

# 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

## درس 5: الشبكات العصبية المتكررة RNN

## الأسبوع الثاني : المعالجة اللغوية الطبيعية NLP

استخدام المعالجة اللغوية الطبيعية Natural language processing والتي تختصر NLP مع التعلم العميق هو من التطبيقات الهامة , والتي تستخدم الكلمات كمصفوفات رأسية vectors , ويمكنك ان تقوم بتدريب الشبكة المتكررة RNN عليها , لتستخدم في العديد من التطبيقات , مثل تحليل المشاعر في الكلام (معرفة هل الكلام سلبي ام ايجابي) , التعرف على الأسماء , او الترجمة الأوتوماتيكية

\* \* \* \* \*

بعدما تعلمنا من الاسبوع الماضي كيفية التعامل مع شبكات الـ RNN و ادوات LSTM , GRU , فسنستخدم هذه الأشياء في تطبيقات هامة في التعامل مع الكلمات , ومنها ما يسمى تضمين الكلمات word embedding .

و يقصد بها قدرة الخوارزم علي معرفة الكلمات المناسبة لبعضها البعض , وما الذي يتم استخدامه حينما يتم اختيار كلمة معينة , واختيار كلمات الذكور و الاناث دون تداخل , وكذلك تجنب اي استخدامات غير مناسبة مع العرق او الدين او النوع .

وكل هذا يعتمد أولا علي فهم : عرض الكلمات word representation

و لفهم هذا المصطلح , علينا ان نتذكر فكرة مصفوفة الكلمات .

ذكرنا أن هناك قاموس بكل الكلمات , وليكن يحتوي علي 10 الاف كلمة , فإذا كانت كلمة man لها ترتيب رقم 5391 في القاموس , فستكون المصفوفة الخاصة بها هي :

Man  
(5391)

0
0
0
0
⋮
1
⋮
0
0

و هي ما تسمى مصفوفة الـ one-hot , لذا يقال عن هذه المصفوفة ان اسمها O5391 حيث حرف الـ O يشير الي one-hot و الرقم هو ترتيب رقم 1 في المصفوفة .

و نفس الحال في باقي المصفوفات :

Man	Woman	King	Queen	Apple	Orange
(5391)	(9853)	(4914)	(7157)	(456)	(6257)
$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \\ 0 \\ 0 \\ 0 \\ 0 \end{bmatrix}$	$\begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix}$

فكل مصفوفة فيهم يكون رقمها O ثم ترتيب رقم 1 فيها .

و المشكلة من هذه المصفوفات , ان الشبكة لن تتمكن من معرفة الكلمات المرتبطة ببعضها البعض عبر الارقام , فالشبكة لن تتمكن من معرفة ان كلمتي , king queen هما اقرب لبعضهما من كلمتي apple , queen

و حتي لو كان هناك جملة هي :

I love eating orange or -----

فلن يتمكن الخوارزم من الربط بين البرتقال و التفاح لاختيارها هنا . و ستكون فرصة اختيار كلمة apple مساوية لكلمة king و هي بالطبع غير سليمة .

ولأن هذه الفكرة غير مناسبة , فسنقوم باستخدام فكرة أخرى تسمى عرض الخصائص featurized representation

ونقوم الفكرة علي , عمل جدول ضخم الأعمدة فيه هي كل الكلمات الموجودة في القاموس , أما الصفوف فهي جميع الخصائص التي يمكن أن تتواجد في اي شئ في الدنيا

فقد يكون الجدول مثل هذا :

	Man 5391	Woman 9853	King 4914	Queen 7157	Apple 456	Orange 6257
Gender						
Royal						
Age						
Food						

Size						
Color						
Cost						

علي أن تكون الأعمدة فيها 10 الاف عمود , بينما الصفوف بها كل الخصائص الممكنة لأي كلمة , وقد تكون مثلا 300 خاصية .

و جميع النتائج تكون بين الـ 1 و الـ -1 , و تكون برقم يتناسب معها , فمثلا الكلمة الأولى :

	Man 5391	Woman 9853	King 4914	Queen 7157	Apple 456	Orange 6257
Gender	-1					
Royal	0.01					
Age	0.03					
Food	0.09					
Size	..					
Color	..					
Cost	..					

وباقى الكلمات قد تكون :

	Man 5391	Woman 9853	King 4914	Queen 7157	Apple 456	Orange 6257
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.7	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97
Size	..	..	..	..	..	..
Color	..	..	..	..	..	..
Cost	..	..	..	..	..	..

