# Statistics and Artificial Intelligence

## Lecture 4: Logistic Regression as a Neural Network (II)

**Yixin Wang**

# Goals for Today

- Gradient descent

- Derivatives with a computational graph

- Logistic regression with gradient descent

- Deep Learning Frameworks: TensorFlow, Keras, Google Colab

Next week: focus on TensorFlow

# Binary Classification

**Logistic regression is for binary classification**

*Convert raw image into #'s computers can read*



- Input: Image X (Unroll pixel values into feature vector)

- Output: Label Y takes value 1 (cat) or 0 (non-cat)

*• Each image can be converted to a table*

*• Output label for each item in training data*

# Binary Classification
## Convert images into vector

- $x \in \mathbb{R}^{n_x}, y \in \{0,1\}$



Organize #s into column vector

Blue
Green
Red

| 255 | 134 | 93 | 22 |
| 255 | 134 | 202 | 22 | 2 |
| 255 | 231 | 42 | 22 | 4 | 30 |
| 123 | 94 | 83 | 2 | 192 | 124 |
| 34 | 44 | 187 | 92 | 34 | 142 |
| 34 | 76 | 232 | 124 | 94 |
| 67 | 83 | 194 | 202 |

$$x = \begin{bmatrix} 255 \\ 231 \\ \vdots \\ 255 \\ 134 \\ \vdots \end{bmatrix} \Bigg\} n_x$$

1 (cat) vs 0 (non cat)
binary labels

# Notation

- $(x, y), x \in \mathbb{R}^{n_x}, y \in \{0,1\}$ Paired observation
  image + label

- $m$ training examples: $\{(x^{(1)}, y^{(1)}, \ldots, (x^{(m)}, y^{(m)})\}$ → m training examples

- # training examples: $m = m_{\text{train}}$; similarly, # test examples: $m_{\text{test}}$

- $x \in \mathbb{R}^{n_x \times m}$, x.shape = $(n_x, m)$ convert each image into a col.

- $y \in \mathbb{R}^{1 \times m}$, y.shape = $(1, m)$

one column for each example
- every image has nx dimensions
- have nx rows

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

← m → columns

$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \cdots, & y^{(m)} \end{bmatrix}$

# Notation

- $(x, y), x \in \mathbb{R}^{n_x}, y \in \{0,1\}$

- $m$ training examples: $\{(x^{(1)}, y^{(1)}, \ldots, (x^{(m)}, y^{(m)})\}$

- # training examples: $m = m_{\text{train}}$; similarly, # test examples: $m_{\text{test}}$

- $x \in \mathbb{R}^{n_x \times m}$, x.shape = $(n_x, m)$

    Puthon

- $y \in \mathbb{R}^{1 \times m}$, y.shape = $(1, m)$

Y vector collects label
- a row, m columns
each label is just a scalar

$$X = \begin{bmatrix} | & | & & | \\ x^{(1)} & x^{(2)} & \ldots, & x^{(m)} \\ | & | & & | \end{bmatrix} \updownarrow n_x$$

$$\longleftarrow m \longrightarrow$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & \ldots & , y^{(m)} \end{bmatrix}$$

# Logistic Regression

- Given $x$, want $\hat{y} = P(y = 1 \mid x),$  → $\hat{y}$ is the probability its labelled a, given the image

- $x \in \mathbb{R}^{n_x}, 0 \leq \hat{y} \leq 1$

- Parameters: $w \in \mathbb{R}^{n_x}, b \in \mathbb{R}$

Find values of w + b that make $\hat{y} \to y$

sigmoid

- Output: $\hat{y} = \sigma(\sum_{j=1}^{n_x} w_j x_j + b) = \sigma(w^{\top} x + b)$

sum    weight for every #, add bias



sigmoid takes # to a, small # to 0

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

- If $z$ large, $\sigma(z) \approx \dfrac{1}{1 + 0} = 1$; if $z$ large negative number, $\sigma(z) \approx \dfrac{1}{1 + \text{big-number}} \approx 0$

Sigmoid

# Logistic Regression
## Scalar and matrix version

*$\hat{y}$ very close to $y \Rightarrow$ loss function is small*

*$\hat{y}$ far from $y \Rightarrow$ loss function very large*

*• continue to minimize*

- $d$-dimensional input

- $\hat{y}^{(i)} = \sigma(\sum_{j=1}^{n_x} w_j x_j^{(i)} + b) = \sigma(w^\top \mathbf{x}^{(i)} + b)$ where $\sigma(z^{(i)}) = \dfrac{1}{1 + e^{-z^{(i)}}}$

*convex $\Rightarrow$ easier to minimize*

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

- Loss function / error function: $L(\hat{y}, y) = -(y \log \hat{y} + (1 - y)\log(1 - \hat{y}))$ (<u>convex</u>)

  *Cross entropy loss*

- Cost function:

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

*look at predictions vs labels across whole dataset*

*min*

*convex*

# Poll

- What are the parameters of logistic regression?

