

# 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

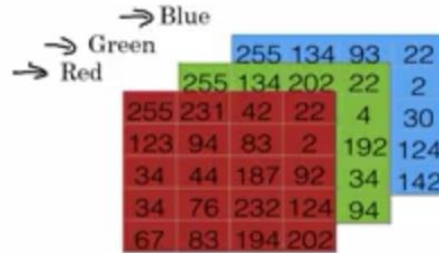
# درس 1: التعليم العميق و الشبكات العصبية

## الأسبوع الثاني : أساسيات الشبكات العصبية

للتعرف علي أساسيات الـ NN علينا التحدث بشكل مستفيض عن الكلاسيكاشن (التصنيف)

التصنيف هي عملية فرز المدخلات الي احد صنفين , غالبا ما تكون 1 و 0 , فمثلا اذا كنا نريد تصنيف الصور , هل بها قطة ام لا , فوجود قطة بها تعني 1 و عدم وجودها تعني 0

و علينا ان نتكلم كيف يتم تخزين الصور , فصورة مثل هذه , تخزن بهذه الطريقة في الحاسب



فالجهاز يقوم بعمل 3 مصفوفات لقيمة كل بيكسيل فيها , كم فيها من الالوان الثلاثة الرئيسية RGB  
 فلو كانت الصورة 100 بيكسيل في 100 بيكسيل , سيكون عندنا 3 مصفوفات ابعاد كل واحدة 100 في 100

بعدها لازم يتم تحويل كل البيانات ديه لفكتور (مصفوفة من بعد واحد) هيكون باسم  $x$  , فهيتم نقل ارقام مصفوفة الاحمر بالترتيب (255 و 231 و 42 و 22 و 123 و ...) بحيث تكون 10 الاف رقم (100 في 100) بعدها يتم نقل ارقام الاخضر , بعدها الازرق , يعني هيكون الفكتور اكس عبارة عن 30 الف صف في عمود واحد

$$x = \begin{bmatrix} 255 \\ 231 \\ 42 \\ \vdots \\ 255 \\ 134 \\ 202 \\ \vdots \\ 255 \\ 134 \\ 93 \\ \vdots \end{bmatrix} \begin{array}{l} \text{red} \\ \text{green} \\ \text{blue} \end{array}$$

وبالتالي هتكون الاكسات مصفوفة 30 الف رقم , والواي هيكون 1 او 0 حسب فيها قطة او لا

وقبل ما نكمل لازم نحدد بعض المصطلحات الهامة :

- مجموع عدد الارقام داخل مصفوفة  $X$  هيكون اسمه  $n$  , واللي هو حاصل ضرب الطول في العرض في 3 (حاليا 30 الف)
- عدد الصور المستخدم للتدريب (training sample) هيكون اسمه  $m$  (سمول)
- وبالتالي انا عندي  $x_1, x_2, x_3$  دلالة عن مصفوفات كل صورة فيهم لغاية  $x_m$
- فيه كمان  $y_1, y_2, y_m$  للغاية  $y_m$  للدلالة علي الناتج لكل صورة هل 1 يعني فيه قطة ولا 0 يعني مفيش
- عشان اجمع كل البيانات مع بعض هيكون فيه مارتكس اسمها  $X$  بالكابيتال , وهي تجمع كل الاكسات مع بعض , بحيث ارقام اول صورة تكون اول عمود , وتاني صورة تاني عمود و هكذا , فهيكون عدد الصفوف هو عدد الارقام لكل صورة ( $n$ ) وعدد العواميد هو عدد الصور للتدريب ( $m$ )

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & & | \end{bmatrix}$$

$\xleftarrow{m}$   $\xrightarrow{n_x}$

اما واي فيبسطة هنرص الارقام جنب بعض , بحيث المصفوفة المجمعـة لـواي  $\gamma$  كـابـيـتـال , هـتـكـون صـف وـاحـد فـي  $m$  عـمـود يـعـنـي مـثـلا  $[1\ 0\ 1\ 1\ 0\ 0\ 0\ 1\ 1]$

$$\gamma = [\gamma^{(1)} \ \gamma^{(2)} \ \dots \ \gamma^{(m)}]$$

\*\_\*\*

## ماذا عن التصنيف ؟

عايزين نعمل لوجيستيك ريجريشن (تصنيف) لمجموعة من الصور (الانبت اكس) عشان يكون الناتج عندنا هل فيها قطعة (واي تساوي 1) او مفيهاش قطعة (واي تساوي 0)

يعني الفورمة هتكون كدة :

$$\text{Given } x, \hat{y} = P(y = 1|x), \text{ where } 0 \leq \hat{y} \leq 1$$

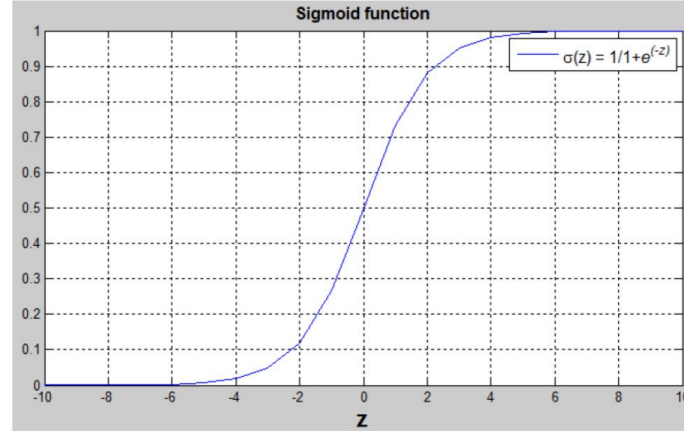
واي هات اللي هي قيمة واي المتوقعة , يعني قيمة الواي اللي الخوارزم هيتنبأ بيها , حيث واي هات هتكون بين الصفر و الواحد , صفر في حالة مفيش اي قطعة , 1 في حالة كل الصور قطط

لاحظ اننا قبل كدة كنا بنقول علي القيمة المتوقعة رمز اتش بدل واي هات

هنقول علي معاملات الفيتشرز (العناصر) اللي هتكون الثيتات رمز  $w$  عشان ترمز للويتس (اللي هي قيمة ثيتات) هتكون واي هات (المتوقعة) تساوي :

$$\hat{y} = \sigma(w^T x + b)$$

رمز السيجما ده معناها ان دالة السيجمويد , اللي دايما بنستخدمها في الكلاسيفيكاشن , عشان بتجيب ناتج من صفر لواحد , واللي بيكون شكلها كدة



