

Logistic Regression

Sargur N. Srihari

University at Buffalo, State University of New York
USA

Topics in Linear Classification using Probabilistic Discriminative Models

- Generative vs Discriminative
 1. Fixed basis functions
 2. Logistic Regression (two-class)
 3. Iterative Reweighted Least Squares (IRLS)
 4. Multiclass Logistic Regression
 5. Probit Regression
 6. Canonical Link Functions

Topics in Logistic Regression

- Logistic Sigmoid and Logit Functions
- Parameters in discriminative approach
- Determining logistic regression parameters
 - Error function
 - Gradient of error function
 - Simple sequential algorithm
 - An example
- Generative vs Discriminative Training
 - Naïve Bayes vs Logistic Regression

Logistic Sigmoid and Logit Functions

- In two-class case, *posterior* of class C_1 can be written as as a logistic sigmoid of feature vector $\phi = [\phi_1, \dots, \phi_M]^T$

$$p(C_1|\phi) = y(\phi) = \sigma(w^T\phi)$$

with $p(C_2|\phi) = 1 - p(C_1|\phi)$

Here $\sigma(\cdot)$ is the logistic sigmoid function

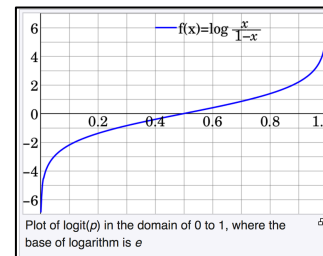
– Known as logistic regression^w in statistics

- Although a model for classification rather than for regression

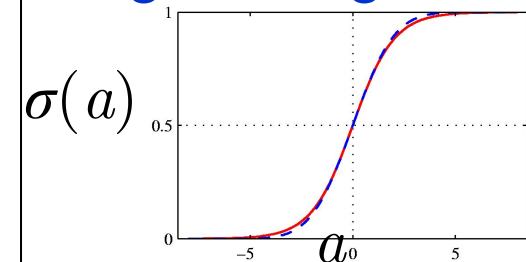
- Logit function:

– It is the log of the odds ratio

- It links the probability to the predictor variables



Logistic Sigmoid



Properties:

A. Symmetry

$$\sigma(-a) = 1 - \sigma(a)$$

B. Inverse

$$a = \ln(\sigma / (1 - \sigma))$$

known as *logit*.

Also known as *log odds* since it is the ratio

$$\ln[p(C_1|\phi) / p(C_2|\phi)]$$

C. Derivative

$$d\sigma / da = \sigma(1 - \sigma)$$

Fewer Parameters in Linear Discriminative Model

- Discriminative approach (Logistic Regression)
 - For M -dim feature space ϕ :
 - M adjustable parameters
- Generative based on Gaussians (Bayes/NB)
 - $2M$ parameters for mean
 - $M(M+1)/2$ parameters for shared covariance matrix
 - Two class priors
 - Total of $M(M+5)/2 + 1$ parameters
 - Grows quadratically with M
 - If features assumed independent (naïve Bayes) still ₅ needs $M+3$ parameters

Determining Logistic Regression parameters

- Maximum Likelihood Approach for Two classes
- For a data set (ϕ_n, t_n) where $t_n \in \{0, 1\}$ and $\phi_n = \phi(\mathbf{x}_n)$, $n = 1, \dots, N$
- Likelihood function can be written as

$$p(\mathbf{t} | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

where $\mathbf{t} = (t_1, \dots, t_N)^T$ and $y_n = p(C_1 | \phi_n)$

y_n is the probability that $t_n = 1$

Error Fn for Logistic Regression

- Likelihood function is

$$p(t | \mathbf{w}) = \prod_{n=1}^N y_n^{t_n} \{1 - y_n\}^{1-t_n}$$

- By taking negative logarithm we get the *Cross-entropy Error Function*

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

where $y_n = \sigma(a_n)$ and $a_n = \mathbf{w}^T \boldsymbol{\phi}_n$

- We need to minimize $E(\mathbf{w})$

At its minimum, derivative of $E(\mathbf{w})$ is zero

So we need to solve for \mathbf{w} in the equation

$$\nabla E(\mathbf{w}) = 0$$

What is Cross-entropy?

- Entropy of $p(x)$ is defined as $H(p) = -\sum_x p(x) \log p(x)$

– If $p(x=1|t)=t$ and $p(x=0|t)=1-t$ then we can write

$$p(x) = t^x (1-t)^{1-x}$$

- Then Entropy of $p(x)$ is $H(p) = t \log t + (1-t) \log (1-t)$

- Cross entropy of $p(x)$ and $q(x)$ is defined as

$$H(p, q) = -\sum_x p(x) \log q(x)$$

– If $q(x=1|y)=y$ then $H(p, q) = t \log y + (1-t) \log (1-y)$

- In general $H(p, q) = H(p) + D_{KL}(p||q)$

– where

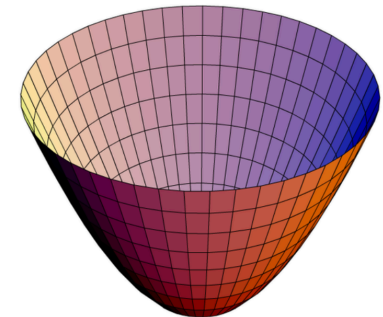
$$D_{KL}(p, q) = -\sum_x p(x) \log \frac{p(x)}{q(x)}$$

