

Week 2 Vectorization

In the deep learning era vectorization is a key skill.

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

$$w = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad x = \begin{bmatrix} \vdots \\ \vdots \\ \vdots \end{bmatrix} \quad w \in \mathbb{R}^{n_x} \quad x \in \mathbb{R}^{n_x}$$

Non-vectorized:

```

Z = 0
for i in range(n_X):
    Z += w[i] * x[i]
Z += b

```

this is going to be really slow

Vectorized

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

in contrast, we just transpose x directly
this is going to be much faster

```
import numpy as np
```

```
a = np.array([1, 2, 3, 4])
```

```
import time
```

```
a = np.random.rand(1,000,000)
b = np.random.rand(1,000,000)
```

```

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(1,000,000):
    c += a[i] * b[i]
toc = time.time()

```

```

print(c)
print("for loop: " + str(1000 * (toc - tic)) + "ms")

```

Vectorized version ~ 1.5ms

Non-vectorized ~ 500ms

Non-vectorized version runs about 300 times faster 300x

Both CPUs and GPUs have parallelization instructions. "SIMD" instructions Single Instruction Multiple Data. If you use numpy functions you can take much better advantage of parallelism.

True of both but GPUs are remarkably good at SIMD calculations. CPUs are not too bad though.

More vectorization Examples

The key rule of thumb: avoid explicit for-loops whenever possible.

It's not always possible to avoid for-loops but you want to look for every way to.

Want: compute vector u as the product of matrix A times some vector v .

$$u = Av$$

$$u_i = \sum_j A_{ij} v_j$$

$$u = np.zeros(n, 1)$$

for i...

for j...

$$u[i] += A[i][j] * v[j]$$

non-vectorized (slow)

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

vectorized (way faster)
eliminates 2 for loops

Want: apply the exponential operation on every element of a matrix/vector

$$v = \begin{bmatrix} v_1 \\ \vdots \\ v_n \end{bmatrix}$$

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

non-vectorized

$$u = np.zeros((n, 1))$$

for i in range(n):

$$u[i] = math.exp(v[i])$$

vectorized

import numpy as np

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

numpy has many vectorized functions

$$np.log(v)$$

$$np.abs(v)$$

$$np.maximum(v, 0) \leftarrow \text{element-wise max (compared to 0 here)}$$

$$v ** 2 \leftarrow \text{element-wise square of each element of } v$$

$$1 / v \leftarrow \text{element-wise inverse}$$

In our logistic regression example we had 2 for-loops

$$J = 0, dw1 = 0, dw2 = 0, db = 0$$

→ for i = 1 to m:

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

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

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

$$dz^{(i)} = a^{(i)}(1 - a^{(i)})$$

$$\rightarrow \text{for } i = 1 \text{ to } n_x: dw_j += x_j^{(i)} dz^{(i)}$$

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

$$J = J/m, dw_j = dw_j/m \dots, db = db/m$$

Let's look at the second for loop. Make dw a vector $dw = np.zeros((n_x, 1))$

$$dw = dw/m$$

Removed the for loop over the features n_x

Still have for loop over all the training examples but you can also get rid of this with vectorization.

Vectorizing Logistic Regression

Implement a single elevation of gradient descent with respect to an entire training set without using a single explicit for loop.

Forward prop:

$$\begin{array}{ll} \text{1st example } (1) & \text{2nd example } (2) \\ z^{(1)} = w^T x^{(1)} + b & z^{(2)} = w^T x^{(2)} + b \\ a^{(1)} = G(z^{(1)}) & a^{(2)} = G(z^{(2)}) \end{array} \dots m$$

$$X = \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} \in \mathbb{R}^{n \times m}$$

recall that we stacked our inputs in an (n_x, m) matrix

construct a $(1, m)$ matrix, a row vector of z ,

$$[z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}] = w^T X + [b \ b \ \dots \ b] \quad \underbrace{\text{(1 x m) row vector}}_{\text{Can be expressed as}}$$

$$\begin{aligned} w^T \begin{bmatrix} | & | & | \\ x^{(1)} & x^{(2)} & \dots & x^{(m)} \\ | & | & | \end{bmatrix} &= w^T X = [w^T x^{(1)} \ w^T x^{(2)} \ \dots \ w^T x^{(m)}] + [b \ b \ \dots \ b] \\ &= [w^T x^{(1)} + b \ w^T x^{(2)} + b \ \dots \ w^T x^{(m)} + b] \quad \underbrace{\text{1 x m vector}}_b \end{aligned}$$

first element is exactly definition of $z^{(1)}$, second element $z^{(2)}$, and so on.

