

Binary Classification:

Logistic regression is an algorithm for binary classification

How does our computer store images?

- If we have a 64×64 pixel image, our computer stores 3 red-green-blue matrices that are 64×64 (rows x columns) corresponding to three colors pixel intensity values.

- A feature vector made from these matrices would be a single column with $64 \times 64 \times 3 = 12288$ values.

\sim Dimension of the feature vector
 $x = w = 12288$ (represents the input features)

- In a binary image classifier, our goal is to be able to input a feature vector (x) and get a 1 or a 0 as an output, where 1 represents cat or non-cat respectively.

Unwrapping an image into a feature vector means converting the RGB values into a suitable feature vector

Notation

- Single training example: (x, y)
- Where $x \in \mathbb{R}^n$, $y \in \{0, 1\}$
- Element of an n -dimensional feature vector

- m training examples (training sets): $\{(x^1, y^1), (x^2, y^2), \dots, (x^m, y^m)\}$
- $m = m_{train}$
- $m_{test} = m_{test}$ examples

- Abbreviated way to write matrices:

$$\text{Capital } X = \left[\begin{array}{c|c|c} & & \\ \hline x^{(1)} & x^{(2)} & x^{(3)} \\ \hline & \cdots & \\ \hline \end{array} \right] \quad \begin{matrix} n \times \text{feature vectors} \\ \text{columns} \end{matrix}$$

$X \in \mathbb{R}^{n \times m}$

$\sim \# \text{ of training examples}$

\sim in Python this comes out to be:
 X shape (n, m)

$n \times m$
 dimensional matrix

- abbreviated for the output:

$$\mathbf{y} = [y^{(1)} \ y^{(2)} \ y^{(3)} \dots \ y^{(m)}]$$

$$\mathbf{y} \in \mathbb{R}^{1 \times m}$$

L in Python: $\mathbf{y}.shape = (1, m)$

Logistic Regression:

columns
matrix

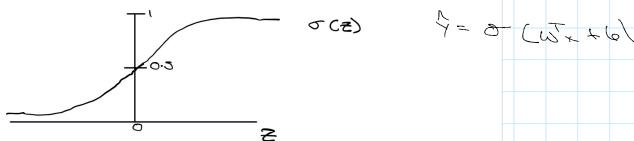
- Given x , want $\hat{y} = P(Y=1|x)$
If x is the car feature, we want \hat{y} to tell us the probability that x is a car feature.
- $x \in \mathbb{R}^n$
 x is an n -dimensional vector given that:
parameters: $w \in \mathbb{R}^n$, $b \in \mathbb{R}$ ~ real-number
 y ~ binary
- Given an input x and the parameters w and b how do we generate \hat{y} ?
- traditional linear regression: $\hat{y} = w^T x + b$
This does not work for binary classification

we want a value between 0 and 1

because you want \hat{y} to be the chance that $y=1$

This value can be greater than 1 or negative which doesn't make sense for a probability

- sigmoid function σ for binary classification



$$\sigma(z) = \frac{1}{1+e^{-z}}$$

If z is large $\sigma(z) = \frac{1}{1+e^{-z}} \approx 1$ ~ the exponential factor

If z is small $\sigma(z) = \frac{1}{1+e^{-z}} \approx 0$ ~ the exponential factor

When programming a neural net, it is better to keep w and b as separate parameters and not in the same vector (as in other

the exponential factor
decays
grows
making
the denominator
large

covered

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

Logistic Regression Cost Function

- a cost function is used to train the parameters ω and b of logistic regression

logistic

$$\bullet \hat{y} = \sigma(\omega^T x + b) \quad \sigma = \frac{1}{1 + e^{-x}}$$

$\hat{y}^{(i)}$ describes data associated with the i -th training example

- loss (error) function:

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

This is not done in logistic regression because when you come to learn the parameters, you learn that the optimization problem becomes non-convex.