Gradient of Error Function

Error function

$$E(\mathbf{w}) = -\ln p(t | \mathbf{w}) = -\sum_{n=1}^N \left\{ t_n \ln y_n + (1 - t_n) \ln(1 - y_n) \right\}$$

where $y_n = \sigma(\mathbf{w}^T \phi_n)$



Using Derivative of logistic sigmoid $d\sigma/da = \sigma(1-\sigma)$

Gradient of the error function

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N \underbrace{(y_n - t_n) \phi_n}_{\text{Error} \times \text{Feature Vector}}$$

Error x Feature Vector

Contribution to gradient by data point n is error between target t_n and prediction $y_n = \sigma(\mathbf{w}^T \phi_n)$ times basis ϕ_n

Proof of gradient expression

Let $z = z_1 + z_2$

where $z_1 = t \ln \sigma(w\phi)$ and $z_2 = (1 - t) \ln[1 - \sigma(w\phi)]$

$$\frac{dz_1}{dw} = \frac{t \sigma(w\phi) [1 - \sigma(w\phi)] \phi}{\sigma(w\phi)}$$

$$\frac{d\sigma}{da} = \sigma(1 - \sigma)$$

and

Using $\frac{d}{dx}(\ln ax) = \frac{a}{x}$

$$\frac{dz_2}{dw} = \frac{(1 - t) \sigma(w\phi) [1 - \sigma(w\phi)] (-\phi)}{[1 - \sigma(w\phi)]}$$

9

Therefore $\frac{dz}{dw} = (\sigma(w\phi) - t) \phi$

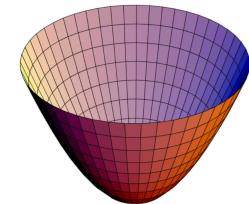
Simple Sequential Algorithm

- Given Gradient of error function

$$\nabla E(\mathbf{w}) = \sum_{n=1}^N (y_n - t_n) \phi_n \quad \text{where } y_n = \sigma(\mathbf{w}^T \phi_n)$$

- Solve using an iterative approach

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_n$$



- where

$$\nabla E_n = \underbrace{(y_n - t_n) \phi_n}$$

Error x Feature Vector

Takes precisely same form as
Gradient of Sum-of-squares
error for linear regression

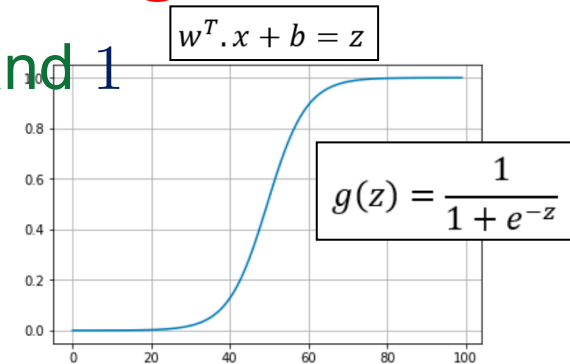
Samples are presented one at a time in
which each each of the weight vectors is updated

Python Code for Logistic Regression

Sigmoid function to produce value between 0 and 1

```
def sigmoid(z):
    return (1 / (1 + np.exp(-z)))
```

Prediction
 $s(z) = p$



Loss and Cost function

Loss function is the loss for a training example

Cost is the loss for whole training set

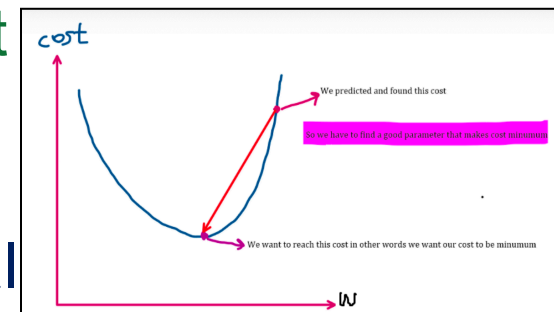
$$L(p, y) = -(y \log p + (1 - y) \log(1 - p))$$

Updating weights and biases

p is our prediction and y is correct value

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

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



Finding db and dw

Derivative wrt $p \rightarrow$ Derivative wrt z .

