معناها الكلمة رقم t في العنصر رقم i

فلو كانت $x^{(7)}_{(12)}$ فهي معناها الكلمة الثانية عشر في العنصر السابع من عناصر العينة .

ونفس الفكرة في y فلو كان لدينا $y^{(16)}_{(7)}$ تساوي 0 فهذا معناه ان الكلمة السابعة في العينة السادسة عشر هي ليست اسم

و لأن كل عنصر من عناصر العينة قد يحتوي عدد كلمات مختلف , فهناك رمز :

$$T_x^{(i)}$$

وهو الذي يدل علي عدد كلمات العنصر رقم i

فلو كان لدينا $T_x(20) = 16$ فهو معناها ان عدد كلمات العينة رقم 20 هي 16 كلمة .

كذلك $T_y(14) = 22$ معناها ان مصفوفة مخرج العينة 14 هي 22 رقم اصفار وواحيد .

و كل هذه التفاصيل هي في قلب ما يسمى NLP او معالجة اللغات الطبيعية Natural Language Processing

اذن كيف يقوم الخوارزم بقراءة هذه الجملة و التعرف عليها ؟

أول خطوة , هو ان تقوم انت كمبرمج بتكوين قاموس كبير خاص بك , يحتوي علي كل كلمات اللغة الإنجليزية التي قد تستخدمها (ليس بالضروري كل الكلمات بالفعل , ولكن كل الكلمات التي ستستخدمها هنا) .

فبفرض اننا سنقوم باستخدام 10 الاف كلمة مشهورة , فيتم رصهم جميعا في قاموس كبير بحيث يكون لكل كلمة رقم محدد يبدأ من 0 و تزيد :

a	0
aaron	1
.	.
.	.
and	367
.	.
harry	4075
.	.
potter	6830
.	.
zulu	10000

و قد تقوم انت باستخدام قاموس جاهز بالفعل , او قد تقوم بصناعة قاموس خاص بك , في حالة كانت القواميس الموجودة لا تفي باستخداماتك , فمثلا لو كنت تستخدم مصطلحات طبية عديدة و ليست موجودة في القواميس المتوفرة , فيمكنك عمل كود بسيط يبحث في ملايين السطور المتوفرة لديك , لايجاد كل الكلمات و مسح المتشابه , و رصها بالترتيب و إعطائها أرقام , ونفس الفكرة لو كنت ستقوم بعمل قاموس للغة أخرى (مثل العربية) .

و أحيانا يتم عمل خانة تسمى <UNK> و هي دلالة عن كلمة unknown و هي تكون للكلمات الغير موجودة في القاموس لديك

ثم يقوم الخوارزم بعمل ما يسمى : مصفوفة 1 الساخنة one-hot vector , وهي المصفوفة ذات العمود الواحد , التي تكون كل قيمها باصفار و ما عدا قيمة واحدة فقط , وتكون هذه القيمة برقم (1) بالتحديد

وهي أن يتم عمل مصفوفة عمود واحد , وبها 10 الاف صف (او عدد الصفوف الذي سيتم تحديده لطول القاموس الخاص بك) , وان تكون كل القيم اصفار , عدا القيمة الموجودة في مكان الكلمة نفسها تكون بواحد

فيقوم الخوارزم باستبدال كلمة **harry** في اول الجملة بالمصفوفة :

$$\begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 4075$$

حيث تكون كل القيم اصفار , عدا الصف رقم 4075 تكون بواحد

كذلك تم استبدال كلمة **potter** بالمصفوفة , التي تكون كلها اصفار عدا الصف الـ 6830 يكون بواحد :

$$\begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \leftarrow 6830$$

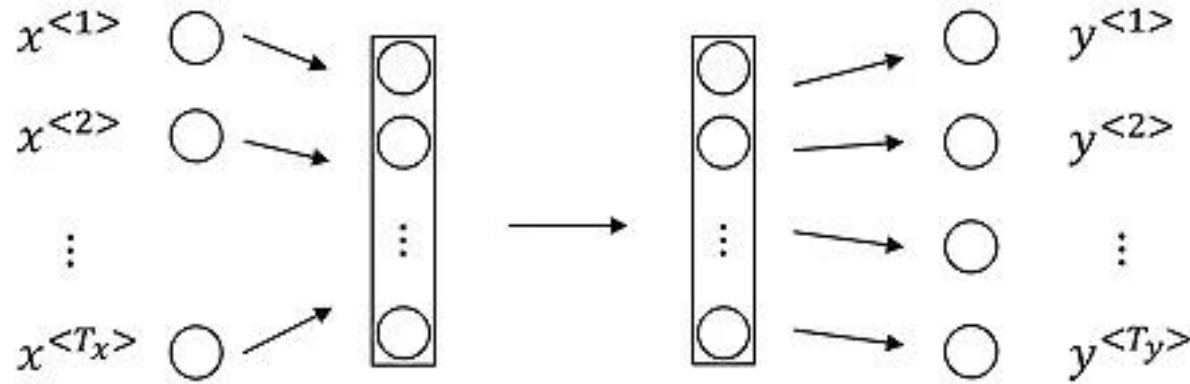
و هكذا في باقي الكلمات :

Diagram illustrating the bit-plane extraction process for a 10-bit input. The input is shown as a sequence of bits $x^{<1>}, x^{<2>}, x^{<3>}, \dots, x^{<4>}, \dots, x^{<9>}$. These bits are mapped to columns of a 10x10 matrix. The first column is labeled '4075', the second '6830', and the fourth '10,000'. The matrix is labeled 'One-bit' at the bottom.

* * * * *

و للتعرف علي الـ RNN علينا أولا ان نتعرف علي الشكل التقليدي الذي لن يتمكن من معالجة نموذج التتابع مثل هذا . .

إذا تناولنا المصفوفات الخاصة بكل كلمة من الكلمات , وقمنا بإدخالها في شبكة عصبية تقليدية مثل هذه , للوصول الي قيم y تحدد الاسماء و غيرها , فهي ستكون غير فعالة لسببين هامين . .



أولا أن عناصر العينة ليس لها اطوال ثابتة , فقد تكون الجملة الأولى 10 كلمات , والثانية 15 و الثالثة 6 و هكذا , ولأن كل كلمة لها مصفوفة , فكل عنصر من العناصر سيكون له عدد مصفوفات يتناسب مع عدد الكلمات , و بالتالي سيكون عدد الـ x مختلف من عنصر لآخر , وهو ما يصعب جدا عملية الـ NN العادية .

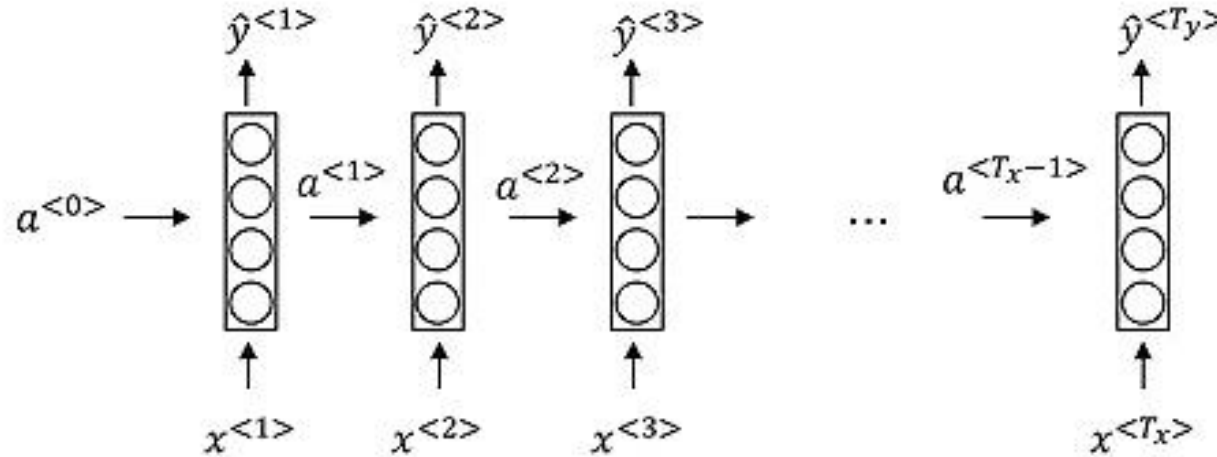
وإذا قمنا بالحساب علي العدد الاقصى للكلمات , فستظهر مشكلة الكفاءة القليلة , والتي ستنتج عن وجود عناصر كثيرة بها فجوات عديدة سيتم ملئها باصفار

السبب الثاني ان اصلا طريقة التعليم غير منطقية , فإذا كان في الجملة الأولي هناك اسماء اشخاص في الكلمة الأولي و الرابعة و السادسة , و في الجملة الثانية في الكلمة السابعة , وفي الجملة الثالثة في الكلمة الخامسة , فترتيب الكلمة ليس منطقيا في تحديد هل هذا اسم ام لا .