    - W, an $n_x$ dimensional vector, and b, a real number.

    - W, an identity vector, and b, a real number.

    - W and b, both real numbers.

    - W and b, both $n_x$ dimensional vectors.

$W \Rightarrow$ vector of weights

# Logistic Regression Cost Function
## Cost function / Loss function

- Loss function measures how good the prediction is, relative to the true label.

- Loss function / error function (defined for a **single** training example):
$$L(\hat{y}, y) = -(y \log \hat{y} + (1 - y)\log(1 - \hat{y}))$$

*During training we minimize the cost function*

- Training: Make the cost function as small as possible

- Cost function (defined for the **whole** training set): cost for parameters:

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

# Logistic Regression Cost Function
## Cost function / Loss function

- Loss function measures how good the prediction is, relative to the true label.

- Loss function / error function (defined for a **single** training example):
$$L(\hat{y}, y) = -(y \log \hat{y} + (1-y)\log(1-\hat{y}))$$

- Why this loss function? Convexity!

L(w)

1
0

we start somewhere + try to
move towards minimum

· Have to decide when we get
to min

· easier w/ convex shape

· don't want to end in local
minima

w

# Gradient Descent for Training Neural Network

**Logistic regression** Find w + b such that cost function is minimized

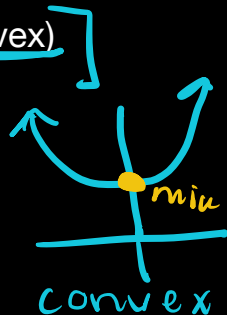- Recap: $\hat{y}^{(i)} = \sigma(\sum_{j=1}^{n_x} w_j x_j^{(i)} + b) = \sigma(w^\top \mathbf{x}^{(i)} + b)$ where $\sigma(z^{(i)}) = \dfrac{1}{1 + e^{-z^{(i)}}}$

- $\underbrace{J(w,b)}_{\text{cost}} = \dfrac{1}{m}\sum_{i=1}^{m} L(\hat{y}^{(i)}, y^{(i)}) = -\dfrac{1}{m}\sum_{i=1}^{m}(y^{(i)}\log\hat{y}^{(i)} + (1 - y^{(i)})\log(1 - \hat{y}^{(i)}))$

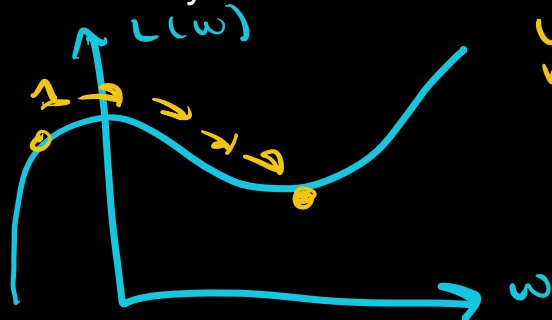

initialize anywhere in space
• calculate gradient + go down
Follow as we move towards a flat gradient
• end at global minimum

https://math.stackexchange.com/questions/1582452/logistic-regression-prove-that-the-cost-function-is-convex

# Gradient Descent for Training Neural Network

- Gradient Descent:

  - Initialize; *Figure this out based on the derivative*

  - Take a step in the steepest downhill direction at each iteration

  - Repeat until the algorithm converges

    *get to global min*



$J(w, b)$ plotted over $w$ and $b$

# Gradient Descent

$J(w)$

derivative

Set w to some # (0 or random)

Repeat

$$w \leftarrow w - \alpha \frac{dJ(w)}{dw}$$

learning rate

converge

$$\frac{dJ(w)}{dw}$$

calculate derivative
at each point
→ move new point
→ calculate new point

Learning rate decides
how much we move
at each step

we just trace the curve
downward

$$J(w, B) \quad w \leftarrow w - \alpha \frac{dJ(w, b)}{dw}$$

$$b \leftarrow b - \alpha \frac{dJ(w, b)}{dw}$$

$$\frac{\partial J(w, b)}{\partial w}$$ → dw

Partial
Notation,
means
same thing

$$\frac{\partial J(w, b)}{\partial b}$$ → db

Preliminary Draft.
Please do not distribute.