و بالتالي بدلا من أن يكون لـ vector الخاص بكلمة man هو :

0  
0  
0  
1  
0  
0

حيث رقم 1 ترتيبه 5391 , سيكون الفيكتور هو :

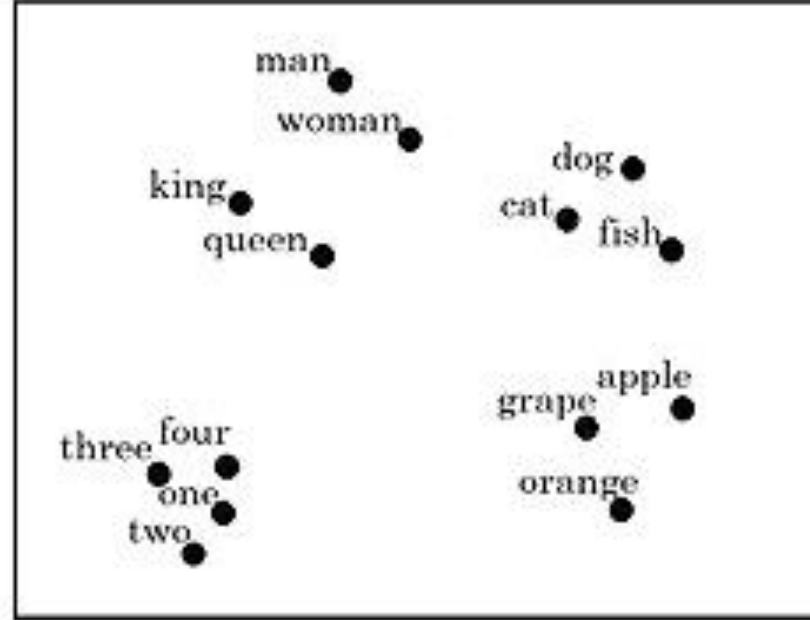
-1  
0.01  
0.03  
0.09  
..  
..

فيكون طول الفيكتور 300 رقم (بناء علي عدد الخصائص المستخدم ) . و تكون هذه الأرقام أكثر تعبيراً عن مضمون الكلمة نفسه , وتقوم بمساعدة الخوارزم علي الربط بين الكلمات المتشابهة . لأن كلا من الفيكتور الخاص بالكلمات المتشابهة سيكون له ارقام قريبة من بعضها البعض .

ووقتها يسمى الفيكتور بدلاً من O5391 سيكون e5391 و حرف الـ e دلالة علي embedding

ثم يمكن رسمها بيانياً كمخطط من بعدين هكذا :





ثم استخدام أحد وسائل التقسيم مثل t-SNE والتي تقوم بتقسيمها لمجموعات متقاربة من بعضها بحيث تجمع man , woman معا , ثم king , queen , ثم تجمع بين الاربعة معا , ثم تجمعهم كلهم مع باقي الحيوانات علي اعتبار انهم كائنات حية و هكذا



سنقوم الآن بتطبيق فكرة جدول الخصائص , لمعرفة كيفية استخدامها جيدا .

فإذا كان لدينا خوارزم لمعرفة اذا كانت الكلمة هي اسم شخص ام لا , وكان لدينا جملة في عينو التدريب :

Sally Johnson is an apple farmer

و قمنا بإعلام الخوارزم ان الكلمتين الأولى و الثانية قيمة  $y$  تساوي 1 (اسماء اشخاص) بينما قيمة  $y$  للباقي هي صفر ( ليسوا اشخاص )

فإذا قمنا في عينة الاختبار بإعطائه مثال وهو :

Robert Lin is an orange farmer

فيكون سهلا علي الخوارزم معرفة ان كلمتي orange & apple هما من نفس النوع و هما قريبتان من بعضهما , وبالتالي قد يتعرف علي ان اول كلمتين هما اسمين .

والميزة في طريقة جدول الخصائص , انه حتي الكلمات نادرة الاستخدام و التي قد لا تتواجد في عينة التدريب, ستكون معروفة .

فلو كان هناك جملة :

Robert Lin is a durian cultivator

الـ durian هي فاكهة غير منتشرة وموجودة أكثر في جنوب شرق اسيا , وكلمة cultivator هي بمعنى مزارع لكنها غير شائعة , فميزة فكرة جدول الخصائص انه سيتمكن من معرفة ان كلمة durian هي من نفس نوعية orange & apple حتى لو لم يرها من قبل في التدریب

و هناك ثلاث خطوات أساسية لتطبيق هذا التكنيك :

1. القيام بتعلم تضمين الكلمات word embedding learn و الذي سيكون عبر التعامل مع كميات هائل من الكلمات , والتي قد تصل لـ 100 مليار كلمة , و التعلم منها لمعرفة خصائص كل كلمة (وقد نستغني عن هذه الخطوة في حالة تواجد تعلم مسبق و تحميله و العمل به)
2. تطبيق هذا التعلم في عينة التدريب في كمية قليلة من الكلمات , وليكن 100 الف
3. عمل تعديل في التعلم في الخطوة الأولى , من نتيجة تطبيق الخطوة الثانية , و هذه غالبا تتم اذا ما كانت العينة في الخطوة الثانية كبيرة .

\* \* \* \* \*

كما رأينا في فكرة تضامن الكلمات word embedding و جدول الخصائص , ان لها عدد كبير من المميزات , و منها ما يسمى مسببات التشابه analogy reasoning وهي التي تختص بإيجاد المساحات المشتركة بين الكلمات بعضها البعض .