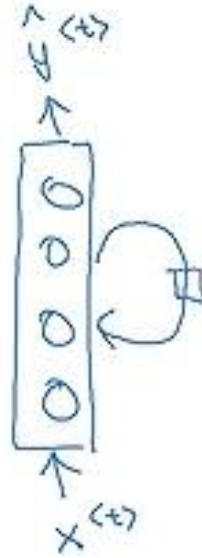
و من هنا كان لابد من اختيار تصميم مختلف للشبكات العصبية , الذي يتلائم مع هذا النوع من البيانات , وهو الشبكات المتكررة RNN

و تعتمد فكرة الـ RNN علي تناول أول كلمة في النص $x<1>$, وإدخالها في دالة لتحديد هل هي اسم ام لا , و من ثم تقوم الدالة بإخراج رقم $y<1>$ و الذي يحدد هل هي صفر او واحد , و حينما تقوم الشبكة بتناول الكلمة الثانية $x<2>$, تقوم باستخدام دالة الـ activation الناتجة من معالجة الكلمة الأولي (لا تعتمد عليها 100 % لكن تستخدمها بشكل مساعد او جزئي) , وكأن الدالة التي تقوم بتقييم الكلمة الثانية $x<2>$ ستكون معتمدة بشكل ما علي نتيجة الكلمة الأولي .

ثم يتكرر الأمر مع معالجة الكلمة الثالثة $x<3>$, والتي تعتمد علي نتيجة معالجة الكلمة الثانية (والأولي بشكل ما) , وهكذا حتي نصل للكلمة الأخيرة .



وكان الـ RNN يحدث فيها إعادة استخدام لدالة الـ activation كمدخل في كل مرة , لذا تسمى الشبكة المتكررة , و لذا يتم رسمها بهذا الشكل أحيانا :



و غالبا ما تكون مصفوفة البداية (التي يتم ادخالها لمعالجة أول كلمة) هي مصفوف اصفار , او ارقام عشوائية بسيطة , و تسمى $a_{<0>}$.

ومن خصائصها الهامة , ان معاملات كل مصفوفة من المصفوفات التي تعالج كل كلمة , هي نفس المعاملات , و تأخذ رمز Wax , كذلك معاملات دالة الـ activation لكل مصفوفات الكلمات تكون هي نفسها و يرمز لها Waa , أيضا معاملات المخرجات y تكون هي هي Wya (و هذا الأمر سنشرحه بالتفصيل بعد قليل)

ومن أهم عيوب الـ RNN أن توقع كل كلمة يعتمد علي الكلمات السابقة لها فقط , لأن معالجة $x_{<5>}$ تعتمد علي الكلمات الاولى و الثانية و الثالثة و الرابعة , ومشكلة هذا الأمر أن فعليا في النصوص قد يكون استنتاج ان هذا اسم شخص يعتمد علي كلمات تالية لها . .

فمثلا في جملة :

I think that Donald Trump at the white house is a reckless president for the US

فحينما ستقوم الـ RNN بمعالجة هذه الجملة , ستقوم بالربط بين كلمتي دونالد ترامب وبين الكلمات السابقة لها , بينما بعدها كلمات مهمة مثل , white , house , president , US و هي كلمات دلالية شديدة الأهمية و مرتبط ب اسم رئيس .

و سيتم معالجة هذا الأمر بشكل ما في ما يسمى الشبكة المتكررة ثنائية الاتجاهات bidirectional recurrent neural network و يرمز لها BRNN

و لتتبع العملية الحسابية في الـ RNN سنجد أن خطوة البداية تكون من مصفوفة $a_{<0>}$ والتي تكون بقيمة اصفار .

ثم لحساب اول اكتيفاشن $a_{<1>}$ تكون المعادلة :

$$a_{<1>} = g (W_{aa} a_{<0>} + W_{ax} x_{<1>} + b_a)$$

أي أنها مصفوفة W_{aa} (التي قلنا أنها ستكون ثابتة لكل الشبكة) , مضروبة في قيمة $a^{<0>}$ (قيمة الاكتيفاشن السابقة لها وهي حاليا الأصفار) , مجموعة علي حاصل ضرب W_{ax} و هي كذلك ثابتة لكل الشبكة مع قيمة $x^{<1>}$ الخاصة بالكلمة نفسها , ثم معامل انحراف b_a وهو كذلك ثابت للشبكة , وكل هذا يدخل في دالة $g()$ والتي غالبا يتم استخدام \tanh أو ReLU لها .

ثم ننتقل لحساب قيمة y هات , والتي تكون بالمعادلة :

$$\hat{y}^{<1>} = g (W_{ay} a^{<1>} + b_y)$$

وهي حاصل ضرب W_{ay} (كذلك ثابتة للشبكة) في الاكتيفاشن الحالية لها (التي تم حسابها من المعادلة السابقة) , مجموعة علي معامل انحراف ثابت b_y , وكل هذا يدخل في دالة $g()$, والتي غالبا تكون سيجمويد او اي نوع من انواع السوفت ماكس , و ذلك لأن حساب \hat{y} هي عملية تصنيف **classification** فنريد قيمة 0 , 1

لا تنس أن دالة g الأولى غير الثانية , حتي ان عدد من الباحثين يكتبون $g1()$, و $g2()$

فتكون الصيغة العامة :

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

كذلك لا تنس أن المصفوفات الثوابت (والتي يقوم الخوارزم بحسابها) هي كلا من :

W_{ax} , W_{ay} , W_{aa} , b_x , b_y

كلا من البيئات معروفة , أما قيم الـ w الثلاثة فكلا منها له مكان مرتبط به , ويسمي باسمه

وهنا يجب أن نقوم بعملية حسابية هامة بهدف التبسيط . .

في المعادلة الأولى نري صيغتها غير مريحة و هي :

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

حيث أن هناك مصفوفتين معاملات W_{aa} , W_{ax} و قيمتين يتم ضربهم فيها و هي , $x^{<t>}, a^{<t-1>}$ فسنقوم حالا بدمجهم بشكل ذكي . .

إذا افترضنا ان مصفوفة $a^{<t-1>}$ هي 100 صف في عمود واحد , فيجب أن تكون مصفوفة W_{aa} هي $100*100$ حتي انها حينما تضرب فيها تقوم بانتاج 100 في 1 (وهي القيمة التي ستخرج لنا في $a^{<t>}$)

كذلك إذا كانت مصفوفة $x^{<t>}$ هي 10 الاف صف في عمود واحد (لا تنس أنها مرتبطة بالقاموس) . فجيب أن تكون مصفوفة Wax هي 100×10000 حتي إذا تم ضربها فيها تكون 100 في 1

أي أن :

$$a^{<t>} = g(W_{aa}a^{<t-1>} + W_{ax}x^{<t>} + b_a)$$

هي :

$$100 \times 1 = 100 \times 100 \times 100 \times 1 + 100 \times 10000 \times 10000 \times 1 + 100 \times 1$$

فرغبة في التبسيط , سنقوم بلصق مصفوفتي Waa , Wax بشكل أفقي hstack, بهذا الشكل :

بحيث تكون مصفوفة واحدة , لها 100 صف , ولها 10100 عمود معا , و يكون اسمها W_a

ثم نقوم بلصق مصفوفتي $a^{<t-1>}$, $x^{<t>}$ بشكل رأسي vstack هكذا :

$$[a^{<t-1>}, x^{<t>}] = \begin{bmatrix} a^{<t-1>} \\ x^{<t>} \end{bmatrix} \begin{matrix} \updownarrow 100 \\ \updownarrow 10100 \end{matrix} \begin{matrix} \updownarrow 10100 \end{matrix}$$

بحيث تكون المصفوفة الجديدة عمود واحد , و 10100 صف .