# Gradient Descent

*Almost all training algorithms are variants of gradient descent*

- Learning rate: Control how big a step we take at each iteration

  *How we trace curve*
- Derivative: <u>Slope</u> of the function; the direction to go downhill

- Code: Often use '<u>dw</u>' to represent the derivative dJ(w)/dw

  - The amount you want to update for w

# Question

- True or false. A convex function always has multiple local optima.

# Intuition about Derivatives

$f(a) = 3a$

- Slope / Derivative: If I nudge $a$ by a tiny little bit, I expect $f(a)$ to move three times as large as the nudge I gave $a$.



$$\frac{df(a)}{da} = 3$$

# Intuition about Derivatives

$f(a) = 3a$

- Formal definition through infinitesimal nudge.

- This function $f(a) = 3a$ has constant slope.

# Intuition about Derivatives

$f(a) = a^2$

# More Derivative Examples

How does the value change if we nudge
the function

# Computation Graph

- The computations of a neural network are organized in terms of

  - a forward pass or a forward propagation step, in which we compute the output of the neural network,

  - followed by a backward pass or back propagation step, which we use to compute gradients or compute derivatives.

- The computation graph explains why it is organized this way.

# Computational Graph

## Left-to-right pass

$J(a, b, c) = 3(a + bc)$

$a = 5$
$b = 3$
$c = 2$

$J(a, b, c) = 33$

$\underbrace{\phantom{a + bc}}_{u}$
$\underbrace{\phantom{a + bc}}_{v}$
$\underbrace{\phantom{a + bc}}_{J}$

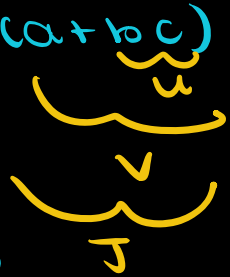Computers approach this calculation more "step by step"

$u = bc$
$v = a + u$
$J = 3 \cdot v$

$\}$ computational graph

every neural net has a computational graph

$a = 5 \longrightarrow$

$b = 3 \longrightarrow$

$c = 2 \longrightarrow$

$u = bc$

6

$v = a + u$

11

$\longrightarrow$

$J = 3v$

33

# Computing Derivatives with a Computation Graph

## Right-to-left pass (Chain rule) $a \rightarrow v \rightarrow J$

- Backpropagation $\dfrac{\mathrm{d}J}{\mathrm{d}a}$



$a = 5$

$b = 3$

$c = 2$

$u = bc$ — 6

$v = a + u$ — 11

$J = 3v$ — 33

# Computing Derivatives with a Computation Graph

## Right-to-left pass (Chain rule) $b \to J, c \to J$

- Backpropagation: $\dfrac{\mathrm{d}J}{\mathrm{d}b}$, and $\dfrac{\mathrm{d}J}{\mathrm{d}c}$

# Questions?

# Logistic Regression Gradient Descent
## Logistic regression recap

- $z = \sum\limits_{j}^{n_x} w_j x_j + b, w^\top x + b$

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

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

# Logistic Regression Gradient Descent
## Logistic regression derivatives

# Logistic Regression Gradient Descent

## Logistic regression on m examples

- $J(w, b) = \dfrac{1}{m} \sum\limits_{i=1}^{m} L(a^{(i)}, y^{(i)})$

- $a^{(i)} = \hat{y}^{(i)} = \sigma(z^{(i)} = \sigma(w^{\top} x^{(i)} + b)$

# Logistic Regression Gradient Descent
## Logistic regression on m examples

$J = 0$ ; $dw_1 = 0$ ; $dw_2 = 0$ ; $db = 0$

For $i = 1$ to $m$

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

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

$\quad J += -\left[ y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)}) \right]$

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

$\quad dw_1 += x_1^{(i)} dz^{(i)}$

$\quad dw_2 += x_2^{(i)} dz^{(i)}$  $\quad n = 2$

$dw_3$
$\vdots$
$dw_n$

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

# Questions?

# Deep Learning Frameworks

- Caffe/Caffe2

- CNTK

- DL4J

- Keras

- Lasagne

- mxnet

- PaddlePaddle

- TensorFlow

- Theano

- Torch

# Deep Learning Frameworks

- Choosing deep learning frameworks

- - Easeofprogramming(development and deployment)

- - Runningspeed

- - Truly open (open source with good governance)

The materials in this course are adapted from materials created by Alexander Amini, Alfredo Canziani, Justin Johnson, Andrew Ng, Bhiksha Raj, Grant Sanderson and the 3blue1brown channel, Rita Singh, Ava Soleimany, and Ambuj Tewari.