Local and Global Optima

- the loss function is used to define how good our \hat{y} is when the true label is

To \hat{y} it measures the performance of our predictions on the given

loss function used in logistic regression: J(\hat{y})

$$J(\hat{y}, y) = -y \log(\hat{y}) + (1-y) \log(1-\hat{y})$$

This is used in optimization

why this makes sense? because which is

If $y=1$: $J(\hat{y}, y) = -\log \hat{y}$

want $\log \hat{y}$ large, want \hat{y} large

if $y=0$: $J(\hat{y}, y) = -\log(1-\hat{y})$

want $\log(1-\hat{y})$ large, want \hat{y} small

so that the loss is low

- Cost Function Convex Function

$$J(\omega, b) = \frac{1}{m} \sum_{i=1}^m J(\hat{y}^{(i)}, y^{(i)})$$

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

The cost function is the cost of your parameters

Convex

4:15

We want $-\gamma \log(\hat{y}-y)$

Q: Small as possible

• when

training set

$\{(x^{(1)}, y^{(1)}), \dots, (x^{(m)}, y^{(m)})\}$

We want function prediction

$\hat{y}^{(i)}$

approximately to true

true label $y^{(i)}$

• for the

loss function

we don't minimize the square error

because it results

in a non-convex

shape where gradient

descent can't

find the global

optima

Local

We want this error

to be as small as

possible

$J(\hat{y}) \geq 0$

$J(\hat{y}) = 0.0001$ is large loss

negative

Used to make loss positive as \hat{y} moves 0-1 which would give negative loss

Now what we do for the entire training set

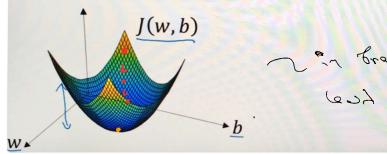
- in training the model, we are trying to find parameters w and b such that the cost function is minimized

- **Discussion:** The loss function computes the error for a single training example. The cost function is the average loss function for the whole training set.

Gradient Descent

- the gradient descent algorithm is used to train/learn the w and b parameters on training sets. ~ making the $J(w, b)$ cost function as small as possible

Want to find w, b that minimize $J(w, b)$



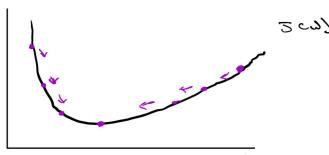
~ in practice, we can have higher dimensional loss functions than just one for n-dimensional feature.

We want to find the value (w, b) such that we minimize the cost function J . **Correct** (minimum)

(min) constraint is a straight downward convex function

- the fact that the shape of our cost function is convex, as seen above, is one of the huge reasons we like it for logistic regression.

- for gradient descent we start at a point and work our way down until we reach the optimum value. **Step size / learning rate**



direction (lower rate of change (downward))

$$\text{Update } \tilde{w} : \tilde{w} = w - \left[\frac{\partial J(w)}{\partial w} \right]$$

we do this until the algorithm converges

α : learning rate: controls how big a step we take each iteration

will be explosive

$\alpha < 1$: it will never stop

Writing code $w := w - \Delta w$

```

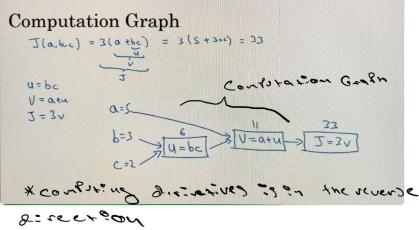
for gradient descent J(w,b)
w := w - α  $\frac{\partial J(w,b)}{\partial w} \approx \Delta w$ 
b := b - α  $\frac{\partial J(w,b)}{\partial b} \approx \Delta b$ 
* in code

```

* 2 lines
intuitive
symbolic vs
when differentiating
in respect to
2 variables.