فحينما يتم ضرب المصفوفة الجديدة W_a في المصفوفة الجديدة $[a^{<t-1>}, x^{<t>}]$ تكون :

$$W_a * [a^{<t-1>}, x^{<t>}] = 100 * 10100 \times 10100 * 1 = 100 * 1$$

تكون نفس ابعاد التصور القديم هذا , ونفس قيمه كذلك :

$$W_{aa}a^{<t-1>} + W_{ax}x^{<t>}$$

و بالنسبة للمعادلة الثانية

$$\hat{y}^{<t>} = g(W_{ya}a^{<t>} + b_y)$$

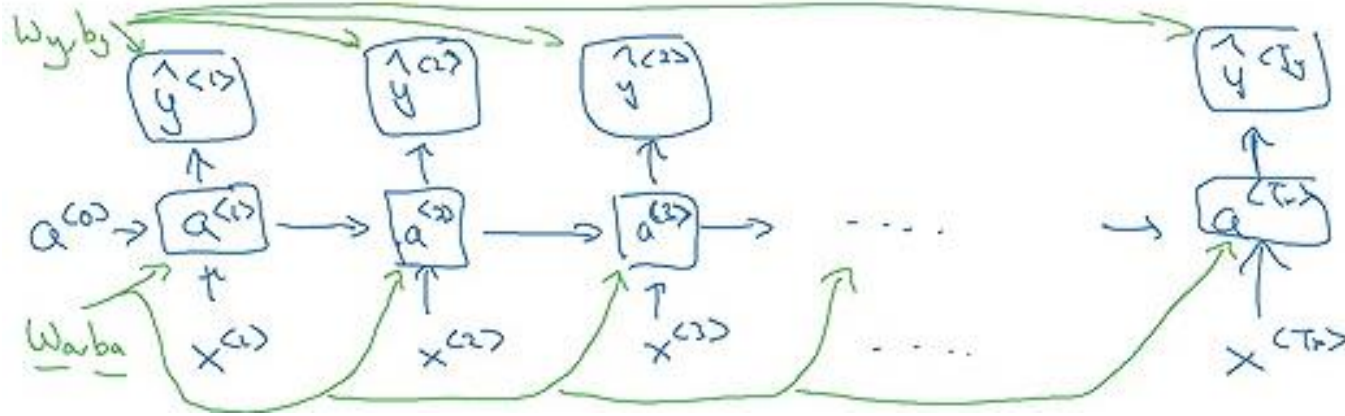
سنقوم فقط بتغيير اسم Wya الي Wy مما يجعل المطلوب الكلي في الخوارزم ايجاد Wa, Wy, ba, by , بشكل ايسر .

* * * * *

حسنًا ماذا عن المسار العكسي للبيانات ؟ back propagation

و بالطبع مهمتها الأساسية تحديد قيم الأوزان W_a , W_y , b_a , b_y التي ستقوم بصياغة الشبكة .

نبدأ أولاً برسم المسار الأمامي بشكل تفصيلي :



فقيمة a_0 , x_1 تستخدمان معاً لإنتاج a_1 والتي تساعد x_2 في إنتاج a_2 وهكذا , وكل هذا يتم عبر المعاملات W_a, b_a

ولإنتاج قيمة y_1 نحتاج لقيمة a_1 وذلك عبر المعاملات W_y, b_y والتي ستستخدم كذلك في باقي الوايات .

هذا هو المسار الأمامي , فماذا عن معادلة الخطأ loss function ؟

لا تنس أن معادلة الخطأ شديدة الأهمية حتي يتم صياغة المسار الخلفي علي اساس تقليلها . .

لأن هذا الخوارزم يقوم في النهاية بتوقع كل كلمة هل هي اسم ام لا , فستكون معادلة الخطأ شديدة التشابه مع معادلة خطأ التصنيف logistic regression , فتكون المعادلة لكل قيمة y هي :

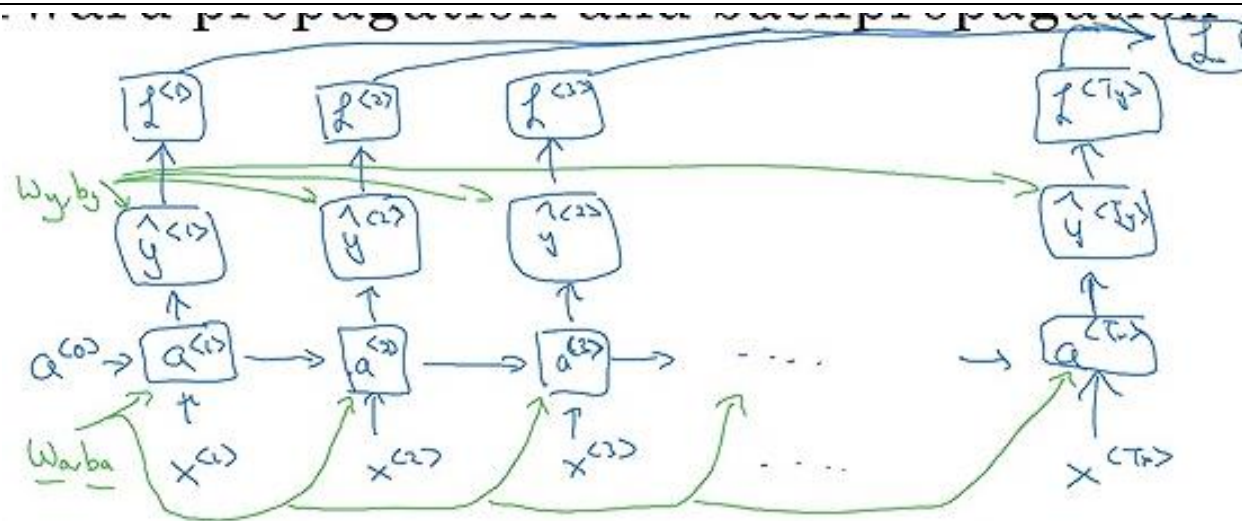
$$\mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>}) = - y^{<t>} \log \hat{y}^{<t>} - (1 - y^{<t>}) \log (1 - \hat{y}^{<t>})$$

و تكون المعادلة الكلي للنص بالكامل , هي مجموع قيم أخطاء كل كلمة :

$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{|x|} \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

حيث السمشن هنا يكون الي Tx وهو عدد كلمات النص

و بها يتم حساب قيمة الخطأ في كل كلمة كالتالي :

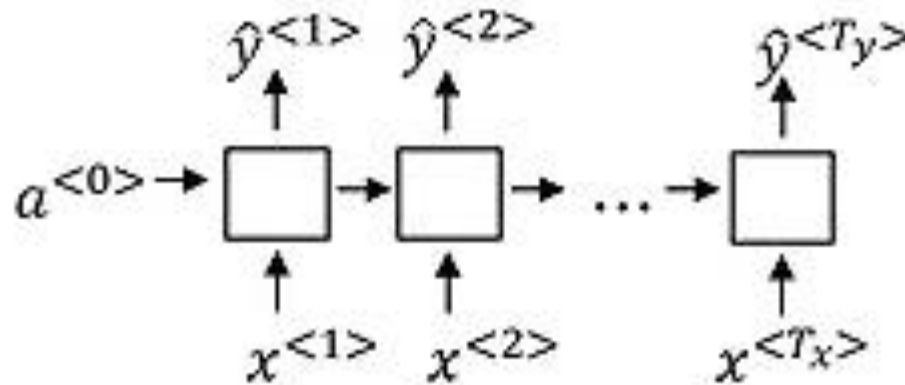


و لاستخدام المسار الخلفي back propagation , علينا فقط تتبع الخطوط بشكل عكسي من اليمين الي اليسار كما في الاسهم الحمراء .