قيمة الزى هي الدالة  $(Wt x + b)$  , والسيجمويد هتكون 1 علي 1 زائد اكسبونينشيل ناقص الزى

$$s = \sigma(w^T x + b) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

وبالتالي لما تكون الزى قليلة سالب , هتكون الإي قيمة كبيرة , ووحد علي كبير يبقى صفر , و العكس , لما تكون إي قيمة كبيرة , هتكون الإي صغيرة جدا , فهتساوي واحد , ولما تكون زي بصفر , الإي هتكون بواحد يعني نص

تعالى نتكلم عن صيغة الدبليو و البى

الدبليو هي قيمة ثيتات و المفروض ان القيمة ديه كلها  $w^T x + b$  تساوي , ثيتا زيرو في اكس زيرو + ثيتا 1 في اكس 1 + ثيتا 2 في اكس 2 وهكذا

ثيتا زيرو هي نفسها قيمة  $b$  و اكس زيرو بواحد , فثيتا زيرو في اكس زيرو تساوي البى , والبى نفسها  $b$  اختصار لكلمة bias و هو القيمة اللي كنا بنحطها في التنبؤ عشان الانحراف

بينما الدبليو هي مصفوفة عمود واحد و فيها صفوف بعدد الثيتات , فلما اعمل ترانسبوز ليها , هتكون صف واحد و فيها عوامد بعدد الثيتات , ولما اضربها في الاكس اللي هي عمود واحد و فيها صفوف بعدد الاكسات , تعطيني قيمة ثيتا 1 في اكس 1 + ثيتا 2 في اكس 2 وهكذا

$$x_0 = 1, \quad x \in \mathbb{R}^{n_x+1}$$
$$\hat{y} = \sigma(\Theta^T x)$$
$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_{n_x} \end{bmatrix} \left\{ \begin{array}{l} b \leftarrow \\ w \leftarrow \end{array} \right.$$



## معادلة الخطأ للتصنيف LR Cost function for LR

الأول عايزين نوضح كام حاجة مهمة  
دايما السوبر سكريببت (الرقم اللي بيكون فوق الرمز) , بيوضح رقم العينة اللي باتكلم عنها , بينما السب سكريببت (الرقم اللي تحت) بيوضح باتكلم عن انه عنص  
داخل نفس العينة

يعني لما يكون عندي إكس فوق فيه 3 , وتحت 5 , يبقى باتكلم علي العنصر الخامس في العينة الثالثة , يعني إم تساوي 3 , وإن تساوي 5

إذن لو بصينا في المعادلة ديه

$$\hat{y}^{(i)} = \sigma(w^T x^{(i)} + b), \text{ where } \sigma(z^{(i)}) = \frac{1}{1 + e^{-z^{(i)}}}$$

معناها ان (واي هات أي) يعني قيمة واي المتوقعة للعينة أي (العينة رقم 5 مثلا) هي سيجمويدي ديليو ترانسبوز (ديلوي ملهاس علاقة بالأي لانها ثابتة) , مضروبة في  
إكس أي (يعني كل اكسات العينة الخامسة) , + بي اللي هي ثابتة

و بالتالي مصطلح زي أي , المقصود بيه قيمة ديليو ترانسبوز, مضروبة في إكس أي, + بي

كمان ممكن نقول إن فيه واي هات لكل واحد من العينة , واللي احنا نفسنا تكون بتساوي واي الحقيقية لنفس العينة , والمعطيات هي اكسات ووايات لكل واحد من العين  
كل أكس فيها هي في الحقيقة مصفوفة اكسات العناصر

Given  $\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$ , we want  $\hat{y}^{(i)} \approx y^{(i)}$

وبالتالي معادلة الخطأ ستكون كالتالي , رمز إل اللي يشير لـ lose سيكون نص فرق قيمة واي هات وواي الحقيقية مربع  
فلو كانت واي المتوقعة زي واي الحقيقية بالضبط (التوقع سليم) , ستكون إل تساوي صفر  
وكل ما تبعد واي المتوقعة عن واي الحقيقية , تزداد قيمة الخطأ

لاحظ إن واي الحقيقية دايما برقم صحيح صفر او واحد (ممكن يكون فيه 2 و 3 لو كان تصنيف لاكثر من شئ) , بينما واي المتوقعة ممكن تكون رقم عشري

$$L(\hat{y}^{(i)}, y^{(i)}) = \frac{1}{2}(\hat{y}^{(i)} - y^{(i)})^2$$

بس المعادلة ديه محدش بيستخدمها حاليا لانها بتعمل صعوبة كبيرة في الوصول للقيمة الاوبتيمم للثينات , فهنستخدم المعادلة اللوغاريتمية ديه :

$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

وقبل ما نفسرها , هنتكلم عن خطوة سابقة عشان نفهمها :

في حالة كانت واي (القيمة الحقيقية) تساوي 1 , فالاحتمالية هتساوي واي هات (القيمة المتوقعة) , عشان لو كانت واي هات تساوي كمان , فالاحتمالية تساوي 1 يعني  
بتساويها (ترو) , ولو كانت هات تساوي صفر فالاحتمالية صفر كمان لانها مش بتساويها (فولس)

وفي حالة واي تساوي صفر , فالاحتمالية تساوي 1 ناقص واي هات , عشان لو كانت هات المتوقعة تساوي 1 فالاحتمالية صفر (فولس) لانها مش بتساوي واي , ولم كانت هات بصفر , فالاحتمالية بواحد (ترو) لانها بتساويها

$$\begin{aligned}\text{If } y = 1: & \quad p(y|x) = \hat{y} \\ \text{If } y = 0: & \quad p(y|x) = 1 - \hat{y}\end{aligned}$$

و ديه ممكن نصيغها بالمعادلة ديه :

$$p(y|x) = \hat{y}^y (1-\hat{y})^{(1-y)}$$

واي هات اس واي , مضروبة في 1 ناقص واي هات اس 1 ناقص واي

في حالة واي تساوي 1 , هتكون واي هات اس 1 , و 1 ناقص واي هات هتختفي

وفي حالة واي تساوي صفر هيختفي الجزء الاولي , وهيتبقى 1 - واي هات

ولو عملت لوج للطريفيين , هتلاقي ان الاس هيبقي جنب اللوج بحيث تعمل لنا المعادلة اللي بنتعامل معاها:

$$L(\hat{y}^{(i)}, y^{(i)}) = -(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))$$