و لفهم مني التشابه , دعنا نلقي نظرة علي الجدول السابق شرحه :

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

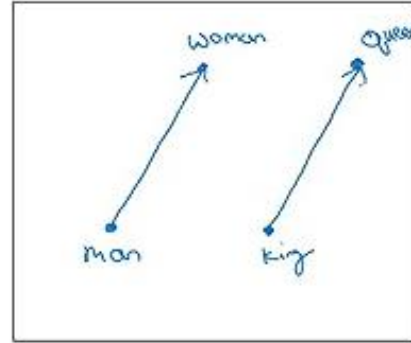
لو كان لدينا سؤال معين عن : علاقة الرجل بالمرأة و مثل علاقة الملك بمن ؟

$$e_{man} - e_{woman} \approx e_{king} - e_?$$

فالطريقة الممكن استخدامها , عبر طرح مصفوفتين الرجل و المرأة معا في البداية , فتكون مصفوفة الرجل - مصفوفة المرأة ستساوي تقريبا :

2  
0  
0  
0

ثم عمل تجربة طرح مصفوفة الملك من باقي المصفوفات , اختيار من يكون الفارق بينهما هونفس الفارق , وستكون الأرقام التقريبية مع الفارق بين الملك و الملكة وبالتالي لو تم رسم النقاط الأربعة في جراف من بعدين , و تم إيصال الفيكتور من man to woman , ثم من king to queen هتجد ان الفيكتور الواصل بينهم (رياضيا يمثل الطرح) , قيمة متشابهة هنا و هنا



و هذا يتم بسهولة في بايثون , بإيجاد دالة تأتي بالتشابه الأكبر كالتالي :

find word w , where max sim(  $\mathbf{e}_w$  ,  $\mathbf{e}_{king} - \mathbf{e}_{man} + \mathbf{e}_{woman}$  )

فقيمة الملك + الرجل - المرأة , ستكون قريبة جدا من الكلمة المطلوبة و هي الملكة .

ومعادلة التشابه تسمى معادلة cosine و تكون قيمتها الرياضية :

$$\text{Sim}(u, v) = \frac{u^T v}{\|u\| \|v\|}$$

و السبب في تسميتها بالـ  $\cos$  أن هذا القانون يحسب قيمة الـ  $\cos$  بين متجهين , كما أن لو كان المتجهين متطابقين (الزاوية 0) تكون الـ  $\cos$  تساوي 1 , لو كانا متعامدين (الزاوية 90 ) يكون الـ  $\cos$  يساوي 0 , ولو كانا متضادين (الزاوية 180 ) تقوم القيمة سالب 1 , وهو ما يتوافق مع معنى التشابه  $\cos$  similarity

و من الأمثلة المتنوعة لدالة التشابه :

Man – woman	=====	boy – girl
Ottawa – Canada	=====	Nairobi – Kenya
Big- bigger	=====	Tall – Taller
EN – Japan	=====	Ruble-Russia

\* \* \* \* \*

حينما يتم تدريب خوارزم خاص بالتحليل اللغوي , وبالتحديد في مجال تضمين الكلمات word embedding تظهر لنا ما يسمى مصفوفة التضمين embedding matrix , فما هي . .

كنا ذكرنا بالفعل أنه يجب تحضير مصفوفة ضخمة عن كل الكلمات المستخدمة , وان يكون لكل كلمة عدد كبير من الخصائص , مثل هذه :

	Man (5391)	Woman (9853)	King (4914)	Queen (7157)	Apple (456)	Orange (6257)
Gender	-1	1	-0.95	0.97	0.00	0.01
Royal	0.01	0.02	0.93	0.95	-0.01	0.00
Age	0.03	0.02	0.70	0.69	0.03	-0.02
Food	0.09	0.01	0.02	0.01	0.95	0.97

و عدد أعمدة هذه المصفوفة هي نفس عدد الكلمات في القاموس , وليكن 10 الاف , بينما عدد صفوفها هي عدد الخصائص التي تم تحديدها , وليكن 300 خاصية .

و هذه المصفوفة تسمى E بالكابيتال , وتكون عامة لكل الكلمات .





إذن كيف يتم التدريب ؟

إذا ما كان لدينا جملة مثل هذه , ونريد توقع الكلمة القادمة :

I    want    a    glass    of    orange    \_\_\_\_\_.

حيث لكل كلمة فيهم رقم محدد , فتكون :

I	want	a	glass	of	orange	_____.
4343	9665	1	3852	6163	6257	

فبيدأ الأمر , بإيجاد مصفوفات one-hot لكل كلمة علي حدة :

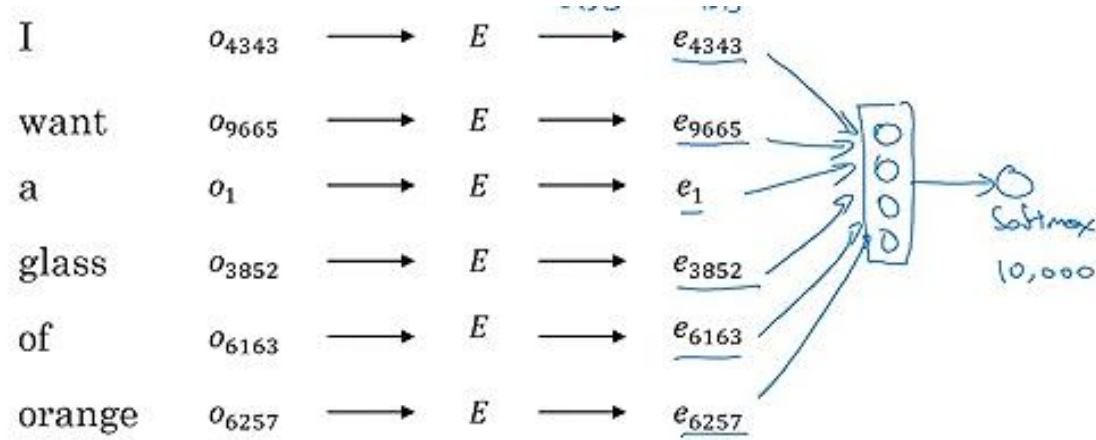
I	$o_{4343}$
want	$o_{9665}$
a	$o_1$
glass	$o_{3852}$
of	$o_{6163}$
orange	$o_{6257}$