Computation Graph

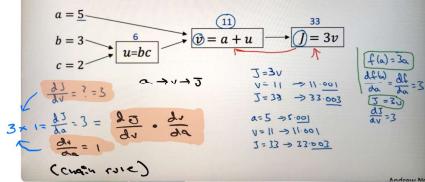
- The way neural networks are organized is a forward propagation step in which the outputs are calculated followed by a backward propagation step used to calculate gradients / derivatives to fine tune parameters.



right to left pass

Computing Derivatives

Computing derivatives

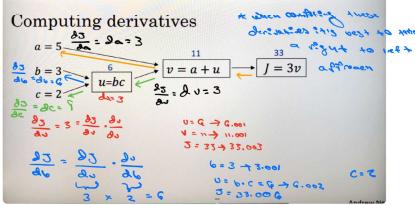


- By changing all those changes first from v to a and therefore to J

- Final output variable (J) = $\frac{\partial J}{\partial v}$ var ("d v")
* Python notation

$$\frac{\partial J}{\partial v} \sim \text{code syntax}$$

name \rightarrow context \rightarrow and



Logistic Regression Gradient Descent

$$\begin{aligned}
 & z = w_1x_1 + w_2x_2 + b \\
 & \hat{y} = \sigma(z) = \frac{1}{1+e^{-z}} \\
 & \text{Loss function (example)} \\
 & \frac{\partial \text{loss}}{\partial z} = -\frac{y}{\hat{y}} + \frac{1-y}{1-\hat{y}}
 \end{aligned}$$

In logistic regression we want to update these parameters (w_1, w_2) in order to reduce the loss function $J(w_1, w_2, b)$

Gradient descent on m examples

$$\begin{aligned}
 \text{Cost function: } J(w, b) &= \frac{1}{m} \sum_{i=1}^m \text{loss}(z^{(i)}, y^{(i)}) \\
 \rightarrow z^{(i)} &= w^T x^{(i)} + b = w_1 x_1^{(i)} + w_2 x_2^{(i)} + b \\
 \frac{\partial}{\partial w_1} J(w, b) &= \frac{1}{m} \sum_{i=1}^m \frac{\partial}{\partial w_1} \text{loss}(z^{(i)}, y^{(i)}) \\
 &\quad \downarrow \text{one training example}
 \end{aligned}$$

This would give us the overall gradient when we can sum up all gradients.

The Overall Gradient $\frac{\partial J}{\partial w_1}$

if what we computed before

$$\frac{\partial \text{loss}}{\partial z} = \frac{\partial z}{\partial w_1} \cdot \frac{\partial \text{loss}}{\partial z} = x_1 \cdot \frac{\partial \text{loss}}{\partial z}$$

$$\frac{\partial \text{loss}}{\partial w_2} = \frac{\partial z}{\partial w_2} \cdot \frac{\partial \text{loss}}{\partial z} = x_2 \cdot \frac{\partial \text{loss}}{\partial z}$$

$$\frac{\partial \text{loss}}{\partial b} = \frac{\partial z}{\partial b} \cdot \frac{\partial \text{loss}}{\partial z} = 1 \cdot \frac{\partial \text{loss}}{\partial z}$$

Update by averaging the loss function for all m different training examples

use to implement gradient descent

Concrete algorithm of this...

$$J=0, \frac{\partial J}{\partial w_1}=0, \frac{\partial J}{\partial w_2}=0, \frac{\partial J}{\partial b}=0 \quad \begin{array}{l} \text{these are being used} \\ \text{as accumulators at the start} \end{array}$$

* execute for loop over training set (minimizes)

(1) For $i=1$ to m for loop over training set to calculate all derivatives for example and then add them up