سنتناول الان عدد من أنواع الشبكات المتكررة RNN و التي في عدد منها يكون عدد المدخلات T_x لا يساوي عدد المخرجات T_y

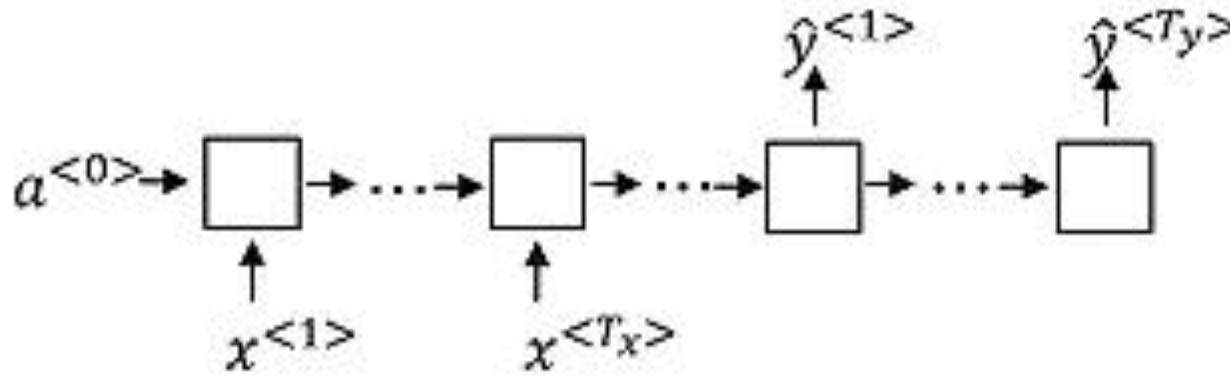
اذا تذكرنا التطبيقات الذي ذكرناها في بداية الملف , سنجد أن هناك نماذج حينما يكون المدخل لا شيء (صناعة موسيقي) , او نماذج حينما يكون المخرج رقم واحد (الكلام يتحول لتقييم) , او نماذج حينما يكون المدخل و المخرج بنفس الرقم (تحديد الاسماء) , او نماذج حينما يكون المدخل و المخرج غير معروف في رقمهم (الترجمة)

فأحد النماذج المستخدمة يسمى many-to-many و الي يكون له العديد من المدخلات يصنع العديد من المخرجات , مثل المثال السالف شرحه في الفيديوهات السابقة (اختيار الاسماء) , حيث أن هناك مدخلات عديدة , وكل منها يكون لها مخرج معين :

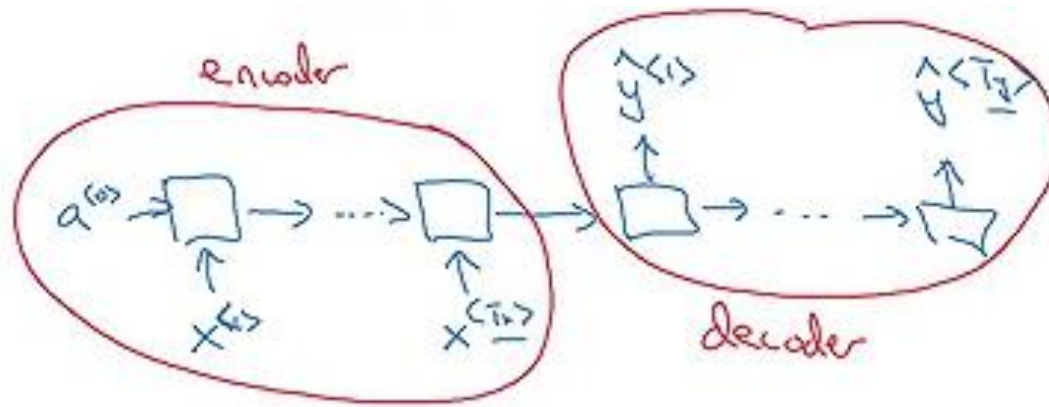


و هناك شكل اخر من نفس النوع , لكن حينما يكون عدد المدخل T_x لا يساوي عدد المخرج T_y و هي مثلا في حالة الترجمة , فلو تم ترجمة جملة فرنسية للانجليزية , فلن تكون عدد الكلمات متساوي و كذلك لن تكون كل كلمة انجليزية صانعة للكلمة الفرنسية بالمقابل .

و تكون الرسة الخاصة بها :



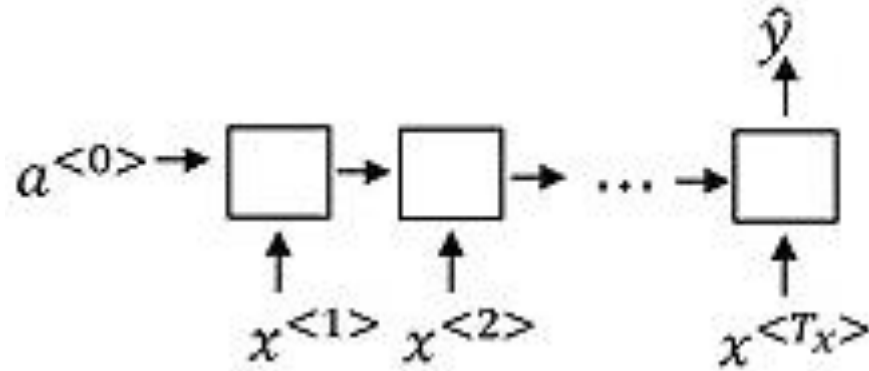
حيث يقوم الخوارزم اولا بتناول كل المدخلات و معالجتها جيدا , ثم يقوم بصناعة المخرجات في النهاية , وكان الجزء الأول هو القارئ encoder و الجزء الثاني هو الكاتب decoder



نموذج آخر يسمى : many-to-one و الذي يكون له مدخلات عديدة , لكن مخرج واحد .

مثل لو كان لدينا خوارزم يقوم بفحص التعليقات علي فيلم معين , لمعرفة تقييم الشخص لها و فالتعليق هو نص له كلمات عديدة , بينا التقييم هو رقم واحد (من 1 الي 10 مثلا)

فلو كان التعليق المكتوب : there is nothing interested in this movie , فتكون الشبكة المتكررة :



حيث ان لكل كلمة مدخل معين , و بدلا من وجود مخارج لكل كلمة , يكون هناك مخرج واحد في النهاية يعتمد علي كل المداخل .

كما أن هناك one-to-many مثل خوارزميات انتاج الموسيقى , حيث تقوم بإعطائه رقم محدد (نوع الموسيقى المطلوب) , ويقوم هو بإنتاج العديد من المخرجات بهذا الشكل

سنتناول الان ما يسمى نموذج اللغة language model و الذي هو أحد التطبيقات الهامة في التعامل مع اللغة , ويتم بنائها بالـ RNN .

و للفهم الكامل لمعني language model يمكن أن نقول , أنه قدرة الآلة علي توقع الكلمة التي قلتها صوتيا , عبر فهم المعني الكامل حولها , وهو الذي يكون صعبا حينما تكون الكلمة المنطوقة لها اكثر من spelling

فمثلا لو قال شخص :

the apple and pear salad

أو قال

the apple and pair salad

فممكن الممكن لكمة pear ان تكتب pair لأنها تقريبا نفس النطق , فعلي عاتق الخوارم ان يعلم اي كلمة مقصودة عبر فهم المعاني المحيطة .

ونفس الفكرة في الفارق بين كلمات toe , too , two , to خاصة مع اختلاف لهجة من يقولها .

و يبدأ الأمر عبر تحديد احتمالية probability للكلمة المنطوقة , وسط سياق الجملة . وبالتالي يقوم الخوارزم بتحديد مدي احتمالية صحة كل جملة .

و تكون الخطوة الأولى في بناء الخوارزم , هو إحضار كمية كبيرة من بيانات التدريب corpus of English text و كلمة corpus معناها كمية كبيرة من البيانات .

بعدها نقوم بعملية اسمها tokenization يعني الترميز, وهي عملية ان كل كلمة من كلمات العينة اعمل لها رمز y .

فمثلا جملة :

cats average 15 hours of sleep a day

اعمل للكلمة cats رمز $y<1>$ و كلمة average $y<2>$, وهكذا .

مع ملاحظة ان حرف a او an او or او & سيكون ليها كمان رمز مستقل , بل ان احيانا علامات الترقيم (, . / * + -) هي كمان يكون لها رمز مستقل .

و أحيانا يتم عمل في نهاية الجملة ما يسمى <EOS> يعني end of sentence توضع في النهاية , حتي يعلم الخوارزم ان الجملة قد انتهت بالفعل , و اذا ما تم استخدامها فهي ايضا تأخذ رمز مستقل

و في حالة ن كان هناك كلمة ليست في القاموس (الكلمات الاكثر استخدام) وقتها يقوم الخوارزم باستبدالها برمز محدد هو <UNK> وهي اختصار unknown للدلالة علي كلمة غير معروفة في القاموس .