Just as X is what you obtained by stacking your training examples horizontally, we're going to define capital $Z = [z^{(1)} \ z^{(2)} \ \dots \ z^{(m)}]$

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

b is a $(1, 1)$ \mathbb{R} real number. When you do this python automatically takes the real number b and expands it out to a $1 \times m$ row vector. Called broadcasting in Python.

$$\text{Now we want: } [a^{(1)} \ a^{(2)} \ \dots \ a^{(m)}] = A = G(Z)$$

So instead of looping over M training examples we can just use these two formulas for capital Z and capital A to compute all of the lowercase z 's and a 's all at the same time. Very efficient way to compute all of the activations, a 's.



Vectorizing Logistic Regression's Gradient Computation

$$\text{we want } \rightarrow d_2^{(1)} = a^{(1)} - y^{(1)} \quad d_2^{(2)} = a^{(2)} - y^{(2)} \quad \dots \text{and so on for all training examples.}$$

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

(i, m) dims matrix

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

$$dZ = A - Y = \begin{bmatrix} a^{(1)} - y^{(1)} & a^{(2)} - y^{(2)} & \dots & a^{(m)} - y^{(m)} \end{bmatrix}$$

previously we had gotten rid of one for loop (d_{w_1}, d_{w_2}, \dots) and replaced it with just d_w . Now we still have d_{ws} and d_{bs} for each training example

$d\omega = 0$ ← Initialize →

$$dw^+ = x^{(1)} dz^{(1)}$$

$$dw^+ = x^{(2)} d\bar{z}^{(2)}$$

$$dw : x^{(m)} \mapsto z^{(m)}$$

then divide by m

$$dw / = m$$

$$d_b = 0$$

$$db+ = dz^{(1)} \\ db+ = dz^{(2)}$$

$$db + dz^{(2)} = 0$$

$$ds^2 = dz^{(m)}$$

11

$$db / = m$$

now we want to vectorize this too.

We are really just doing

$$d\omega = \frac{1}{m} (X dZ^T)$$

$$= \frac{l}{m} \begin{bmatrix} | & & | \\ x^{(1)} & \dots & x^{(m)} \\ | & & | \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ \vdots \\ dz^{(m)} \end{bmatrix}$$

\vec{t} transposed $1 \times m$ (row vector)

$$= \frac{1}{m} \left[x^{(1)} dz^{(1)} + \dots + x^{(m)} dz^{(m)} \right]_{n \times 1}$$

putting this all together:

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

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

$$A = G(\mathbb{Z})$$

$$dZ = A - Y$$

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

$$d\hat{b} = \frac{1}{m} \text{np. sum}(dZ)$$

this is both forward prop and back prop;
computes the predictions and the derivatives
for all training examples
without using a for loop

then you would perform the gradient descent update

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

$$b := b - \alpha db$$

which would be one step of gradient descent for logistic regression

You still need an outer for loop for multiple steps of gradient descent.

Broadcasting in Python

example

| | 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 | 94.0 | 0.9 |
| | | | | 59 cal |

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

Calculate % of calories from carbs, proteins, fats

$$\frac{\text{calories from carb}}{\text{total}} = \frac{56}{59} = 94.9\%$$

So what you really want to do is sum all four columns of the matrix to get total sums of cal. Then divide throughout the matrix so as to get the % of cal from C, P, F, for each of the 4 foods

The question is

Can you do this without an explicit for loop?

Yes:

Set the 3,4 matrix = A

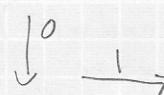
Sum the columns $\text{cal} = A.sum(\text{axis}=0)$

In second line percentage = $(A / \text{cal}).reshape(1,4) * 100$

redundant call
for clarity

$\text{sum}(\text{axis}=0)$ means you want to sum vertically

$\text{axis}=1$ ← sums horizontally



reshape is a constant time call so don't hesitate to use it

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

We use this type of broadcasting with parameter b.
Broadcasting works with both column and row vectors.