ثم ضرب المصفوفة الكبرى E في فيكتور الـ one-hot لكلا منهم , لإيجاد الـ embedding matrix لكل كلمة هكذا :

I	$o_{4343}$	$\longrightarrow$	$E$	$\longrightarrow$	$e_{4343}$
want	$o_{9665}$	$\longrightarrow$	$E$	$\longrightarrow$	$e_{9665}$
a	$o_1$	$\longrightarrow$	$E$	$\longrightarrow$	$e_1$
glass	$o_{3852}$	$\longrightarrow$	$E$	$\longrightarrow$	$e_{3852}$
of	$o_{6163}$	$\longrightarrow$	$E$	$\longrightarrow$	$e_{6163}$
orange	$o_{6257}$	$\longrightarrow$	$E$	$\longrightarrow$	$e_{6257}$

بحيث تكون كل embedding matrix لكل كلمة e مصفوفة 300 صف في عمود واحد .

ثم نتناول هذه المصفوفات و ندخلها معا في شبكة عصبية (طبقة واحدة أو أكثر) و منها نخرج لـ softmax نقوم بالاختيار من وسط 10 الاف اختيار (عدد كلمات القاموس) هكذا :



و عبر بناء هذا الخوارزم , وتدريبه جيدا مع عشرات الالاف من النماذج , سيقوم بتحديد قيم لثبات الخاصة بالشبكة و سوفت ماكس , لتوقع اي كلمة قادمة .  
و غالبا ما يتم تحديد عدد معين من الكلمات ليتم علي اساسها توقع الكلمة التالية , فمثلا انت تقوم بتحديد أنني ساستخدم فقط 4 كلمات سابقة للكلمة الفارغة لاحدد هذه الكلمة , فتكون في هذا المثال :

a      glass      of      orange      \_\_\_\_\_  
1      3852      6163      6257

و أهمية هذا التحديد , هو التمكن من بناء الخوارزم بشكل صحيح , لأن عدد المعاملات  $b$  ,  $w$  التي ستكون في الشبكة العصبية معتمدة علي عدد المدخلات inputs , وغير منطقي ان تكون مفتوحة

فإذا كان لدينا 4 كلمات سيتم استخدامها , ولكل كلمة 300 رقم , فيكون عدد المدخلات هو 1200 , ويتم رصهم فوق بعضهم البعض , فيكون هناك 1200 ثيتا سيتم تدريبها في الشبكة المطلوبة .

و في حين هذا التصور هو الشكل الرسمي او المثالي , فهناك تصورات أخرى قد تكون أكثر بساطة , وفي نفس الوقت بنفس الفعالية الخاصة بها .

فمثلا لو كان لدي جملة مثل :

I want a glass of orange ----- to go along with my cereal

وبفرض كانت كلمة juice هي المجهولة و المطلوب استنتاجها , فالشكل الاصلى ان نختار الكلمات الأربعة الأخيرة هكذا :

a glass of orange -----

إلا أنه يمكن عمل المزيد من الدقة (والوقت) باختيار كلمات اربع قبل و كلمات اربع بعد مثل :

a glass of orange ----- to go along with

و قم يتم اختيار كلمة واحدة فقط قبلها :

orange 

وقد يتم اختيار كلمة محددة قبلها , ليس بالضرورة السابقة لها , مثل كلمة glass

\* \*\*\*\* \* glass \*\* \*\*\*\*\* -----

وهو التكنيك المسمى skip gram والذي سنتكلم عنه لاحقا

\* \* \* \* \*

و الفكرة هنا تسمى word2vec و القائمة علي المبدأ الذي سبق ذكره وهو : skip gram

و تقوم الفكرة علي اختيار الكلمة التي سنستدل بها علي الكلمة الناقصة , و تسمى كلمة المضمون context و الكلمة الناقصة نفسها و نسميها كلمة الهدف target .

ثم نتناول الجملة و ليكن :

I want a glass of orange juice to go alone with my cereal

لاحظ ان كلمات المضمون و الهدف قد تختلف , فقد تكون قريبة من بعضها او تالية بعيدة عن بعضها البعض , لكن غالبا ما تكون في نطاق 10 كلمات , اذ نادرا ما تؤثر كلمات في مسافة ابعد من هذا

فيتم إيجاد مصفوفة المضمون embedding matrix الخاصة بكل كلمة , وهي التي لها رمز e مثل e6257 , وهي التي تأتي من ضرب مصفوفة الخصائص E في مصفوفة الـ one-hot .

ثم يتم حساب قيمة  $P(t | c)$  أي حساب احتمالية كلمة الهدف t اعتمادا علي معلومية كلمة المضمون c , وتكون بالقانون :

$$p(t|c) = \frac{e^{\theta_t^T e_c}}{\sum_{j=1}^{10,000} e^{\theta_j^T e_c}}$$

نبدأ بشرحه جزء جزء .

