# Python and Vectorization

## Vectorization



The vectorize version took 1.5 milliseconds. The explicit for loop and non-vectorize version took about 400, almost 500 milliseconds. The non-vectorize version took something like 300 times longer than the vectorize version. With this example you see that if only you remember to vectorize your code, your

code actually runs over 300 times faster. Let's just run it again. Just run it again. Yeah. Vectorize version 1.5 milliseconds seconds and the for loop.

So 481 milliseconds, again, about 300 times slower to do the explicit for loop. If the engine x slows down, it's the difference between your code taking maybe one minute to run versus taking say five hours to run.

And when you are implementing deep learning algorithms, you can really get a result back faster. It will be much faster if you vectorize your code. Some of you might have heard that a lot of scaleable deep learning implementations are done on a GPU or a (graphics processing unit). But all the demos I did just now in the Jupiter notebook where actually on the CPU. And it turns out that both GPU and CPU have parallelization instructions.

They're sometimes called SIMD instructions. This stands for a single instruction multiple data. But what this basically means is that, if you use built-in functions such as this np.function or other functions that don't require you explicitly implementing a for loop.

> The rule of thumb to remember is whenever possible, avoid using explicit for loops.

# More Vectorization Examples

## Vectors and matrix valued functions

Say you need to apply the exponential operation on every element of a matrix/vector.

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

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

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

Andrew Ng

## Logistic regression derivatives

$J = 0,$ ~~$dw1 = 0,$ $dw2 = 0$~~, $db = 0$ $\qquad dw = np.zeros((n_x, 1))$

→ for i = 1 to m:
$$z^{(i)} = w^T x^{(i)} + b$$
$$a^{(i)} = \sigma(z^{(i)})$$
$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$
$$dz^{(i)} = a^{(i)} - y^{(i)}$$

for $j=1...n_x$
for $dw_j +=$

$$dw_1 += x_1^{(i)} dz^{(i)}$$
$$dw_2 += x_2^{(i)} dz^{(i)}$$
$n_x = 2$ $\qquad dw += x^{(i)} dz^{(i)}$
$$db += dz^{(i)}$$

$J = J/m,$ ~~$dw_1 = dw_1/m,$ $dw_2 = dw_2/m$~~, $db = db/m$

$$dw \mathrel{/}= m.$$

Andrew Ng

---

# Vectorizing Logistic Regression

## Vectorizing Logistic Regression

$$z^{(1)} = w^T x^{(1)} + b \qquad z^{(2)} = w^T x^{(2)} + b \qquad z^{(3)} = w^T x^{(3)} + b$$
$$a^{(1)} = \sigma(z^{(1)}) \qquad a^{(2)} = \sigma(z^{(2)}) \qquad a^{(3)} = \sigma(z^{(3)})$$

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

$(n_x, m)$

$\mathbb{R}^{n_x \times m}$

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

$$Z = \begin{bmatrix} z^{(1)} & z^{(2)} & \cdots & z^{(m)} \end{bmatrix} = w^T X + [b\ b \cdots b] = \begin{bmatrix} w^T x^{(1)} + b & w^T x^{(1)} + b & \cdots & w^T x^{(m)} + b \end{bmatrix}$$

$1 \times m$

$1 \times m$

"Broadcasting"

$$Z = np.dot(w.T, X) + b$$
$(1,1)$ $\mathbb{R}$

$$A = \begin{bmatrix} a^{(1)} & a^{(2)} & \cdots & a^{(m)} \end{bmatrix} = \sigma(Z)$$

Andrew Ng

## Vectorizing Logistic Regression

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

$$dZ = [dz^{(1)} \; dz^{(2)} \cdots dz^{(m)}] \quad \leftarrow$$
$$\underset{1 \times m}{}$$

$$A = [a^{(1)} \cdots a^{(m)}]. \qquad Y = [y^{(1)} \cdots y^{(m)}]$$

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

$$\rightarrow dw = 0$$
$$\begin{bmatrix} dw \mathrel{+}= x^{(1)} dz^{(1)} \\ dw \mathrel{+}= x^{(2)} dz^{(2)} \\ \vdots \\ dw/ = m \end{bmatrix}$$

$$\begin{bmatrix} db = 0 \\ db \mathrel{+}= dz^{(1)} \\ db \mathrel{+}= dz^{(2)} \\ \vdots \\ db \mathrel{+}= dz^{(m)} \\ db/ = m. \end{bmatrix}$$

$$db = \frac{1}{m} \sum_{i=1}^{m} dz^{(i)}$$
$$= \frac{1}{m} \, np.\text{sum}(dZ)$$

$$dw = \frac{1}{m} X \, dz^{T}$$
$$= \frac{1}{m} \begin{bmatrix} x^{(1)} \cdots x^{(m)} \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$
$$= \frac{1}{m} \begin{bmatrix} x^{(1)} dz^{(1)} + \cdots + x^{(m)} dz^{(m)} \end{bmatrix}$$
$$\underset{n \times 1}{}$$

The left side is non-vectorize approach, and the right side is the vectorize approach.