وديه بتغطي الحالتين للواي

في حالة واي تساوي 1 , هيختفي التيرم الثاني , وهتكون معادلة الخطأ هي سالب لوج الواي المتوقعة , و ده هيجعل قيمة الخطأ صفر لو كانت واي المتوقعة بواحد , وهيجعل الخطا كبير لو كانت واي المتوقعة ابعد عن 1

$$\text{If } y^{(i)} = 1: L(\hat{y}^{(i)}, y^{(i)}) = -\log(\hat{y}^{(i)})$$

وفي حالة واي تساوي صفر , هيختفي الجزء الاول , وهتكون سالب لوج 1 - القيمة المتوقعة , ولو كانت واي المتوقعة تساوي صفر , هتكون لوج 1 , يعني الخطا صفر , ولو كانت ابعد عن الصفر قيمة الخطا بتزداد

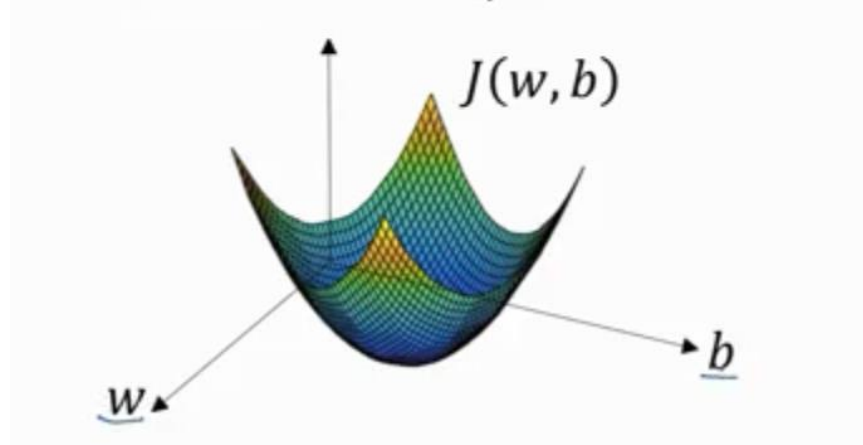
$$\text{If } y^{(i)} = 0: L(\hat{y}^{(i)}, y^{(i)}) = -\log(1 - \hat{y}^{(i)})$$

كل اللي فات ده كانت قيمة الخطأ لعينة واحدة من العينات , فمادنا عن الخطأ الكلي , ده اللي هيكو اسم معادلة الكوست , اللي بنرمز ليه برمز جي J , واللي هيكور مجموع قيم اخطاء كل عينة , مقسوم علي عددها إم يعني هيكو كدة

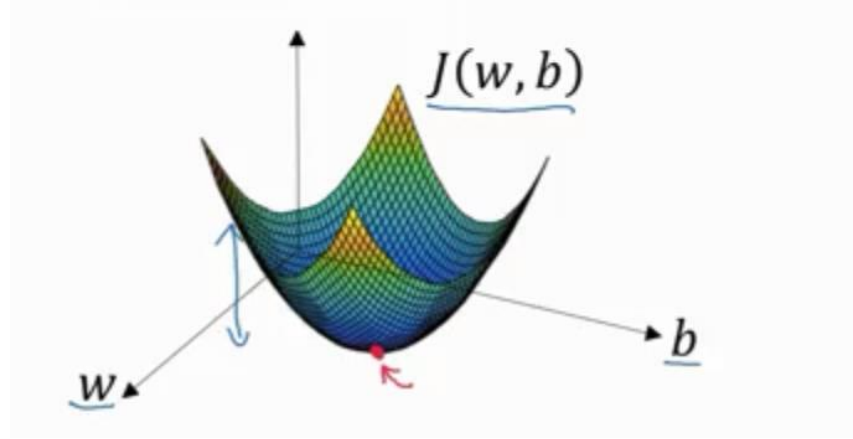
$$J(w, b) = \frac{1}{m} \sum_{i=1}^m L(\hat{y}^{(i)}, y^{(i)}) = -\frac{1}{m} \sum_{i=1}^m [(y^{(i)} \log(\hat{y}^{(i)}) + (1 - y^{(i)}) \log(1 - \hat{y}^{(i)}))]$$

## الانحدار التدريجي Gradient Descent

فاكرينالانحدار التدريجي من محاضرات التنبؤ , واللي بيكون شكله كدة



قلنا من شوية ان الـ  $w$  هي الاوزان يعني مصفوفة الثبتات , والـ  $b$  هي قيمة الثبتا صفر , وقيمة الـ  $J$  هي قيمة الخطأ , أي إن قيمة كلا من الدبليو و البي بتحدد قيم الـ  $J$  , وبالتالي احنا بنبحث عن قيم دبليو و بي اللي بتعمل قيمة جي اقل ما يكون



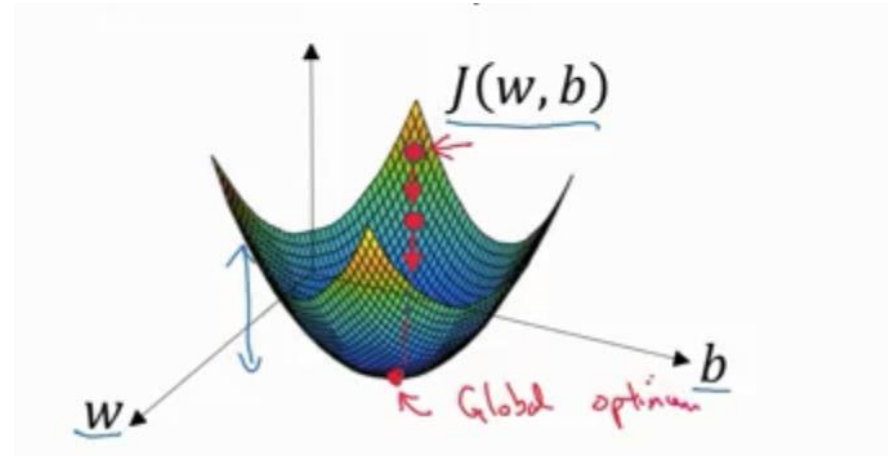
ولازم تكون المعادلة الخاصة بالجي تكون مقعرة , عشان يكون ليها حد ادني واضح اعرف اوصله , وده السبب اللي خلانا ناخذ المعادلة اللوغاريتمية لان نهايتها واضحة وبتكون مقعرة كدة



بينما فيه معادلات تانية مش بتكون مقعرة , فمش هينفع اجيب منها الحد الادني زي كدة



عشان كدة انا ممكن ابدأ بارقام مبدئية للدليو و البي , واعمل تكرار Iteration اكثر من مرة وكل مرة هجيب قيمة للجى اقل , لغاية لما اجيب الـ global minimum



طيب ماذا عن المعادلة الرياضية , ازاي باقدر اوصل للقيمة الدنيا في الـ J ؟

ببساطة بابدأ افرض قيمة دبليو باي قيمة , واستخدم المعادلة ديه عشان اوصل للقيمة الدنيا :