أولا علي اليسار القيمة :

$$p(t|c)$$

هي احتمالية الحصول علي الكلمة الهدف  $t$  بناء علي معلومية الكلمة المضمون  $c$

بالنسبة للبسط :

$$e^{\theta_t^T e_c}$$



فحرف e الاسفل هو الـ exponential اي 2.71 , و الأس الخاص به هو ثيتا الخاصة بالكلمة الهدف t بعد عمل مقلوب لها T مضروبة في ec وهي مصفوفة المضمون embedding matrix الخاصة بالكلمة المضمون c

بينما المقام

$$\sum_{j=1}^{10,000} e^{\theta_j^T e_c}$$

هو مجموع كل expo علي حدة للكلمات العشر الاف , حيث كل واحدة expo لها أس هو الثيتا الخاصة بها مقلوبة T مضروبة في ec وهي الكلمة المضمون

بالنسبة للثيتا t هي الأوزان الخاصة بكل كلمة في الجملة المعطاة للخوارزم

والفكرة تقوم علي قسمة بسط علي مقام حيث البسط يوضح مدي ارتباط كلمة المضمون بكلمة الهدف , بينما المقام هو مجموع كل العلاقات بين كلمة المضمون و باقي الكلمات في اللغة .

أما معادلة الخطأ المطلوب تقليلها هي :

$$\mathcal{L}(\hat{y}, y) = - \sum_{i=1}^{10000} y_i \log \hat{y}_i$$

مع ملاحظة ان قيمة  $y$  هنا هي الناتج , وهي مصفوفة الـ one-hot الخاصة بالكلمة الناتجة .

وهذا التكنيك يسمى skip gram , لأنه قد يتخطى عدد من الكلمات قبل اختيار الكلمة المناسبة , فليس شرط ان تكون تالية او سابقة لها

و من عيوب هذه الطريقة الوقت المطلوب لها , حيث ان المقام يحتاج لعشر الاف عملية حسابية , و تكون المشكلة اعظم اذا ما كان القاموس يحتوي علي عدد ضخم و ليكن مليون كلمة .

فيكون أحد الحلول استخدام التصنيف الهرمي hierarchical softmax classifier , وهو الذي يجعل الأمر اسرع .

ويبدأ عبر جعل الكلمات العشر الاف في قسمين كل قسم 5 الاف , فيقوم الخوارزم بفحص هل هي في القسم الاول ام الثاني , فإذا كان في الاول يقوم بفحص هل هي من الـ 2500 الاول ام الثانية , وهكذا حتي نصل للكلمة المطلوبة , في عد اقل من العمليات , وبوقت اسرع .



كنا قد رأينا في التكنيك السابق skip gram أنه فعال في توقع الكلمات , لكن يعيبه البطء الشديد , سنتعرف الآن علي تكنيك العينة السالبة negative sample و الذي يعطي نفس الكفاءة و في وقت أقل .

و تقوم الفكرة علي عمل جدول , يضم ثلاث عواميد , الأول هي كلمة المضمون context التي سنتعامل معها , الثاني ستكون كلمات سنقوم باختبار مدي ارتباطها بكلمة المضمون , والثالث هو نتيجة هذا الارتباط

وتكون الكلمة الأولى في الـ word هي الكلمة المرتبطة بكلمة المحتوي , لذا تكون قيمة الـ target تساوي 1 , بينما يتم اختيار عدد اخر من الكلمات بشكل عشوائي من القاموس , بحيث تكون قيمة الـ target تساوي 0 , هكذا

Context	Target	y
Orange	juice	1
Orange	king	0
Orange	book	0
Orange	the	0
Orange	of	0

و يكون عدد الكلمات العشوائية هو k و الذي يكون من 5 الي 20 كلمة في العينات ذات العدد القليل , بينما من 2 – 5 في الأعداد الكبيرة .

ثم يتم إعطاء هذه البيانات للخوارزم ليتدرب عليها , علي اعتبار ان أول عمودين هما المدخلات x بينما العمود الثالث هو المخرج y

ويسمي هذا التكنيك العينة السالبة , لأننا نقوم باختيار كلمات غير مرتبطة بكلمة المضمون وذات ارتباط سلبي بها

و هنا تكون المعادلة الخاصة بالاحتمالية :

$$P(y=1 \mid \mathbf{c}, \mathbf{t}) = \sigma(\Theta_t^T \mathbf{e}_c)$$

وهو الذي يعني أننا نأتي بثباتي ، ويتم عمل مقلوب لها ، ثم ضربها في **ec** ثم عمل سيجمويد لكل هذا .