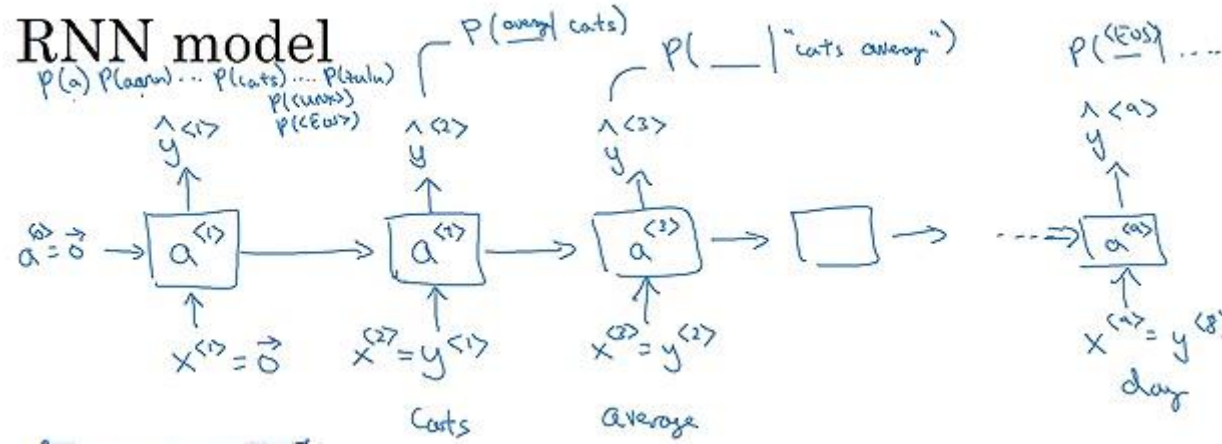
ولتتبع مسار الخوارزم , في الجملة :

cats average 15 hours of sleep a day

في البداية ان المدخل هو a_0 و التي هي مصفوفة كلها اصفار , وهي تدخل في الوحدة $a<1>$ والتي أيضا تأخذ مدخل هو : $x<1>$, ومهمة هذه الوحدة ان تقوم بيجاد احتمالية الكلمة الأولى و هي cats , وهذه الاحتمالية هي $\hat{y}<1>$. اي واي هات المتوقعة

ويتم حساب الاحتمالية عبر استخدام softmax به 10 الاف مخرج , وهو كلمات القاموس المستخدم .

وقيمة الاحتمالية تعتمد علي فرصة أن تأتي هذه الكلمة في البداية من قلب الكلمات العشر الاف في القاموس



ثم نذهب للوحدة الثانية , وهي التي كون المدخل فيها هو $x^{(2)}$ والتي هي نفسها قيمة $y^{(1)}$, وهي الكلمة الأولى بالفعل cats .

لا تنس أن قيمة \hat{y} هي الكلمة المحتملة , بينما y هي القيمة الحقيقية (كلمة cats هنا)

اذن قيمة $x^{(2)}$ هي نفسها $y^{(1)}$, والتي ستدخل في الوحدة $a^{(2)}$ لتتوقع الكلمة التالية و هي average , ولكن التوقع سيكون بناء علي معلومية الكلمة الأولى cats فسيكون هو التوقع المشروط probability with condition والذي يكتب هكذا :

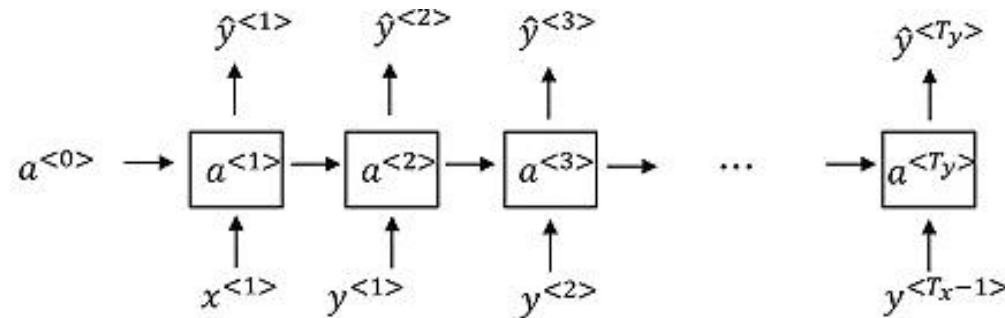
$P(\text{average} \mid \text{cats})$

و يتم تكرار الأمر في الوحدة التالية , لتوقع الكلمة الثالثة fifteen بمعلومية الكلمتين السابقتين الاولى و الثانية فتكون هكذا :

$P(\text{fifteen} \mid \text{average} , \text{cats})$

ويستمر الأمر حتي نصل للكلمة الأخيرة EOS والتي نريد أن نتوقعها بناء علي الكلمات السابقة كلها . فكل الخوارزم يتناول الكلمات من اليسار لليمين , و يقوم بتدريب الثبتات عليه .

والصورة الكاملة هنا :



وتكون معادلة الخطأ لكل كلمة علي حدة هي :

$$\mathcal{L}(\hat{y}^{<t>}, y^{<t>}) = - \sum_i y_i^{<t>} \log \hat{y}_i^{<t>}$$

بينما معادلة الخطأ الكلية تكون مجموع نسب الأخطاء السابقة

$$\mathcal{L} = \sum_t \mathcal{L}^{(t)}(\hat{y}^{(t)}, y^{(t)})$$

و بالتالى إذا ما اردنا حساب احتمالية كتابة جملة من 3 كلمات هي :

cats average 15

ولأنها احتمالية متراكمة , فتكون الاحتمالية الكلية هي حاصل ضرب الاحتماليات الثلاث , فتكون :

$$P(y_{<1>}, y_{<2>}, y_{<3>}) = P(y_{<1>}) P(y_{<2>} | y_{<1>}) P(y_{<3>} | y_{<1>}, y_{<2>})$$

* * * * *

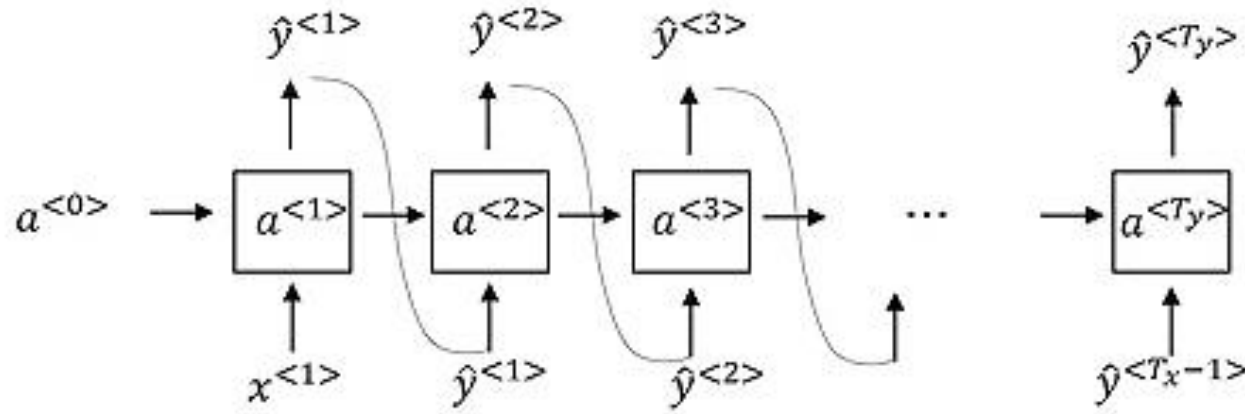
بعدها يتم عمل تدريب للخوارزم , يمكنك التأكد من مدي كفاءته عمر تطبيق ما يسمى متتبع عينة الرواية `sampling novel sequences`

والمقصود به انه يمكن توقع كلمات من لا شيء , بل وكتابة سلسلة كاملة منها علي هيئة مقال او قصة .

و يبدأ الأمر عبر ادخال قيمة `a0` الصفرية , ثم `x<1>` التي ستكون عشوائية او صفرية , فيقوم الخوارزم (بعد التدريب) باستنتاج الكلمة الأولى , و سيكون الاستنتاج بناء علي ما هي الكلمات المشهورة بانها تاتي اولا , فغالبا سيختار كلمة مثل `the , a , an , once` , حسب قيمة `y`

لا تنس أن واي هات هي قيمة واحدة من 10 الاف من السوفت ماكس التي سيختارها من القاموس .

ثم نأخذ هذه الكلمة `y<1>`, ونقوم بإدخالها في الوحدة الثانية `a<2>` والتي ستكون هي بالفعل قيمة `x<2>` ليقوم باستنتاج الكلمة التالية بعد بمعلومية الكلمة الأولى `the` و سيختار كلمة و ليكون `girl` .



و يتم استمرار العملية ليتم استنتاج الكلمة الثالثة بناء علي الاولى و الثانية , ثم الرابعة بناء علي الكلمات الثلاث الأولى , وهكذا .