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

$$\hat{y}^{(i)} = \sigma(z^{(i)}) \quad \text{Prediction}$$

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

in addition to $\frac{\partial J}{\partial w_1} = \hat{y}^{(i)} - y^{(i)}$
 one example $\frac{\partial J}{\partial w_1} += x_1^{(i)} \frac{\partial \text{loss}}{\partial z} \quad n=2$ for loop
 $\frac{\partial J}{\partial w_2} += x_2^{(i)} \frac{\partial \text{loss}}{\partial z}$
 $\frac{\partial J}{\partial b} += \frac{\partial \text{loss}}{\partial z}$

required over n iterations

$$J = \frac{1}{n} \sum_{i=1}^n J_i$$

$$(2) \left[\begin{array}{l} \frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial w_1} \\ \frac{\partial J}{\partial w_2} = \frac{\partial J}{\partial w_2} \\ \frac{\partial J}{\partial b} = \frac{\partial J}{\partial b} \end{array} \right] \quad \text{we are computing averaged so divide by } n$$

$$\text{In the end } \frac{\partial J}{\partial w_1} = \frac{\partial J}{\partial w_1}$$

$$\text{and } \begin{aligned} u_1 &:= u_1 - \alpha \partial u_1 \\ u_2 &:= u_2 - \alpha \partial u_2 \\ b &:= b - \alpha \partial b \end{aligned} \quad \text{These will be visited in the end}$$

* this concrete optimization improvement
Just one step in gradient descent so we
would need to repeat everything here
multiple times to take multiple steps of
gradient descent.

- there are two weaknesses to this way of implementing the algorithm:
 - L+20 for root explicit
 - Cutting one and another one for features / predictors
 - Lexicist + for roots in deep learning works
 - algorithms run less efficiently
 - L recolorization techniques let us get rid of these explicit for roots

Vectorization

$$\text{What is vectorization?}$$

$\vec{z} = \omega_1 x_1 + \dots + \omega_n x_n$

$\vec{w} = \begin{bmatrix} w_1 \\ \vdots \\ w_n \end{bmatrix} \times \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} \in \mathbb{R}^{n_w} \times \mathbb{R}^{n_x}$

Non-vectorized:

$\vec{z} = 0$

for $i \in \text{range}(n-w)$:

$\vec{z} = \vec{w}[i] * \vec{x}[i]$

$\vec{z} = b$

Vectorized:

$\vec{z} = \underbrace{\vec{w} \cdot \vec{x}}_{\text{dot product}} + b$

$\vec{z} = \text{np.dot}(\vec{w}, \vec{x}) + b$

\vec{z} usually library
in Python



The screenshot shows a Jupyter Notebook interface with the title "vectorization demo (Last Checkpoint: 3 minutes ago [unseen changes])". The code cell contains Python code demonstrating vectorized operations:

```
File Edit View Insert Cell Kernel Widgets Help
In [1]: %pylab inline
import numpy as np
from time import time

tic = time()
arr = np.random(1000000)
arr[::2] = arr[::2]*2
toc = time()

print("Time taken: ", toc-tic)

# Using list comprehension
tic = time()
arr = [x*2 for x in arr]
toc = time()

print("Time taken: ", toc-tic)

# Using np.multiply
tic = time()
arr = np.multiply(arr, arr)
toc = time()

print("Time taken: ", toc-tic)

# Using np.multiply_inplace
tic = time()
np.multiply_inplace(arr, arr)
toc = time()

print("Time taken: ", toc-tic)
```

- avoid explicit for and if loops whenever possible

Vectorizing logistic regression's gradient computation

Vectorizing Logistic Regression

Writing a dedication on the first page

Vectorizing Logistic Regression:

Working a reduction on the first example

forward:

$$z^{(1)} = \omega^T x^{(1)} + b$$

$$a^{(1)} = \sigma(z^{(1)})$$

Second example

$$z^{(2)} = \omega^T x^{(2)} + b$$

$$a^{(2)} = \sigma(z^{(2)})$$

}

weights have to
do this m times
for m examples.