و الميزة في هذا التكنيك عن السابق , ان عدد البيانات التي يتم تدريب الخوارزم عليها يكون عدد قليل وهو  $k+1$  (اي عدد العناصر الغير سليمة + 1 وهو العنصر السليم )

فبدلاً من جعل الخوارزم يتدرب على 10 آلاف كلمة بوقت رهيب ، اجعله يتدرب على كلمات عشوائية غير مرتبطة به ، هو ما يجعل الوقت اقل

و لاختيار الكلمات العشوائية , يتم استخدام صيغة معقدة نوعا لكن مؤثرة وهى :

$$P(w_i) = \frac{f(w_i)^{3/4}}{\sum_{j=1}^{10,000} f(w_j)^{3/4}}$$

\* \* \* \* \*

نتعلم الان تكنيك مختلف في نفس مجال NLP و هوا ما يسمى تكنيك GloVe , وعلي الرغم من ان انشاره و شهرته اقل من تكنيك word2vec إلا انه ابسط و له استخدامات عديدة

و اسمه GloVe هو اختصار :

Global Vector for word representation

و تبدأ الفكرة بالتعرف علي مصطلح  $X_{ij}$

و تعريفه هو عدد المرات التي ظهرت فيها الكلمة  $i$  في نطاق الكلمة  $j$  , وكلمة في نطاق المقصود بها في خلال العدد المحدد من الكلمات و ليكن 10 كلمات .

كأن الكلمة  $i$  تعبر عن الهدف  $t$  بينما الكلمة  $j$  هي المضمون  $c$  , وهو ما يجعل هذا التكنيك مبنيعلي نفس اساس التكنيك السابق , وهو إيجاد العلاقة بين كلمتي المضمون و الهدف .

و غالبا ما يكون  $X_{ij}$  هي نفسها  $X_{ji}$  , حيث انه لم تم استبدال مكاني كلمتي الهدف و المضمون ستكون القيمة نفسها .

و يكون الهدف هو تقليل القيمة بين ثيتا  $i$  (مقلوب) مضروبة في  $e_j$  , وبين لوج قيمة  $X_{ij}$  , لان اقتراب القيمتين من بعضهما يعني ان الكلمتين مرتبطتان ببعضيهما

$$(\theta_i^T e_j - \log X_{ij})^2$$

ويتم عمل الصيغة الكاملة :

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})(\theta_i^T e_j - \log X_{ij})^2$$

مع اعتبار ان علامتي sum في البداية لمسح جميع الكلمات و علاقاتها بباقي كلمات القاموس .

أما الجزء الخاص بالدالة  $f(X_{ij})$  فله مهمتين .

الأولي انه سيساوي صفر في حالة كانت  $X_{ij}$  تساوي صفر , وقيمة  $X_{ij}$  تساوي صفر معانها عدم وجود اي كلمة هدف في نطاق كلمة المضمون .

و أهمية جعل الـ  $f(X_{ij})$  تساوي صفر لتجنب المشكلة الرياضية الناتجة من حساب log للصفر , والذي سيساوي سالب مالانهاية , فوجود الصفر في البداية يجعل القيمة بصفر و يجنبنا المشكلة الرياضية .

كما ان الدور الثاني بها متعلق بتحديد اوزان للكلمات حسب اهميتها و انتشارها في اللغة الانجليزية , فكلمات منتشرة مثل ( he , she , a , of , and ) يكون لها اوزان كبيرة , بينما الكلمات النادرة يكون لها اوزان اقل , وهو ما يجعل الخوارزم اكثر دقة في التعامل معها

أخيرا يتم وضع معامل bias لكلا من  $i, j$  والذي سيكون هكذا:

$$\text{minimize } \sum_{i=1}^{10,000} \sum_{j=1}^{10,000} f(X_{ij})(\theta_i^T e_j + b_i - b'_j - \log X_{ij})^2$$

\* \* \* \* \*



نتعرف الان علي تطبيقات NLP وهو ما يسمي تصنيف الانطباع sentiment classification

كلمة sentiment تعني مشاعر , و يقصد بها هنا : التعرف علي انطباعات و آراء كاتب التعليق , من كتابته .

فلو كان هناك الاف التعليقات علي منتج معين او فيلم ما , فنريد عمل ترجمة لها لارقام , لمعرفة كاتب هذا التعليق يقصد تقييمه بكم من 5 مثلا . او علي الأقل هل هو مؤيد لو رافض له .

كما ان له استخدامات سياسية , في تناول اعداد هائلة من التويتات , وتحليلها لمعرفة نسبة كم منها مؤيد لترامب و كم منها رافض له .

و هنا امثله لها :

$x$	$y$
The dessert is excellent.	★★★★☆
Service was quite slow.	★★☆☆☆
Good for a quick meal, but nothing special.	★★★☆☆
Completely lacking in good taste, good service, and good ambience.	★☆☆☆☆

و أكبر تحدي في التعامل معها , أننا لا نملك كمية كبيرة من بيانات التدريب لتدريب الخوارزم عليها , فيندر ان تجد اكثر من 10 الاف تعليق عن كتاب او مطعم معين , و هذا رقم متواضع يصعب الاعتماد عليه .