و حتي نهاية الجملة ستأتي عبر استنتاج كلمة $\langle \text{EOS} \rangle$, او يمكن عمل الأمر يدويا عبر تحديد حد معين للكلمات برقم محدد .

و هذا الأمر سيجعل الخوارزم يقوم بصياغة جمل كاملة و علامات ترقيم و كل شيء , سواء كانت رواية او مقال او تحليل او خبر .

كل ما سبق كان يسمى النموذج اللغوي للكلمات word level language model بينما هناك مستوي ادق و اعرق يسمى character level language model وهو الخاص بالتحليل بالحروف .

و يقصد به أنه يتم تحويل الجملة ليس لكلمات , و لكن لحروف , فإذا كانت نفس الجملة :

cats average 15 hours of sleep a day

فيكون حرف c هو واي 1 , وحرف a هو واي 2 و هكذا , فيكون القاموس لا يحتوي علي 10 الاف كلمة , ولكن علي الحروف الخاصة باللغة الانجليزية + الارقام + علامات الترقيم + نفس الحروف كابيتال + رموز هامة مثل \$#%\$

و ميزة هذه الفكرة انها ستتجنب مشكلة الكلمات غير المعلومة , فكل حرف يتم تحليله , و عيبه ان عدد الوحدات اكبر بكثير و هذا سيجعل التحليل بطئ للغاية و بالتالي مكلف .

و في الوقت الحالي ينذر استخدام هذا النموذج , لكن قد يتم استخدامه في المستقبل حينما تكون اجهزة الكمبيوتر اسرع .

و لا تنس ان الخوارزم سيقوم بالتنبؤ بجمل , اعتمادا علي قيم ثيتا التي تم تدريبه عليها , وبالتالي لو تم تدريبه علي أدب او اخبار او رياضة او فن او غيره , فسيقوم بكتابة جمل في نفس السياق

هنا نصوص تم التنبؤ بها من نوع الأخبار

News

President enrique peña nieto, announced
sench's sulk former coming football langston
paring.

"I was not at all surprised," said hich langston.

"Concussion epidemic", to be examined.

The gray football the told some and this has on
the uefa icon, should money as.

و هنا من نوع أدب شكسبير

Shakespeare

The mortal moon hath her eclipse in love.

And subject of this thou art another this fold.

When besser be my love to me see sabl's.

For whose are ruse of mine eyes heaves.

* _ * _ * _ * _ * _ * _ * _ * _ * _ * _ * _ * _ *

نتكلم الان هم مشكلة شائعة في الشبكات المتكررة RNN و هي ما يسمى اختفاء الاشتقاقات Vanishing gradients in RNN

و للتعرف علي معناها , تعالي نلقي نظرة علي جملتين :

the cat which ate fish & drink the fresh milk , was brown
the cats which ate fish & drink the fresh milk , were brown

كما هو واضح ففي الجملة الأولى قطة واحدة فتم استخدام was بينما في الجملة الثانية قطط جمع فتم استخدام were .

و المشكلة في أن الفاصل المكاني بين كلا من الكلمتين هنا و هناك كبير , وهذا سيجعل من الصعب علي الخوارزم ان يقوم بالربط بين الكلمة المؤثرة (cat) والكلمة الناتجة (was)

كما أننا في RNN نقوم بعمل backpropagation و بالتالي ستظهر مشكلة اختفاء الاشتقاقات vanishing gradients و هي التي تكلمنا عنها في الكورس الأول , حينما ذكرنا الشبكات شديدة العمق (100 طبقة مثلا) فأثناء عمل المسار الخلفي فالاشتقاقات تتلاشي حينما تصل للطبقات الأولى , مما يجعل تدريب الطبقات الأولى و بالتالي تعديل قيم الثبتات شئ غير سهل .

و لأن فكرة الـ RNN الخاصة بالبناء اللغوي , تعتمد علي ان كل كلمة لها وحدة و بالتالي طبقة , فتأثير الكلمات المتباعدة سيكون ضئيل .

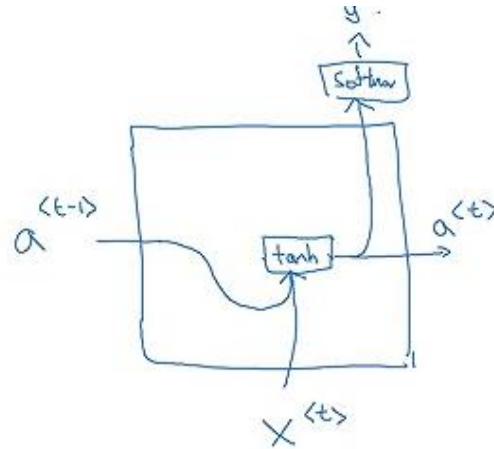
و ليست هذه هي المشكلة الوحيدة في الشبكات العميقة , فكما ان الاشتقاقات ممكن ان تختفي , فهي ممكن ان تزيد بشكل رهيب في ما يسمى مشكلة exploding gradients

نتكلم الآن عن أحد الحلول لمشكلة اختفاء الاشتقاقات , وهي باستخدام تقنية وحدة البوابة المتكررة Gated recurrent unit و التي باختصار GRU .

وإذا تذكرنا المعادلة الخاصة بالـ RNN تكون :

$$a^{<t>} = g(W_a[a^{<t-1>}, x^{<t>}] + b_a)$$

والتي يمكن تصميم رسم بياني مبسط لها وهو :



و علينا الان ان نتعرف علي قيمة جديدة و هي : خلية الذاكرة memory cell و التي يرمز لها بالرمز c .

و في هذا المثال ستكون بالضبط مساوية لدالة الـ activation لذا ستكون هي نفسها قيمة a لكن في تطبيقات أخرى سنتناولها ستكون مختلفة . .
إذن ستكون معادلتها :

$$\tilde{c}^{<t>} = \tanh(W_c[c^{<t-1>}, x^{<t>}] + b_c)$$

والتي تتطابق مع معادلة a , مع التأكيد اننا نستخدم دالة تاننش هنا كدالة $g(x)$, ايضا لا تنس ان الـ c هنا تسمى c tilde لأن فوقها علامة ~
كما سنتناول قيمة أخرى و هي جاما , والتي تم تسميتها بهذا لتشابه حرف الـ G في كلمة $gate$ و هي البوابة , المشتقة من اسم التكنيك GRU
و تكون معادلته :

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

و هنا الدالة سيجمويد بدلا من التاننش .
و لأن دالة سيجمويد تخرج لنا قيمة تتراوح بين الصفر و الواحد , فأغلب القيم التي تكون لجاما هنا تكون كبيرة للغاية (1) , او صغيرة للغاية (0) .
وهنا نقوم الـ c بالربط بين الكلمة المطلوبة was و الكلمة البعيدة cat دون الوقوع في مشكلة اختفاء الاشتقاقات , وتستخدم المعادلة :

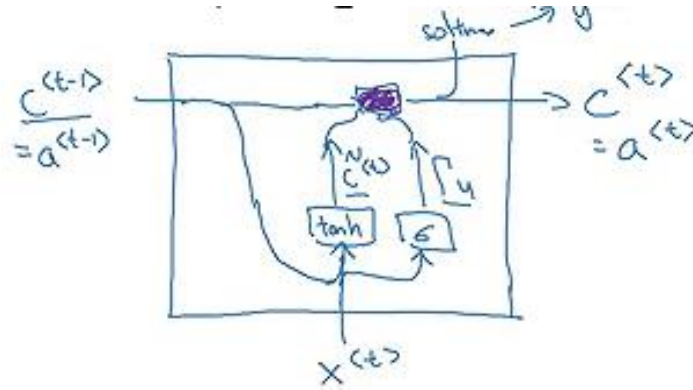
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

وهي تشبه معادلة التصنيف القديمة

$$J(\theta) = -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_{\theta}(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\theta}(x^{(i)})) \right]$$

و هي تحمل اسم بوابة , لأنها تتحكم في ما قيم ثبثات التي سيتم تعديلها (فتح البوابة) , وما هي القيم التي لن يتم تعديلها (غلق البوابة) .