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

columns

dimensions

training data matrix

how do we calculate all the $z^{(i)}$'s in one shot?

$$\begin{bmatrix} z^{(1)} & z^{(2)} & \dots & z^{(m)} \end{bmatrix} = \underbrace{\omega^T X}_{\text{2 row vector}} + \underbrace{\begin{bmatrix} b & b & \dots & b \end{bmatrix}}_{\text{1xm vector}} = \begin{bmatrix} \omega^T x^{(1)} + b & \omega^T x^{(2)} + b & \dots & \omega^T x^{(m)} + b \end{bmatrix}$$

\boxed{Z}

$Z_1 \quad Z_2 \quad Z_m$

* to implement this in memory,
 $Z = \text{matlab } (\omega^T, X)$ + b

(that's why we're all the same b)

→ broadcasting: transform b into appropriate
vector for computation.

how to compute $a^{(1)}$ through $a^{(m)}$,

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

* implemented in the programming assignment

Vectorizing Logistic Regression's Gradient Computation

$$\begin{aligned} \frac{\partial z^{(1)}}{\partial \omega} &= a^{(1)} - y^{(1)} \\ \frac{\partial z^{(2)}}{\partial \omega} &= a^{(2)} - y^{(2)} \end{aligned} \quad \left. \begin{array}{l} \text{how do we compute this for all } m \text{ training examples} \\ \text{this for all } m \text{ training examples} \end{array} \right\}$$

$$\frac{\partial Z}{\partial \omega} = \left[\frac{\partial z^{(1)}}{\partial \omega} \ \frac{\partial z^{(2)}}{\partial \omega} \ \dots \ \frac{\partial z^{(m)}}{\partial \omega} \right]$$

1xm
matrix

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

$$\underbrace{\qquad\qquad\qquad}_{\frac{\partial Z}{\partial \omega}} = A - Y$$

$$\frac{\partial Z}{\partial \omega} = A - Y$$

Programming Implementation

loop before

$$\Delta \omega = 0$$

$$\Delta b += x^{(1)} \Delta z^{(1)}$$

$$\Delta b += x^{(2)} \Delta z^{(2)}$$

.

$\Delta \omega$:

$$/ = m$$

$$\Delta b = 0$$

$$\Delta b += \Delta z^{(1)}$$

$$\Delta b += \Delta z^{(2)}$$

.

$$\Delta b^{(m)} += \Delta z^{(m)}$$

$$\Delta b / = m$$

Previous implementation (for loop over m examples)

New vectorized implementation

$$\partial b = \frac{1}{m} \sum_{i=1}^m \partial z^{(i)}$$

$$= \frac{1}{m} \cdot \text{np. sum}(\partial z) \quad \text{Implementation}$$

$$\partial w = \frac{1}{m} \cdot X \cdot \partial z^\top = \frac{1}{m} \cdot \begin{bmatrix} 1 & 1 & 1 & 1 \\ x^{(1)}_1 & x^{(2)}_1 & x^{(3)}_1 & \dots & x^{(m)}_1 \\ 1 & 1 & 1 & 1 \end{bmatrix} \begin{bmatrix} \partial z^{(1)} \\ \vdots \\ \partial z^{(m)} \end{bmatrix} = \underbrace{\frac{1}{m} \left[x^{(1)} \partial z^{(1)} + \dots + x^{(m)} \partial z^{(m)} \right]}_{\text{vector}} \quad \text{# of features}$$

Code Implementation of vectorization

$$z = w^\top x + b = \text{np. dot}(w, x) + b$$

$$A = \partial(z)$$

$$\partial z = A - y$$

$$\partial w = \frac{1}{m} X \cdot \partial z^\top$$

$$\partial b = \frac{1}{m} \text{np. sum}(\partial z)$$

all of $A, y, \partial z$
are calculated
as dictated
by convention

$$w := w - \alpha \partial w \quad \text{update to}$$

$$b := b - \alpha \partial b \quad \text{the weighting}$$

* this is a single iteration of gradient

descent for logistic regression

↳ if we need multiple iterations of gradient

descent, we still require update for

loop over the # of iterations