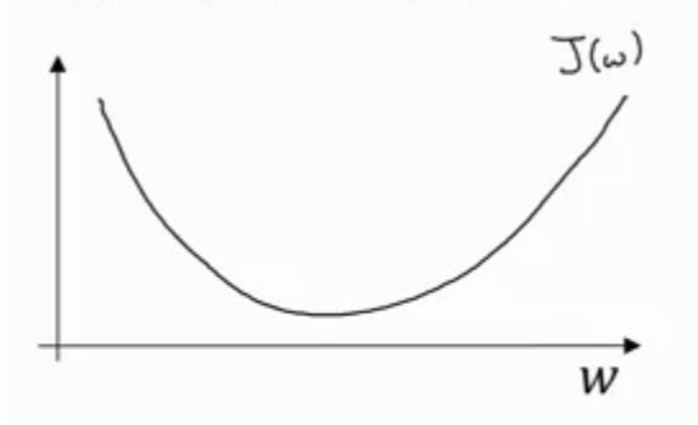
$$\text{Repeat } \left\{ \begin{array}{l} \omega := \omega - \alpha \underbrace{\frac{\partial J(\omega)}{\partial \omega}}_{\text{"dw"}} \\ \end{array} \right. \quad \begin{array}{l} \text{learning rate} \\ \swarrow \end{array}$$

وهي معناه ببساطة ان قيمة دبليو الجديدة , هي القديمة , ناقص حاصل ضرب الفا في تفاضل جي بالنسبة لدبليو

اولا الفا هي الـ learning rate واللي بتحددنا مدي سرعة عمل الـ learning , وكالمعتاد , لو كبيرة هتكون سريعة و مش دقيقة , لو قليلة هتكون دقيقة و بطيئة

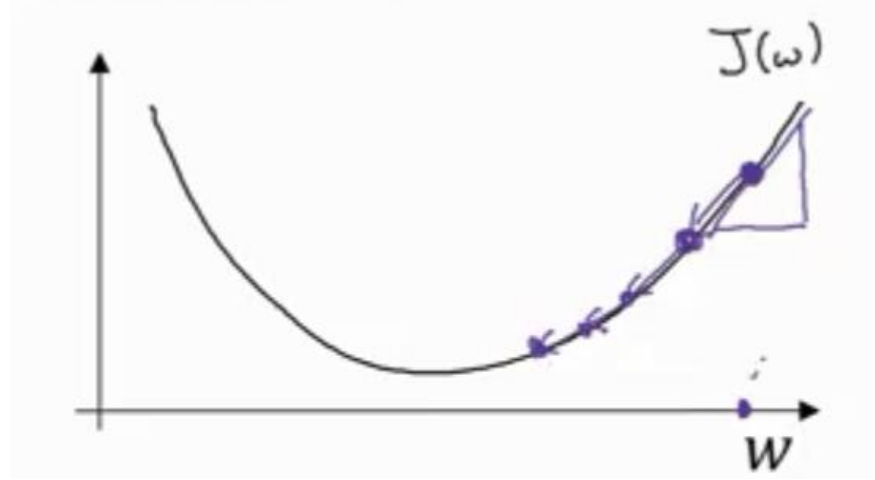
بالنسبة للتفاضل فلازم نشوف الرسم عشان نفهم





لو هنتجاهل بي مؤقتا , ونقول ان الدبليو هي المؤثر الوحيد (منتساش ان دبليو هي مصفوفة الثيتات والبي هي ثيتا 0) , فلو قلنا ان المعادلة اللوغاريتمية هتكون مقعرة كدة , ففيه قيمة للدبليو بحيث تعمل القيمة الدنيا للـ J

فلو فرضنا اني لما اخترت قيمة دبليو كانت علي اليمين , يعني اكثر من قيمة دبليو المثالية زي كدة :

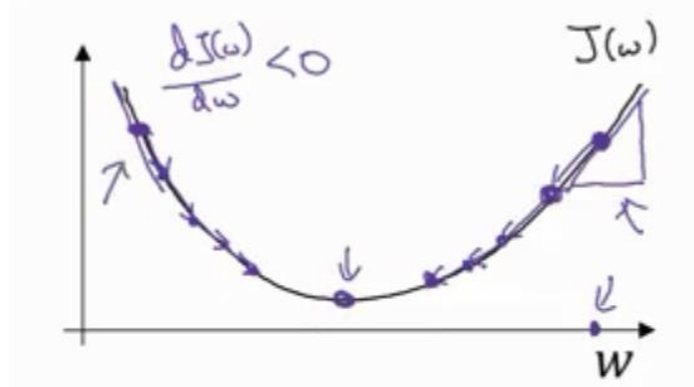


ففي الجانب الايمن , لو عملت اشتقاق للجى , هلاقي ان فيمته هتكون موجبة , وكبيرة شوية , فلما ابص علي المعادلة , هلاقي ان الدبليو الجديدة هتكون القيمة ناقص قيمة التفاضل الكبير شوية مضروب في الفا , يعني قيمة دبليو هتقل شوية , وهنتقل للنقطة اللي علي شمالها

لما اكرر تاني هلاقي ان التفاضل اللي هيتطرح هيكون اقل شوية , لان الميل قل شوية , فالقيمة المطروحة هتكون اقل شوية , وبرضه هدخل شمال

هكرر الموضوع لغاية لما اوصل للحد المثالي للدبليو او قريب منه , هلاقي ان الاشتقاق قريب من الصفر , عشان الخط مستقيم , يعني الدبليو مش هتتغير

ولو كنت اخترت الدبليو علي الشمال زي كدة :



علي الشمال الاشتقاق بالسالب , فلما اضربه في الفا و اطرحه , هلاقي ان دبليو بتزيد معايا مش بتقل , يعني بتدخل يمين , وكل ما ادخل يمين قيمة الاشتقاق بتقل , لغاية برضه لما اوصل للحد الادني

وهيتم في البي نفس اللي حصل في دبليو , مع فارق ان الاشتقاق هيكون بالنسبة لها

$$w := w - \alpha \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \frac{\partial J(w, b)}{\partial b}$$

طب لو فيه اكثر من إكس في المعادلة الاساسية  
مش عايزين ننسي ان المعادلات العادية في حالة وجود اكس هي :

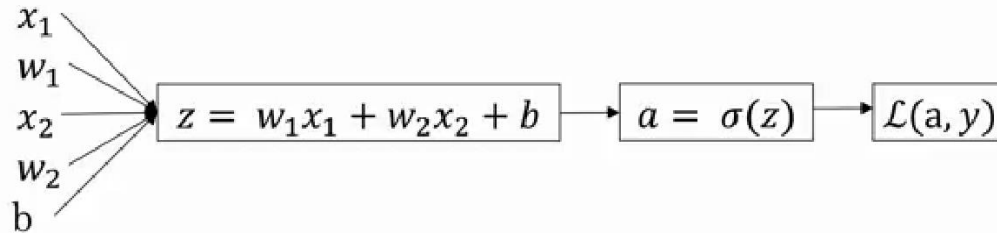
$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