فكأن في المسار الخلفي للشبكة , لدي كلمة cat ستكون قيمة جاما تساوي 1 , أي أن قيمة c (وهي كأنها ثبثا) سيكون لها قيمة ثبثا المعدلة (سي تيلدا) , بينما في باقي الكلمات ستكون قيمة جاما تساوي صفر , فتكون c عبر المعادلة السابقة لها نفس قيمة c السابقة ولا يتم تعديلها update وهو الشئ المطلوب



و كما هو موضح بالرسم , دخل قيمة c السابقة مع قيمة x لمعادلتي c الحالية بدالة تاناش و معادلة جاما بدالة سيجمويد , و ثم تدخل علي المربع الازرق و هي المعادلة الاخيرة المذكورة :

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) + c^{<t-1>}$$

فيكون المخرج هو قيمة y عبر السوفتماكس , وقيمة c الحالية

و غالبا ما يتم إضافة جاما r في المعادلة الأول هكذا :

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

حيث معادلتها هي :

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

* * * * *

كما تكلمنا عن GRU كوسيلة لحل مشكلة المسافات البعيدة بين العوامل الهامة , والتي قوانسنها هكذا :

$$\tilde{c}^{<t>} = \tanh(W_c[\Gamma_r * c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[c^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_r = \sigma(W_r[c^{<t-1>}, x^{<t>}] + b_r)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

$$a^{<t>} = c^{<t>}$$

, نتناول الان وسيلة اخري تسمى الذاكرة طويلة قصيرة المدى long short term memory والتي اختصارها LSTM , وهي أقوى و أكثر كفاءة من GRU

والقوانين الخاصة بها هي :

$$\tilde{c}^{<t>} = \tanh(W_c[a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u[a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f[a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o[a^{<t-1>}, x^{<t>}] + b_o)$$

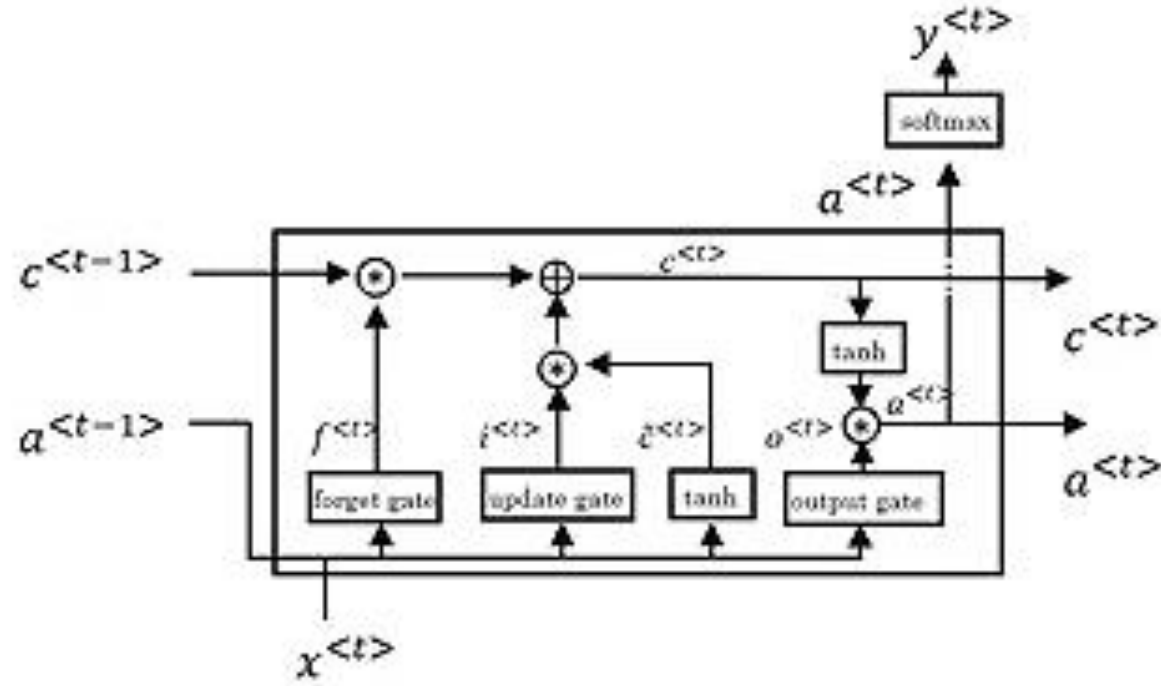
$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

$$a^{<t>} = \Gamma_o * \tanh c^{<t>}$$

حيث أن القيمة u تعبر عن التعديل update بينما القيمة f تعبر عن عدم التعديل و النسيان forget اما القيمة o تعبر عن المخرج output

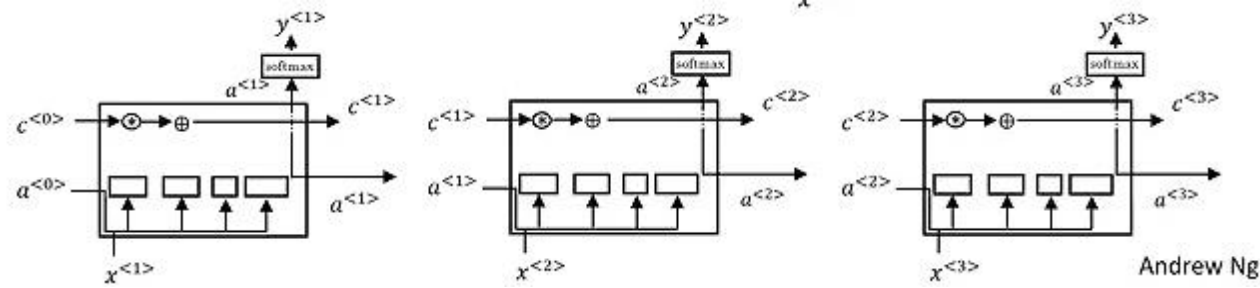
كذلك نلاحظ ان قيمة c لا تساوي قيمة a كما في GRU , كما انه تم فصل جاما للتعديل u عن جاما عدم التعديل f

وتكون الرسمة الخاصة بها هكذا :



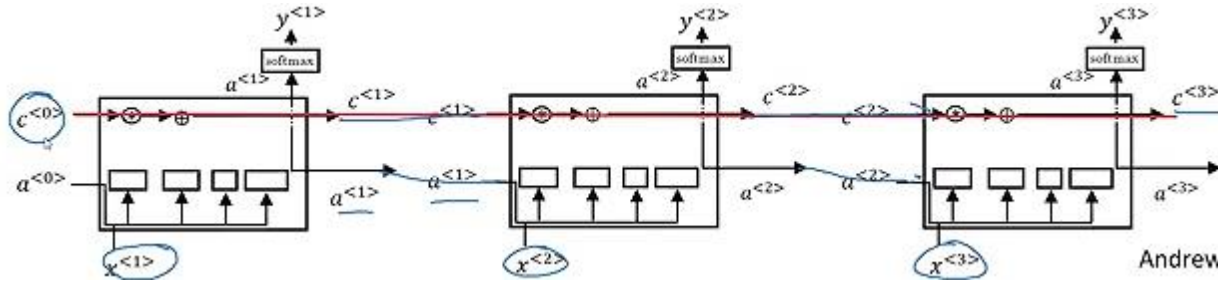
فمن الرسم يتضح ان قيمة $a^{<t-1>}$ (اسفل اليسار) تستخدم لتعديل كلا من forget gate , update gate , output gate كما نري في الجامات الثلاثة في المعادلات , بالإضافة انها تستخدم مع تاننش لإيجاد (سي تيلدا تي)

وكان السلسلة كاملة ستكون :



Andrew Ng

و من عوامل قوة LSTM انه في حالة كان هناك لايا عديدة سيتم نسيانها , فستنتقل قيمة الـ c سريعا من وحدة الي اخري دون تفليلها او غيرها , كما نري في الخط الأحمر



Andrew Ng

و هي ايضا مثل الـ GRU في فكرة ان عنصر من طبقة معينة (كلمة محددة) قد يقوم بالتأثير علي طبقة اخري بعيدة عنها (كلمة اخري) , دون التأثير علي الباقي بالضرورة .

نصل للسؤال الهام , أيهما افضل الـ GRU أم الـ LSTM ؟

لا توجد إجابة لهذا السؤال , إذ أن كلا منهما له مميزات و عيوب ضد الآخر . .