و نبدأ بتناول التعليق المكتوب كمدخل  $x$  بينما عدد النجوم هو المخرج  $y$  هكذا :

The	dessert	is	excellent	★★★★☆
8928	2468	4694	3180	

و تكون الخطوة الأولى بحساب صفوف المضمون embedding matrix لكل كلمة (حاصل ضرب الـ  $E$  في مصفوفة one-hot )

فيكون لدينا 4 مصفوفات (عدد 4 معتمد علي عدد الكلمات ) و كل واحدة عمود واحد مع 300 صف (عدد الخصائص)

ثم نقوم بإيجاد مجموع هذه الارقام , او متوسطها , وإدخالها في شبكة عصبية , لتخرج لنا في سوفت ماكس بقيمة من 5 قيم



وكان هذا النظام مع تدريبه , يقوم بتعليم الخوارزم ان انتشار كلمات مثل كذا و كذا , يعني نتيجة كذا , فيتمكن من قراءتها بسهولة و تقييمها و من أهم عيوب هذه الطريقة , انه يمكن خداعها بسهولة , فمثلا جملة مثل :

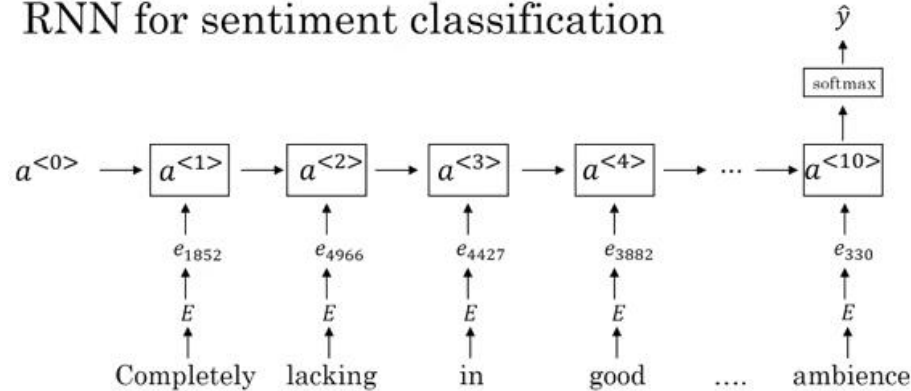
Completely lacking in good taste,  
good service, and good ambience.

من الواضح ان التعليق سلبي تماما , إلا أن انتشار كلمة good , taste , service سيجعل الخوارزم يميل في كفة التقييم العالي , مع أن كلمة lacking قد قلبت المعني .

فمشكلة ايجاد المجموع او المتوسط انه يفحص عدد وجود كلمة ما , بغض النظر عن معناها الإجمالي .

فيكون الحل هنا هو استخدام RNN من نوعية many-to-one حينما نقوم بتغذيتها بكل الكلمات , وتقوم بإخراج واحد مثل هذا :

## RNN for sentiment classification



وميزة الـ RNN هنا ان كل الكلمات تتداخل معا في اتخاذ القرار , و هنا سيتعلم الخوارزم ان not good شئ سلبي , وأن lack of good taste نقاب المعنى , وذلك لأن طبيعة الـ RNN تجعل المخرج يعتمد على جميع المدخل , وليس على متوسط مثل التكنيك السابق .

و الميزة هنا انه حتي لو تم استخدام كلمة سلبية غير منتشرة مثلا (absence of good taste) حتي لو لم يتم استخدامها اثناء التدريب , فاقتراب معناها من معنى lack سيجعل الخوارزم يفهم المقصود .

و الخوارزم سيفهم انها قريبة منها اعتمادا على مصفوفة المضمون embedding matrix و التي تشتمل على الخصائص .

\* \* \* \* \*

الموضوع الأخير هذا الأسبوع , متعلق بعلاج الأخطاء العنصرية في المعالجة اللغوية .

و المقصود بها تجنب اي تصنيف او اختيار , بما يدعم فكرة التفرقة العنصرية , او الصورة النمطية الخاصة بالنوع , او العرق , او الدين , او التوجه الجنسي .

فحديثنا الان ليس عن خطأ الاختيار انه غير صحيح , لكن عن الاخلاقيات الواجب مراعاتها للآلة حينما تقوم بصياغة الجمل .

كمثال بسيط , لو قلنا أن :

رجل == امرأة يساوي ملك == ملكة

فإذا ذكرنا أن : رجل == مبرمج يساوي امرأة == ؟

حينما قامو بتدريب الخوارزم اكتشفو نتائج مخيفة , ان اغلب الاستنتاجات كانت لاشياء مثل : ربة منزل , عاطلة .

كذلك . .

أب == طبيب يساوي أم == ؟

كانت النتيجة : ممرضة . .

و المشكلة هنا ليست مشكلة تكنيكال , لكن أن البيانات التي يتم تدريب الخوارزم عليها هو فيه الكثير من العنصرية , والأزمة أن الكلام الحقيقي الموجود في كل مكان فيه الكثير من العنصرية , ولن نستطيع اختراع بيانات للتدريب من الهواء .

و سنتناول الان فقط مشكلة العنصرية مع النوع : ذكر و أنثي , ولكن فكرة الحل كون نفسها مع اي نوع آخر

و لعلاج المشكلة , نبدأ بأولا باستعراض أماكن عدد من الكلمات المستخدمة و التي تداخل بشكل ما في المشكلة مثل :

doctor , babysitter , boy , girl , he , she , grandfather , grandmother

و نري أن هذه الكلمات قد يتم وضعها هكذا , بناء علي التدريب من كمية كبيرة من البيانات , والتي قد يكون بها بعض العنصرية :



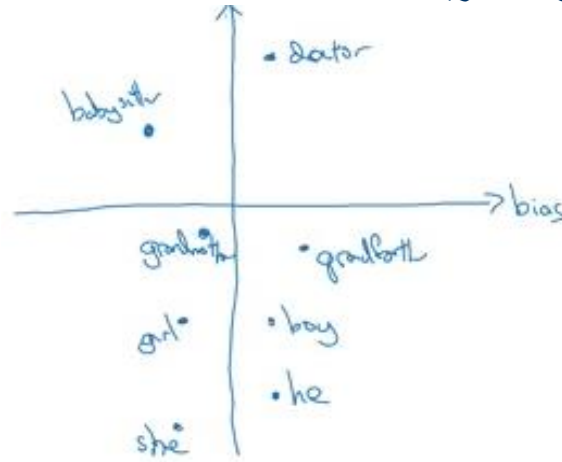
A handwritten list of words, each followed by a small dot, arranged in two columns. The words are: doctor, babysitter, grandfather, grandmother, boy, girl, he, and she. The words are written in a cursive, handwritten style.