ولو فيه اكثر من اكس , هتكون البي هي هي (لأنها قيمة مضافة ثابتة هيتجمعو مع بعض) لكن هيكون فيه اتنين دبليو لاتنين اكس , لان كل اكس هيكون ليها الثبتات الخاصة بيها والي هي الدبليو

فهتكون كالتالي :



لو هبدا اعمل باك بروباجيشن , يعني ارجاع عكسي , عشان اشوف التغيير في النهايات بيعمل تغيير قد ايه في المداخل

فهذا في التفاضل الاخير . اللي هيكون

dl/da

يعني تفاضل المعادلة النهائية بالنسبة لقيمة ايه , واللي هي قيمة سيجمويد الزى ,, زي ما هي فوق

لو جيت تعمل اشتقاق للمعادلة ديه بالنسبة لايه :

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

هتلاقي ان الواي ثابت بالنسبة لها , و اشتقاق اي لوج هو اشتقاقه فوق علي نفسه, و بالتالي الاشتقاق هيكون :

$$-\frac{y}{a} + \frac{1-y}{1-a}$$

ولو رجعت خطوة ورا , عشان تجيب اشتقاق ايه بالنسبة لزي يعني

da/dz

هيكون قيمته

a(1-a)

وبالتالي عشان اجيب قيمة اشتقاق المعادلة كلها إل بالنسبة لـ  $z$  , هاستخدم قاعدة السلسلة chain rule

$$\frac{\partial \mathcal{L}}{\partial z} = \frac{\partial \mathcal{L}}{\partial a} \cdot \frac{\partial a}{\partial z}$$

يعني هضرب ديه  $-\frac{4}{a} + \frac{1-4}{1-a}$  في ديه  $a(1-a)$

هتعطيني نتيجة مباشرة و هي :

$$a - y$$

و بالتالي النتائج الفرعية لاشتقاق إل بالنسبة لدبليو 1 او 2 او بي , هتكون كدة

$$\frac{\partial \mathcal{L}}{\partial w_1} = "dw_1" = x_1 \cdot dz \quad dw_2 = x_2 \cdot dz \quad db = dz$$

كل اللي فات ده عن عينة واحدة , يعني صف واحد , طيب ماذا عن التطبيق لأكثر من صف , يعني عدد الإيم كبير ؟ ؟  
وقتها المعادلات المعروفة اللي هي كانت كدة

$$z = w^T x + b$$

$$\hat{y} = a = \sigma(z)$$

$$\mathcal{L}(a, y) = -(y \log(a) + (1 - y) \log(1 - a))$$

هتبقى كدة :

$$\begin{aligned} \text{For } i = 1 \text{ to } m \\ z^{(i)} &= w^T x^{(i)} + b \\ a^{(i)} &= \sigma(z^{(i)}) \\ J &= -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})] \end{aligned}$$

يعني كل معادلة هتكون مرتبطة بقيمة بتفاصيل الصف اللي هي معاه  
ولما احي اعمل الاشتقاقات , قيمة الـ  $dz$  اللي هي اصلا  $dj/dz$  اللي هي كانت :

$$a - y$$

هتكون :

$$\underline{dz^{(i)}} = a^{(i)} - y^{(i)}$$

وقيم  $dw_1$  اللي هي اصلا  $dl/dw_1$  واخواتها واللي كانت اصلا :

$$\frac{\partial l}{\partial w_1} = "dw_1" = x_1 \cdot dz \quad dw_2 = x_2 \cdot dz \quad db = dz$$

هتكون كدة

$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}$$

(متنساش ان الإل هي نفسها الجي , بتتكتب كدة او كدة )

بعدها متنساش انك تقسم كل قيم جي و دي دبليو 1 و 2 و بي علي الإم (عدد العناصر) عشان تجيب المتوسط



$$J/m \leftarrow$$

$$\underset{\uparrow}{dw_1/m} ; \underset{\uparrow}{dw_2/m} ; \underset{\uparrow}{db/m}.$$

واخيرا اعمل ابديت لقيم دبليو 1 و 2 و بي , مع طرحها من قيمة الاشتقاق مضروب في الفا (معامل التدريب)

$$w_1 := w_1 - \alpha \underline{dw_1}$$

$$w_2 := w_2 - \alpha \underline{dw_2}$$

$$b := b - \alpha \underline{db}$$

وديه الصورة العامة :

## Logistic regression on $m$ examples

$$J=0; \underline{\Delta w_1}=0; \underline{\Delta w_2}=0; \underline{\Delta b}=0$$

For  $\bar{c} = 1$  to  $m$

$$z^{(i)} = \omega^T x^{(i)} + b$$

$$Q^{(i)} = \sigma(z^{(i)})$$

$$J_t = -[y^{(i)} \log a^{(i)} + (1-y^{(i)}) \log (1-a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_i = x_i^{(i)} dz^{(i)} \quad \uparrow n=2$$

$$dx^{(i)} = x^{(i)} dz^{(i)} \quad \downarrow$$

$$d_b + = dz^{(i)}$$

$$J/m \ll$$

$$\begin{matrix} dw_1/ = m & ; & dw_2/ = m & ; & db/ = m. & \leftarrow \\ \uparrow & & \uparrow & & \uparrow & \end{matrix}$$

$$dw_1 = \frac{\partial J}{\partial w_1}$$

$$\omega_1 := \omega_1 - \alpha \underline{d}\omega_1$$

$$w_2 := w_2 - \alpha \underline{dw}_2$$

$$b := b - \alpha \nabla b$$

Andrew Ng

\* \* \* \* \*

## التوجيه Vectorization

المقصود بها هو التسريع في العمليات الرياضية للاكواد المستخدم عبر استخدام تقنيات محددة , تتعلق بالمصفوفات , اكثر من الدورات (loops) وتسمى فيكتوريزاشن لانها تشير الي المصفوفات

والعملية مهمة لانها بتختصر كميات كبيرة من الوقت , خاصة لما يكون الكود اصلا طويل و محتاج وقت كبير

فمثلا لو انا عايز احسب قيمة زي اللي هي :

$$z = w^T x + b$$

متناساش ان الدبليو و الاكس مصفوفتين كلا منهما عمود واحد , وعدد الصفوف يساوي عدد العناصر إن

$$w = \begin{bmatrix} : \\ : \\ : \end{bmatrix} \quad x = \begin{bmatrix} : \\ : \\ : \end{bmatrix}$$