فالـ GRU لأنه بسيط التكوين , فهو اسرع في التنفيذ , و يتناسب مع الشبكات العميقة , لكن كفاءته ليست مثلي , بينما الـ LSTM العكس , فهو أكثر تعقيدا , و ابطئ في التنفيذ , لكنه ذو كفاءة أعلى

* * * * *

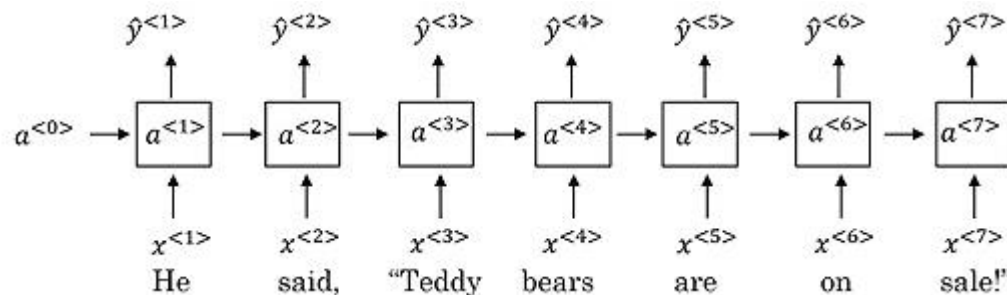
نتكلم الآن عن الشبكة المتكررة ذات الإتجاهين Bidirectional RNN . ويطلق عليها BRNN وهي التي تسمح للشبكة بالتعامل مع الكلمات السابقة و التالية معا , في الإتجاهين .

فكما رأينا في مثال سابق, قد يكون لدينا جملتين :

he said : teddy bears are fine

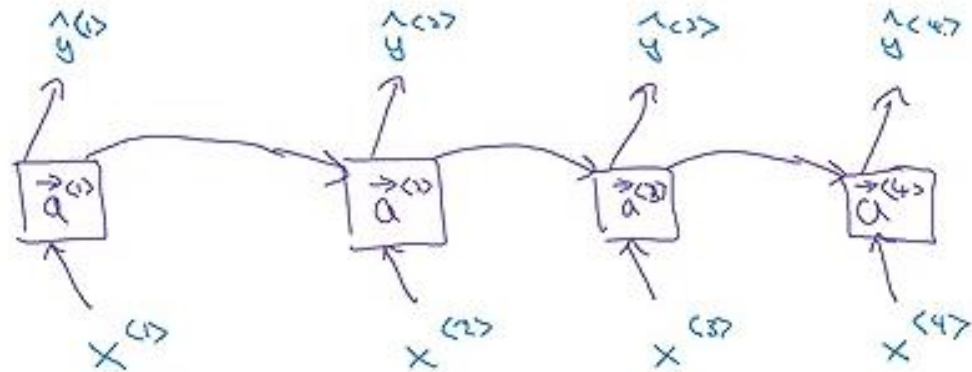
he said : teddy roosevelt was a US president

فالمشكلة ان الكلمات السابقة لكلمة teddy غير كافية لتحديد معناها و بالتالي تحديد الكلمات التالية لها , كما في الرسم :

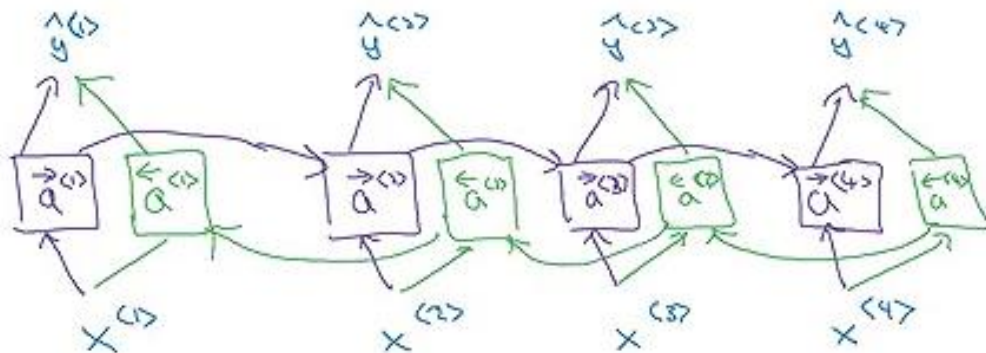


لذا علينا التعامل مع الشبكة ذات الاتجاهين .

فإن كان المسار الطبيعي لشبكة RNN هو هكذا :



فالمسار العكسي لها سيكون كالتالي :



حيث الوحدات الزرقاء تمثل المسار الأمامي , بينما الوحدات الخضراء تمثل الخلفي ..

وبالتالي ستكون معادلة y هي :

$$\hat{y}^{(t)} = g(w_y [\vec{a}^{(t)}, \vec{a}^{(t)}] + b_y)$$

والتي ستجعل قيمة y معتمدة علي كلا من المسار الأمامي و الخلفي معا , و هنا حينما يكون هناك كلمة teddy سيتدرب الخوارزم علي الكلمات السابقة و التالية له

و غالبا ما تكون الوحدات المستخدمة ليس فقط شبكات RNN وحدها بل يكون معها وحدات LSTM او GRU

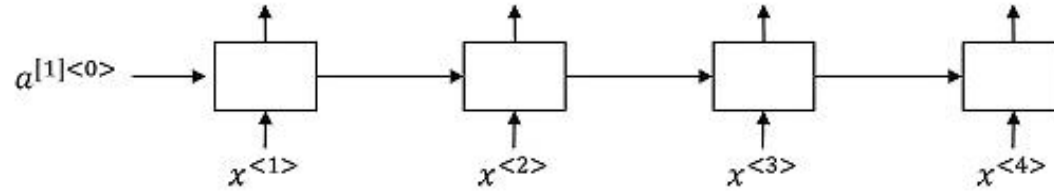
لكن من عيوبها انها تحتاج الي البيانات كاملة لبدء المعالجة , وبالتالي هي لا تتناسب مع البيانات التي تأتي بشكل آني و متواصل , مثل تحويل الصوت القادم الي نصوص

* * * * *

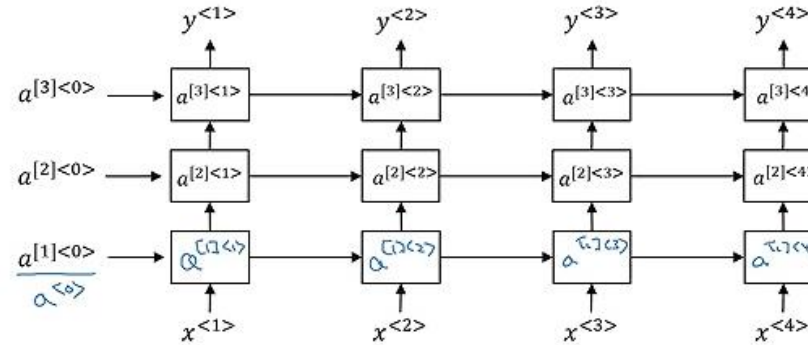
في نهاية الأسبوع الأول , سنتكلم عن الشبكات المتكررة العميقة Deep RNN , والتي تتعلق ببناء انواع مختلفة من النماذج التي رأيناها معا .

فكل نوع من الأنواع السابق ذكرها يمكن ان يعمل وحده جيدا , لكن ماذا عن ربطهم معا ؟

فإذا كانت هناك طبقة محددة من الـ RNN هكذا :



فدمجهم مع يمكن أن يكون هكذا :



فتكون كل قيمة تعتمد علي قيم مجاورة او سابقة لها . .

فمثلا القيمة $a[2]<3>$ وهي الخاصة بالكلمة الثالثة في العينة الثانية . فستكون معادلتها :

$$a^{[2]<3>} = g(w_a^{[2]} [a^{[1]<2>}, a^{[1]<3>}] + b_a^{[2]})$$

حيث ان لها مدخلين هما $a[2]<2>$ وهي الكلمة السابقة لها في نفس العينة , و $a[1]<3>$ وهي نفس ترتيب الكلمة , لكن في العينة السابقة

مع التأكيد علي أن كل طبقة سيكون لها معاملات b , w ثابتة لكل الوحدات في نفس الطبقة , لكن ستختلف بالطبع مع طبقة تالية

و لأن RNN هي معقدة بطبيعتها , فنقوم بتجنب ربط عدد كبير من الطبقات معا لتجنب صعوبة الحساب , لذا يتم احيانا بناء شبكة RNN ثلاثية الابعاد, حيث تكون هذه الطبقات الثلاثة وراءها 3 طبقات اخري , ثم 3 طبقات أخرى و هكذا , مما يسهل الحساب نوعا .

و قد تكون هذه الوحدات RNN عادية , و قد تكون ملحق بها GRU or LSTM لكن نادرا ما يتم تكوين سلسلة طويلة من الـ BRNN لصعوبتها الكبيرة

--*-*-*-*-*-*-*-*-*-*-*-*-*-*

نهاية الاسبوع الأول