و مصدر المشكلة رياضيا , انه مع مقارنة he—she تكون المقارنة قريبة جدا من doctor—babysitter وهو ما يجعل الخوارزم يتخذ قرارات عنصرية , ونريد ان نقوم بتغيير هذا الأمر

و يكون حل المشكلة في خطوات ثلاث .

الأولي : في تحديد اتجاه الإنحراف (العنصرية)

فيكون المحور الأفقي هنا هو محور الإنحراف , لأنه يتحول من اليسار إلي اليمين لما يتماشي مع اختلاف الجنس , بينما المحور الرأسي هو اللا انحراف , لأن الفارق بين كلمتي she , girl ليس متعلق بالإنحراف او العنصرية



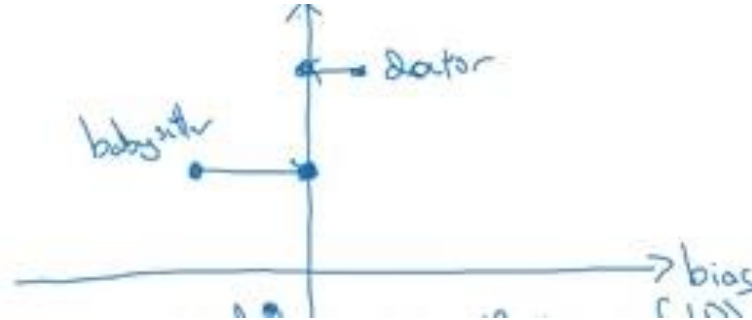
تأتي الخطوة الثانية : في التفريق بين نوعين من الكلمات .

الكلمات الي يكون تعريفها معتمد أساسا علي الفارق في النوع , مثل grandmother , فهي في الأساس مبنية علي أنها أنثي , كذلك she و كذلك : girl و wife

وكلمات أخرى تعريفها لا يعتمد علي النوع , مثل nurse , babysitter , programmer , doctor , فهذه الكلمات ليس لها علاقة بالنوع من الأساس .

فنمسك كلمات النوع الثاني , و نقوم بتحريكها حتي تسقط علي المحور الرأسي , وذلك حتي يكون ليس لها تأثير مع اختلاف النوع

فتتحرك كلمة doctor من اليمين لليسار و كلمة babysitter من اليسار لليمين هكذا

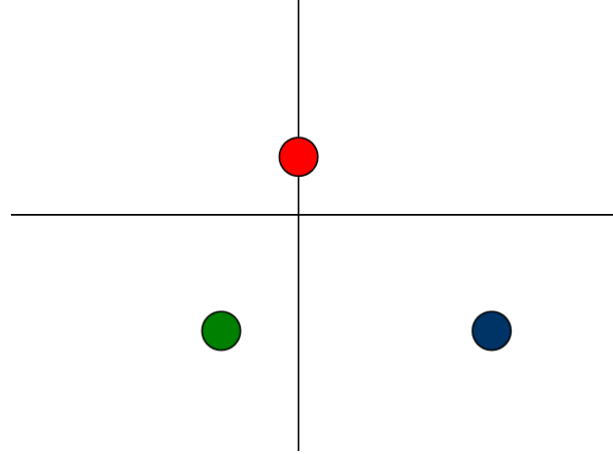


ثم تأتي الخطوة الثالثة وهي تناول الكلمات من النوع الأول (المعتمدة في تعريفها علي النوع ) وتطبيق: تساوي النوعين في المسافة .



لو كانت كلمة he هي الزرقاء , و كلمة she هي الخضراء , وتم بالفعل تغيير مكان كلمة babysitter الحمراء من اليسار حتي تقع علي المحور الرأسي (لتجنب تصنيفها) .

فحتي مع هذا الأمر , لازالت كلمة she اقرب لكلمة babysitter من كلمة he و هو ما يجعل اختيارها اقرب للأنثي منها الي الذكر , فالمشكلة قائمة .  
و لن نستطيع جعل كلا من she , he علي المحور الراسي لان هذا سيحدث خلل في التعامل اللغوي من الاساس , اذ ان التفريق في النوع بينهم شئ هام



فيكون الحل , ضبط مكانيهما , بحيث يكون كلا منهما علي نفس المسافة من المحور الرأسي , مع الإبقاء أن هذا يمين و هذا يسار , وهو ما يجعل المسافة بين كلا من she , he ( الكلمات المعتمدة علي اختلاف النوع ) هي نفسها من كلمة لا تعتمد علي اختلاف النوع