فلو هعملها بالطريقة العادية (غير الموجهة) هتكون

Non-vectorized:

```
z = 0
for i in range(1-x):
    z += w[i,j] * x[i,j]
z += b
```

بينما لو اتعملت بالطريقة الموجهة هتكون سطر واحد :

Vectorized

$$z = \underbrace{np.dot(w, x)}_{w^T x} + b$$

فالتوجيه مش بس اسهل في كتابة الكود , ده كمان اسرع في التنفيذ بكتير

و هنا فيه مثال عملي في بايثون

```
import numpy as np
import time
```

```
a = np.random.rand(1000000)
b = np.random.rand(1000000)
```

```
tic = time.time()
c = np.dot(a, b)
```

```
toc = time.time()
print(c)
print("vectorized version: " + str(1000 * (toc-tic)) + "ms")

c = 0
tic = time.time()
for i in range(1000000):
    c += a[i]*b[i]

toc = time.time()
print(c)
print("non-vectorized version: " + str(1000 * (toc-tic)) + "ms")
```

```
250347.13984556697
vectorized version: 3.0002593994140625ms
250347.1398455659
non-vectorized version: 1409.0805053710938ms
```

تم عمل عملية رياضية وهي ضرب مليون رقم عشوائي في مليون رقم عشوائي , لما عملناها بالطريقة العادية اخدت 1409 ملي ثانية , بينما بالطريقة العادية اخدت 3 ملي ثانية , والنتيجة هي هي , فطبعا في العمليات المعقدة سيكون الفرق باين جدا

وهتلاحظ ان الفرق تقريبا من 300 ل 400 ضعف , و ده معناه ان العملية اللي هتاخذ دقيقة واحد في العملية الموجهة , هتحتاج 5 ساعات في الغير موجهة

ملحوظة :

CPU : Central Process Unit (normal machine)

GPU : Graphical Process Unit (plugged device)

SIMD : Single Instruction Multiple Data

فالأساس هنا انك تحاول تتجنب استخدام (فور لووب) علي قد ما تقدر , طالما فيه بديل ليها

فمصلا لو عايز ا ضرب مصفوفة بمصفوف و عواميد , في فيكتور , اللي هيكون عمود واحد و فيه صفوف زي كدة

$$u = Av$$

فممكن اضربها بالطريقة العادية الطويلة للكود كدة :

```
u = np.zeros(n,1)
for i ...
    for j ...
        u[i] += A[i][j]*v[j]
```

او الاسرع اني استخدم الطريق الموجهة , بكود مباشر كدة :

`A = np.dot(A, v)`

`u = np.dot(A, v)`

نفس الفكرة , لو عايز مثلا اعمل اكسبونينشيل لعدد من العناصر زي كدة

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix} \rightarrow u = \begin{bmatrix} e^{v_1} \\ e^{v_2} \\ \vdots \\ e^{v_n} \end{bmatrix}$$

فبالطريقة العادية :

```
→ u = np.zeros((n,1))  
→ for i in range(n): ←  
    → u[i]=math.exp(v[i])
```

```
u = np.zeros((n, 1))  
for i in range(n):  
    u[i] = math.exp(v[i])
```

او بالطريقة الموجهة :

```
import numpy as np  
u = np.exp(v) ←  
  ^      ^
```

```
u = np.exp(v)
```

ونفس الفكرة لباقي الدوال المشابهة :

$$\begin{array}{l} n p. \log(v) \\ n p. \log(v) \\ n p. \log(v, 0) \\ v \times v^2 \quad v/v \end{array}$$

ولو هنوسع المثال شوية فمثلا في الكود الكبير الخاص التصنيف اللي هو :

## Logistic regression derivatives

```
J = 0, dw1 = 0, dw2 = 0, db = 0
for i = 1 to 'm:
    z(i) = wTx(i) + b
    a(i) = σ(z(i))
    J += -[y(i) log ŷ(i) + (1 - y(i)) log(1 - ŷ(i))]
    dz(i) = a(i)(1 - a(i))
    dw1 += x1(i) dz(i)
    dw2 += x2(i) dz(i)
    db += dz(i)
J = J/m, dw1 = dw1/m, dw2 = dw2/m, db = db/m
```

هتلاحظ ان فيه اثنين فور لوب , الاول فوق المكتوبة , والثانية اللي عند dw اللي هي



$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}$$

وده في حالة ان فيه اكثر من ديليو

فممكن استغني عن الفور هنا عن طريق اني اعمل فيكتور باسم ديليو فوق كدة :

$$dw = np.zeros((n-x, 1))$$

بدل كدة :

$$dw1 = 0, dw2 = 0,$$

ووقتها هيكون الضرب كدة

$$dw += x^{(i)} dz^{(i)}$$

بدل كدة :

$$\begin{aligned} dw_1 &+= x_1^{(i)} dz^{(i)} \\ dw_2 &+= x_2^{(i)} dz^{(i)} \\ db &+= dz^{(i)} \end{aligned}$$

و في القسمة علي إم هتكون قسمة  $dw$  مباشرة , بدل ما اقسم عنصر عنصر

$$dw_1 = dw_1/m, dw_2 = dw_2/m$$

كمثال ثاني اعمق  
لو قلنا في التصنيف اني عايز اجيب زد 1 , زد 2 , وهكذا , هتكون معادلاتهم كدة :

$$\begin{aligned} z^{(1)} &= w^T x^{(1)} + b & z^{(2)} &= w^T x^{(2)} + b & z^{(3)} &= w^T x^{(3)} + b \\ a^{(1)} &= \sigma(z^{(1)}) & a^{(2)} &= \sigma(z^{(2)}) & a^{(3)} &= \sigma(z^{(3)}) \end{aligned}$$

متنساش ان الدبليو هي عبارة عن مصفوفة ثيئات , بالشكل ده

$$\Theta = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \theta_2 \\ \vdots \\ \theta_n \end{bmatrix} \left\{ \begin{array}{l} b \leftarrow \\ w \leftarrow \end{array} \right.$$

و بنسميها رمز دبليو كبير

اما كل اكس من الاكسات , فهنجمعها في مصفوفة  $X$  كبيرة اللي هتكون عواميدا بعدد الإكسات (عدد العناصر  $m$ ) و صفوفها بعدد العناصر  $n$  اللي هي اصلا عدد صفوف اي اكس فيهم

$$X = \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

ويمكن اقول علي المطلوب (الزندات) ان فيه مصفوفة كبيرة اسمها  $Z$  من صف واحد و هي تجمع كل الزندات في عواميد كدة

$$Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(n)}]$$

ساعتها ممكن نقول ان المعادلة النهائية هي :

$$Z = w^T X + b$$

فالدليو ترانسبوز هتكون صف واحد و عواميد فيها الثينات , ولما تضرب في  $X$  اللي هتكون كدة

$$w^T \begin{bmatrix} x^{(1)} & x^{(2)} & \dots & x^{(n)} \\ 1 & 1 & \dots & 1 \end{bmatrix}$$