## Implementing Logistic Regression

for iter in range (1000): $\leftarrow$

```
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⁾ ]    dw += x⁽ⁱ⁾ * dz⁽ⁱ⁾
    db += dz⁽ⁱ⁾
J = J/m,  dw₁ = dw₁/m,  dw₂ = dw₂/m
db = db/m
```

$$Z = w^{T}X + b$$
$$= np.\text{dot}(w.T, X) + b$$
$$A = \sigma(Z)$$
$$dZ = A - Y$$
$$dw = \frac{1}{m} X \, dZ^{T}$$
$$db = \frac{1}{m} \, np.\text{sum}(dZ)$$

$$w := w - \alpha \, dw$$
$$b := b - \alpha \, db$$

Broadcasting turns out to be a technique that Python and numpy allows you to use to make certain parts of your code also much more efficient.

# Broadcasting in Python

In the previous video, I mentioned that broadcasting is another technique that you can use to make your Python code run faster. In this video, let's delve into how broadcasting in Python actually works.

# Broadcasting example

Calories from Carbs, Proteins, Fats in 100g of different foods:

$$A = \begin{array}{c c} & \begin{array}{cccc} \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \end{array} \\ \begin{array}{c} \text{Carb} \\ \text{Protein} \\ \text{Fat} \end{array} & \left[ \begin{array}{cccc} 56.0 & 0.0 & 4.4 & 68.0 \\ 1.2 & 104.0 & 52.0 & 8.0 \\ 1.8 & 135.0 & 99.0 & 0.9 \end{array} \right] \end{array} = A \quad (3,4)$$

59 cal $\quad \frac{56}{59} \approx 94.9\%$

Calculate % of calories from Carb, Protein, Fat. Can you do this without explicit for-loop?

problem:

If you look at this (the first column) column and add up the numbers in that column you get that 100 grams of apple has 56 plus 1.2 plus 1.8 so that's 59 calories.

And so as a percentage the percentage of calories from carbohydrates in an apple would be 56/59, that's about 94.9%. So most of the calories in an apple come from carbs, whereas in contrast, most of the calories of beef come from protein and fat and so on.

So the calculation you want is really to sum up each of the four columns of this matrix to get the total number of calories in 100 grams of apples, beef, eggs, and potatoes. And then to divide throughout the matrix, so as to get the percentage of calories from carbs, proteins and fats for each of the four foods.

So the question is, can you do this without an explicit for-loop? Let's take a look at how you could do that.

```
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]])

print(A)
```

```
[[ 56.    0.    4.4  68. ]
 [  1.2 104.   52.    8. ]
 [  1.8 135.   99.    0.9]]
```

In [7]:
```
cal = A.sum(axis=0)
print(cal)
```

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

In [8]:
```
percentage = 100*A/cal.reshape(1,4)
print(percentage)
```
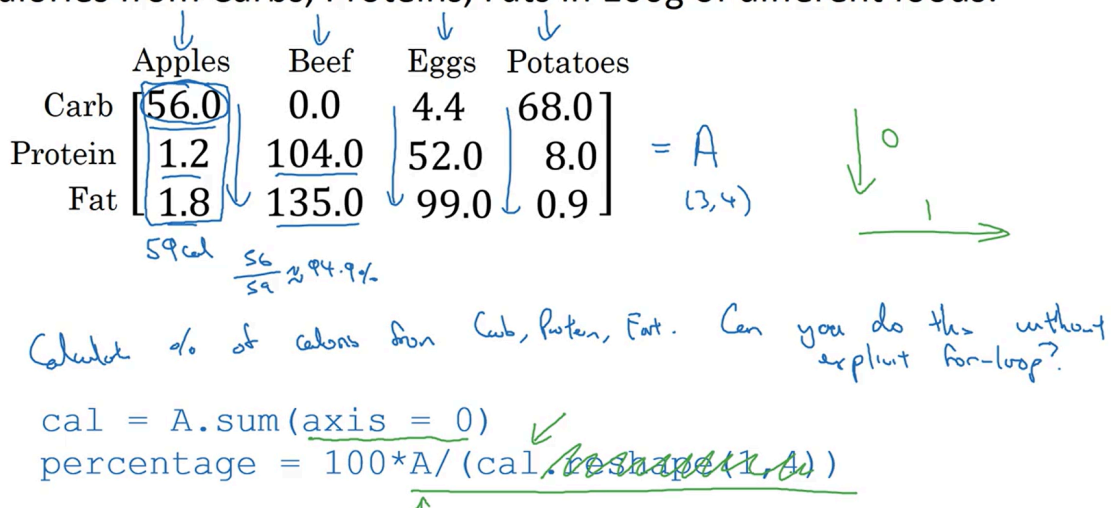
```
[[ 94.91525424  0.          2.83140283  88.42652796]
 [  2.03389831 43.51464435 33.46203346 10.40312094]
 [  3.05084746 56.48535565 63.70656371  1.17035111]]
```

In [ ]:

# Broadcasting example

Calories from Carbs, Proteins, Fats in 100g of different foods:

$$
\begin{array}{c}
 & \text{Apples} & \text{Beef} & \text{Eggs} & \text{Potatoes} \\
\text{Carb} & 56.0 & 0.0 & 4.4 & 68.0 \\
\text{Protein} & 1.2 & 104.0 & 52.0 & 8.0 \\
\text{Fat} & 1.8 & 135.0 & 99.0 & 0.9
\end{array} = A \quad (3,4)
$$

59 cal

$\frac{56}{59} \approx 94.9\%$

Calculate % of calories from Carb, Protein, Fat. Can you do this without explicit for-loop?

```
cal = A.sum(axis = 0)
percentage = 100*A/(cal.reshape(1,4))
```

axis = 0 → sum up vertically

axis = 1 → sum up horizontally

because the shape of 'cal' variable is already (1,4) so that we can remove the 'reshape' term.

other broadcasting example:

# Broadcasting example

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

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

$(m,n)$ $(2,3)$ $\quad (1,n) \leadsto (m,n) \quad (2,3)$

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

$(m,n)$ $\qquad (m,1) \downarrow (m,n)$

# General Principle

$(m,n)$ matrix $\qquad \dfrac{+}{\times} \qquad (1,n) \leadsto (m,n)$

$\qquad\qquad\qquad\qquad\qquad (m,1) \leadsto (m,n)$

$(m,1) \qquad + \qquad \mathbb{R}$

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

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

Matlab/Octave: bsxfun

Before we wrap up, just one last comment, which is for those of you that are used to programming in either MATLAB or Octave, if you've ever used the MATLAB or Octave function bsxfun in neural network programming bsxfun does something similar, not quite the same. But it is often used for similar purpose as what we use broadcasting in Python for. But this is really only for very advanced MATLAB and Octave users, if you've not heard of this, don't

worry about it. You don't need to know it when you're coding up neural networks in Python. So, that was broadcasting in Python.

 I hope that when you do the programming homework that broadcasting will allow you to not only make a code run faster, but also help you get what you want done with fewer lines of code.

# A Note on Python/Numpy Vectors

## Python/numpy vectors

```
a = np.random.randn(5)
        a.shape = (5,)          } Don't use
        "rank 1 array"

a = np.random.randn(5,1)  → a.shape = (5,1)    Column vector ✓

a = np.random.randn(1,5)  → a.shape = (1,5)    row vector ✓

assert(a.shape == (5,1)) ←
        a = a.reshape((5,1))
```

One more thing that I do a lot in my code is if I'm not entirely sure what's the dimension of one of my vectors, I'll often throw in an assertion statement like this, to make sure, in this case, that this is a (5,1) vector. So this is a column vector.

These assertions are really inexpensive to execute, and they also help to serve as documentation for your code. So don't hesitate to throw in assertion statements like this whenever you feel like. And then finally, if for some reason you do end up with a rank 1 array, You can reshape this, a equals a.reshape into say a (5,1) array or a (1,5) array so that it behaves more consistently as either column vector or row vector. So I've sometimes seen students end up with very hard to track because those are the nonintuitive effects of rank 1 arrays. By eliminating rank 1 arrays in my old code, I think my code became simpler. And I did not actually find it restrictive in terms of things I could express in code. I just never used a rank 1 array.

# Explanation of Logistic Regression Cost Function



## Logistic regression cost function

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

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

If $y=1$: $p(y|x) = \hat{y} \quad (1-\hat{y})^0$
$= 1$

If $y=0$: $p(y|x) = \hat{y}^0 \quad (1-\hat{y})^{(1-0)} = 1 \times (1-\hat{y}) = 1-\hat{y}$
$=1$

$\log p(y|x) = \log \hat{y}^y (1-\hat{y})^{(1-y)} = y \log \hat{y} + (1-y) \log(1-\hat{y})$
$= -\mathcal{L}(\hat{y}, y)$

Andrew Ng

So what we've just shown is that this equation is a correct definition for p(ylx).

Now, finally, because the log function is a strictly monotonically increasing function, your maximizing log p(y|x) should give you a similar result as optimizing p(y|x).

And if you compute log of p(y|x), that's equal to log of y hat to the power of y, 1 - y hat to the power of 1 - y. And so that simplifies to y log y hat + 1- y times log 1- y hat, right?

And so this is actually negative of the loss function that we had to find previously. And there's a negative sign there because usually if you're training a learning algorithm, you want to make probabilities large whereas in logistic regression we're expressing this. We want to minimize the loss function.

# Cost on $m$ examples

$$\log p(\text{labels in training set}) = \log \prod_{i=1}^{m} p(y^{(i)} | x^{(i)}) \leftarrow$$

$$\log p(\cdots) = \sum_{i=1}^{m} \log p(y^{(i)} | x^{(i)})$$

$$-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

$$= -\sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Maximum likelihood estimate $\nwarrow$

Cost: (minimize)
$$J(w, b) = \frac{1}{m} \sum_{i=1}^{m} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$