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

2×3
(m, n)

$(1, n) \xrightarrow{\text{python}} (m, n)$

copies the matrix m times

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

(m, n)

$(m, 1)$

\downarrow copies this horizontally 3 times
 (m, n)

General principle

(m, n)
matrix

if you add, subtract, divide $\frac{+}{\times}$ or multiply then
performing element-wise

conversely if you have $(m, 1) \rightarrow (m, n)$

if you have a column vector, and perform these operations against a real number

$$\begin{bmatrix} (m, 1) \\ 1 \\ 2 \\ 3 \end{bmatrix} + 100 = \begin{bmatrix} 101 \\ 102 \\ 103 \end{bmatrix}$$

also works for row vectors

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

$(1, m)$

Read numpy docs for all broadcasting cases.

Broadcasting in python is somewhat like bsxfun in Matlab/Octave.

A note on Python / Numpy Vectors

`a = np.random.randn(5)` will give you a rank 1 python array that is neither a row or a column vector.

`a.shape = (5,)` ← These rank 1 arrays are counter-intuitive and often lead to very hard to find bugs. It is recommended that you don't use them and rather:

`a = np.random.randn(5, 1)` ← This will give you a 5×1 column vector

`a.T` transpose will behave as expected

if you're not sure of the dimensions of your vectors you can throw in asserts like

`assert(a.shape == (5, 1))` ← They're inexpensive and serve as documentation

Explanation of Logistic Regression Cost Function

the prediction $\hat{y} = g(w^T x + b)$ where $g(z) = \frac{1}{1+e^{-z}}$
 we said that we want to interpret $\hat{y} = P(y=1|x)$
 \hat{y} as the probability that $y=1$ given x

so we want our alg to output the chance that $y=1$ for a given set of input features x

$$\text{If } y=1 : P(y|x) = \hat{y}$$

$$\text{If } y=0 : P(y|x) = 1-\hat{y} \quad \leftarrow \begin{array}{l} \text{chance that} \\ y \text{ is equal to zero} \end{array}$$

We want to summarize these two equations

This equation does what we want \rightarrow

$$P(y|x) = \underbrace{\hat{y}^y}_{P} \underbrace{(1-\hat{y})^{(1-y)}}_{P}$$

If $y=1$ then the first term $= \hat{y}$, the second term $= 1 = (1-\hat{y})^{0=(1-1)}$
 $\therefore P(y|x) = \hat{y}$ ← which is exactly what we want

If $y=0$ then the first term $= 1 = \hat{y}^0$, the second term $(1-\hat{y})^{(1-0)=1}$
 $\therefore P(y|x) = 1-\hat{y}$ ← which again is what we want

The above equation is a correct definition for $P(y|x)$

Now, because the log function is a strictly monotonically increasing function, maximizing $\log P(y|x)$ should give you a similar result as optimizing $P(y|x)$ and if you compute $\log P(y|x) =$

$$= \log \hat{y}^y (1-\hat{y})^{(1-y)}$$

$$\begin{aligned} & \text{simplifies to} & & = y \log \hat{y} + (1-y) \log (1-\hat{y}) \\ & & & = -\mathcal{L}(\hat{y}, y) \end{aligned}$$

this is ^T negative because we want to minimize the loss function
 minimizing the loss corresponds to maximizing the probability

Cost on m examples

identically independently distributed

$$P(\text{Labels in training set}) = \prod_{i=1}^m P(y^{(i)} | x^{(i)})$$

You want to carry out maximum likelihood estimation. Maximizing the above is the same as maximizing the log of both sides

$$\log P(\dots) = \sum_{i=1}^m \log \underbrace{P(y^{(i)} | x^{(i)})}_{-\mathcal{L}(\hat{y}^{(i)}, y^{(i)})}$$

a principle of stats called max likelihood estimation means to choose parameters that maximize

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

This justifies the cost we had for logistic regression

By minimizing the cost function J we're really carrying out maximum likelihood estimation with the logistic regression model under the assumption that our training examples were identically independently distributed.

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

This explains why we use the cost function that we do for logistic regression