هتكون النتيجة عبارة عن قيم زي كدة

$$[w^T x^{(1)} \quad w^T x^{(2)} \quad \dots \quad w^T x^{(n)}]$$

لما اجي اجمع عليها قيمة  $b$  فهضيفها لكل عنصر فيها زي كدة

$$[w^T x^{(1)} + b \quad w^T x^{(2)} + b \quad \dots \quad w^T x^{(n)} + b]$$

يعني هتعملي قيم زندات

لاحظ ان قيمة  $b$  هي قيمة مفردة يعني مصفوفة 1 في 1 , بس بايثون لما بيشفوك بتجمع رقم علي مصفوفة , فهيفرد الرقم عشان يكون مصفوفة بنفس حجم اللي بيتجمع عليها , والعملية ديه اسمها broadcasting

وبالتالي المعادلة النهائية هتكون :

$$Z = np.dot(w.T, X) + b$$

ونفس الفكرة في اني اجيب الـ  $a$  عن طريق عمل مصفوفة كبيرة اسمها  $A$  و يكون فيها قيم  $a$  كلها

$$A = [a^{(1)} \ a^{(2)} \ \dots \ a^{(n)}] = \sigma(Z)$$

و ده عن طريق استخدام دالة سيجمويد لكل قيم زي  $Z$  مع بعض , قيم زي كابيتال اللي هي مصفوفة كل قيم زي اللي جبتها

\*\*\*\*\*

طيب ماذا عن ايجاد قيم الاشتقاقات :

عارفين من قبل كدة ان اشتقاق الزى هيكون بالمعادلة :

$$a - y$$

وبالتالي جميع اشتقاقات الزيهات هتكون :

$$dz^{(1)} = a^{(1)} - y^{(1)} \quad dz^{(2)} = a^{(2)} - y^{(2)}$$

فهنعمل تجميع لاشتقاقات الزيهات في  $dZ$  كبيرة , واللي هتحتوي علي  $dz$  صغيرة

$$dZ = [dz^{(1)} \quad dz^{(2)} \quad \dots \quad dz^{(m)}]$$

$1 \times m$

و ممكن نجمع الايهات الصغيرة في  $A$  كبيرة , وكذلك الوايات

$$A = [a^{(1)} \quad \dots \quad a^{(m)}] \quad Y = [y^{(1)} \quad \dots \quad y^{(m)}]$$

فهتكون المعادلة المجمعة ليهم هي :

$$dZ = A - Y = [a^{(1)} - y^{(1)} \quad a^{(2)} - y^{(2)} \quad \dots]$$

طيب ماذا عن باقي الاشتقاقات , دبليو و بي ؟

الـ  $dw$  اساسا هي المفروض مجموع كل اكس في الـ  $dz$  الخاصة بيها, بعدها نقسمها علي  $m$  , يعني كدة

$$\begin{aligned}
 dw &= 0 \\
 dw &+= x^{(1)} dz^{(1)} \\
 dw &+= x^{(2)} dz^{(2)} \\
 &\vdots \\
 dw &= m
 \end{aligned}$$

والـ  $db$  هتكون مجموع الـ  $dz$  مع بعضها و قسمتها علي  $m$

$$\begin{aligned}
 db &= 0 \\
 db &+= dz^{(1)} \\
 db &+= dz^{(2)} \\
 &\vdots \\
 db &+= dz^{(n)} \\
 db &= m
 \end{aligned}$$

فلو بدأنا بالـ  $db$  هنعملها ببساطة مجموع مصفوفة فيها كل الـ  $dz$  مع بعض

$$= \frac{1}{m} \text{ np.sum}(dz)$$

أما الـ  $dw$  فهتكون معادلتها

$$dw = \frac{1}{m} X dz^T$$

متنساش ان الـ  $X$  الكبيرة هي مصفوفة صف واحد و عواميدها هي قيم إكسات لكل عينة من عينات الـ  $m$   
أما الـ  $dz$  الكبيرة فهي كمان مصفوفة صف واحد و عواميدها قيم  $dz$  الصغيرة

فلما تعمل تدوير للـ  $dz$  هتكون عمود واحد بصفوف كتير , اضرب اكس فيها الاول هتتحول للقيم :

$$\underline{x^{(1)} dz^{(1)}} + \dots + \underline{x^{(m)} dz^{(m)}}$$

ولما تقسمها على  $m$  , هتجيبك القيمة النهائية ليها , وبسطر واحد من الكود , ومن غير اي استخدام لفور لوب اللي بتضيع الوقت

وبالتالي العملية كلها اللي هي ديه

```
J = 0, dw1 = 0, dw2 = 0, db = 0
for i = 1 to m:
    z(i) = wTx(i) + b
    a(i) = σ(z(i))
    J += -[y(i) log a(i) + (1 - y(i)) log(1 - a(i))]
    dz(i) = a(i) - y(i)
    dw1 += x1(i) dz(i)
    dw2 += x2(i) dz(i)
    db += dz(i)
J = J/m, dw1 = dw1/m, dw2 = dw2/m
db = db/m
```

هنتحول للكود ده

$$z = w^T X + b$$

$$= n p \cdot \text{dot}(w.T, X) + b$$

$$A = G(z)$$

$$dz = A - Y$$

$$dW = \frac{1}{m} \times dz^T$$

$$db = \frac{1}{m} \text{ np.sum}(d-z)$$

$$w := w - \alpha dw$$

$$b := b - \alpha \Delta b$$

\* \* \* \* \*



## عملية الـ broadcasting

و هي تستخدم في بايثون بالتحديد , عشان اقدر اتعامل مع مجموعة من العمليات مرة واحدة و بكود بسيط بدل ما اعمل عدد كبير من السطور , وده بيتم غالبا باستخدام مكتبات بايثون العديدة

فمثلا لو عندي مصفوفة فيها ارقام زي كدة

	Apples	Beef	Eggs	Potatoes
Carb	56.0	0.0	4.4	68.0
Protein	1.2	104.0	52.0	8.0
Fat	1.8	135.0	99.0	0.9

واللي هي عبارة عن عدد الكالوري اللي في الكربوهيدرات , والبروتين و الدهون , في اربع اكلات معينة

فلو انا محتاج اعمل نسبة مئوية لكل اكلة فيهم , يعني اشوف ال 56 كالوري في الكربوهيدرات نسبتها قد ايه من كل الكالوريز (يعني علي مجموع الـ 56 مع 1.2 , مع 1.8 الخاصة بالتفاح) , ساعتها ممكن اعملها كدة :

```
import numpy as np
```

```
A = np.array([ [56.0, 0.0, 4.4, 68.0],  
               [1.2, 104.0, 52.0, 8.0],  
               [1.8, 135.0, 99.0, 0.9] ] )
```

```
cal = A.sum(axis=0)  
print(cal)
```

```
# [59. 239. 155.4 76.9]
```

```
percentage = 100*A/cal.reshape(1,4)
```

```
print(percentage)
```

```
#[[ 94.91525424  0.          2.83140283 88.42652796]
```

```
# [ 2.03389831 43.51464435 33.46203346 10.40312094]
```

```
# [ 3.05084746 56.48535565 63.70656371 1.17035111]]
```

اللي حصل اني كتبت المصفوفة , و عملت مجموع لكل عمود , بعدها قسمت كل قيمة فيهم علي رقم كل عمود من المصفوفة الثانية

و أمر sum مع كلمة axis=0 معناها اجمع راسيا , ولو قلته تساوي 1 هيجمع افقيا , ولو من غير كلمة اكسي خالص هيجمع كل القيم لقيمة واحدة

بالنسبة لأمر reshape هو مكتوب احتياطي مش أكثر , لان اصلا المصفوفة اللي طلعت حالا كانت هي 1 في 4 , فلو مكتبتهاش مفيش مشكلة

طيب ماذا عن عملية القسمة , والمعروف ان القسمة في المصفوفة بتتم في حالتين , الأول لو هاقسم علي رقم واحد , فبيتم قسمة كل ارقام المصفوفة علي نفس الرقم ده

الحالة الثانية لو هاقسم علي مصفوفة بنفس ابعادها بالضبط , ساعتها بيتقسم كل رقم من المصفوفة الاولى علي الرقم من المصفوفة الثانية

فلو هاقسم المصفوفة ديه

15	17	4	18
25	3	22	14

علي ديه

3	1	2	6
5	1	11	7

هتكون النتيجة كدة

5	17	2	3
5	3	2	2

تعالى نشرح بالتفصيل  
لو جيت اجمع مصفوفة علي رقم , فبايثون اوتوماتيك هيحول الرقم ده لمصفوفة

$$\begin{bmatrix} 1 \\ 2 \\ 3 \\ 4 \end{bmatrix} + \begin{bmatrix} 100 \\ 100 \\ 100 \\ 100 \end{bmatrix} = \begin{bmatrix} 101 \\ 102 \\ 103 \\ 104 \end{bmatrix}$$

فلو جمعنا المصفوفة ديه علي رقم 100 , فهو اوتوماتيك هيجولها لمصفوفة شبيهها , و هتعمل الارقام ديه

كذلك الحال لو تم عمل جمع مصفوفة علي مصفوفة تانية مش مقاسها , فلو هينفع تكبر و تتنسخ عشان تكون قدها هيعمل كدة فمثلا

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + [100 \quad 200 \quad 300]$$

هتكون كدة

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 200 & 300 \\ 100 & 200 & 300 \end{bmatrix} = \begin{bmatrix} 101 & 202 & 303 \\ 104 & 205 & 306 \end{bmatrix}$$

$\downarrow \quad \downarrow \quad \downarrow$   
 $(m,n) \quad (2,3) \quad (1,n) \rightsquigarrow (m,n) \quad (2,3)$

كمان لما نجمع ديه علي ديه

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 \\ 200 \end{bmatrix}$$

هيعملها كدة

$$\begin{bmatrix} 1 & 2 & 3 \\ 4 & 5 & 6 \end{bmatrix} + \begin{bmatrix} 100 & 100 & 100 \\ 200 & 200 & 200 \end{bmatrix} = \begin{bmatrix} 101 & 102 & 103 \\ 204 & 205 & 206 \end{bmatrix}$$

وبالتالي لما يتم عملية رياضية بين مصفوفتين , والمصفوفة الثانية مش متلائمة مع الاولى , بايثون يحولها لو ينفذ

عشان كدة لما قسمنا مصفوفة 3 في 4 , علي مصفوفة 1 في 4 , بايثون مد المصفوفة المقسوم عليها اللي هي 1 في 4 , وخلاها 3 في 4 زي الاول , بحيث يقسم كل رقم علي اخوه , وبالتالي تم قسمة كل رقم من ارقام الكالوري , علي المجموع عشان يجيب النسبة المئوية

\* \* \* \* \*

و كون ان بايثون فيه خاصية ال broadcasting (اللي هي تمديد المصفوفات حسب الشكل المطلوب) فديه ميزة و عيب

ميزة انها بتختصر وقت و جهد في الكتابة , و عيب انها بتخلي البرنامج يعمل حاجات ممكن انت مش مخططها , فهتعمل مشاكل

يعني مثلا لو عملنا الكود ده :

```
import numpy as np
a = np.random.randn(5)
print(a)
# [ 1.2, 2.3, 3.4, 4.5, 5.6 ]
print(a.shape)
# (5,)
print(a.T)
# [ 1.2, 2.3, 3.4, 4.5, 5.6 ]
```

```
Print(np.dot(a,a.T))
4.06
```

هنا نلاحظ ان في الاول عملنا مصفوفة 5 صفوف في عمود واحد , ولما ضربناها في التدوير بتاعها اللي هيكون 1 في 5 , المفروض يعملنا مصفوفة 5 في 5 , لكن بايثون عمل تحويل للنوع ده فورا فعمل الاول 1 في 5 , والثانية 5 في 1 , فالنتج طلع 1 في 1 فلو انت عايز تجبر بايثون علي مسار معين , متستخدمش امر

```
a = np.random.randn(5)
```

لكن

```
a = np.random.randn(5,1)
```

يعني متخليش بايثون يفرضها اوتوماتيك لكن حددهاله واللي هيعملها دلوقتي 5 في 1 , ولما تضربها في تدويرها هيعملك مصفوفة 5 في 5 فالصيغة ديه مش كويسة

```
a = np.random.randn(5)
a.shape = (5,)
"rank 1 array"
```

## الاحسن نستخدم ديه ليفيكتور الصف او العمود

`a = np.random.randn(5,1)` → `a.shape = (5,1)` column vector

`a = np.random.randn(1,5)`  $\rightarrow$  `a.shape = (1,5)` row vector

او ممكن استخدم امر ريشاب لو تم كتابة الفيكتور بشكل مش سليم

```
a = a.reshape((5,1))
```

\* \* \* \* \*