$$\frac{d}{dx} \log(mx) = \frac{(mx)'}{mx} = \frac{m}{mx} = \frac{1}{x}$$

$$\frac{\partial L(p, y)}{\partial p} = \frac{\partial}{\partial p} (-(y \log p + (1 - y) \log(1 - p)))$$

$$\frac{\partial L(p, y)}{\partial p} = \left(\frac{1 - y}{1 - p} - \frac{y}{p} \right)$$

$$\frac{\partial L(p, y)}{\partial z} = \frac{\partial L(p, y)}{\partial p} \frac{\partial p}{\partial z} \quad \dots \text{chain rule}$$

$$\frac{\partial p}{\partial z} = p(1 - p) \quad \text{because of the sigmoid function}$$

$$\frac{\partial L(p, y)}{\partial z} = \left(\frac{1 - y}{1 - p} - \frac{y}{p} \right) (p(1 - p))$$

$$\frac{\partial L(p, y)}{\partial z} = p - y$$

<https://towardsdatascience.com/logistic-regression-from-very-scratch-ea914961f320>

$$\frac{\partial L(p, y)}{\partial b} = \frac{\partial L(p, y)}{\partial z} \frac{\partial z}{\partial b} = dz \frac{\partial}{\partial b} (w^T \cdot x + b)$$

$$\frac{\partial L(p, y)}{\partial b} = db = \frac{1}{\underbrace{m}_{\text{training example}}} \sum_{k=1}^m dz^k$$

db

$$\frac{\partial L(p, y)}{\partial w} = \frac{\partial L(p, y)}{\partial z} \frac{\partial z}{\partial w} = dz \frac{\partial}{\partial w} (w^T \cdot x + b)$$

$$\frac{\partial L(p, y)}{\partial w} = dw = \frac{1}{\underbrace{m}_{\text{training example}}} X dz^T$$

dw

Logistic Regression Code in Python

use sci-kit learn to create a data set.

```
import sklearn.datasets
import matplotlib.pyplot as plt
import numpy as np
X, Y = sklearn.datasets.make_moons(n_samples=500, noise=.2)
X, Y = X.T, Y.reshape(1, Y.shape[0])
epochs = 1000
learningrate = 0.01
def sigmoid(z):
    return 1 / (1 + np.exp(-z))
losstrack = []
m = X.shape[1]
w = np.random.randn(X.shape[0], 1)*0.01
b = 0
for epoch in range(epochs):
    z = np.dot(w.T, X) + b
    p = sigmoid(z)
    cost = -np.sum(np.multiply(np.log(p), Y) + np.multiply((1 - Y), np.log(1 - p)))/m
    losstrack.append(np.squeeze(cost))
    dz = p-Y
    dw = (1 / m) * np.dot(X, dz.T)
    db = (1 / m) * np.sum(dz)
    w = w - learningrate * dw
    b = b - learningrate * db
plt.plot(losstrack)
```

for epoch in range(epochs):

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

$$p = s(z)$$

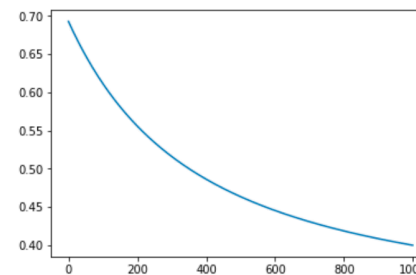
$$dz = p - y$$

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

$$db = \frac{1}{m} \sum_{k=1}^m dz^k$$

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

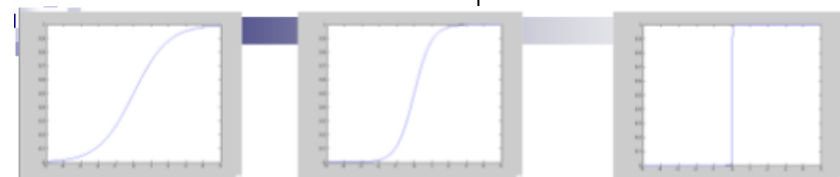
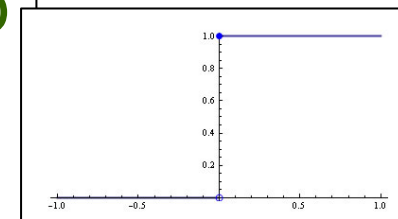
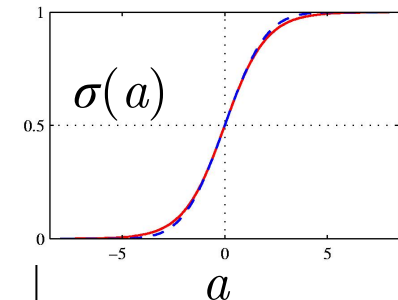
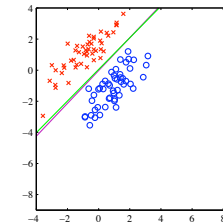
$$w := w - \alpha dw$$



Prediction: From the code above, you find p . It will be between 0 and 1.

ML solution can over-fit

- Severe over-fitting for linearly separable data
 - Because ML solution occurs at $\sigma = 0.5$
 - With $\sigma > 0.5$ and $\sigma < 0.5$ for the two classes
 - Solution equivalent to $a = w^T \phi = 0$
 - Logistic sigmoid becomes infinitely steep
 - A Heavyside step function $\|w\|$ goes to ∞
 - Solution
 - Penalizing wts
 - Recall in linear regression



$$\nabla E_n = -\sum_{n=1}^N \left\{ t_n - w^T \phi(x_n) \right\} \phi(x_n)^T \quad \text{without reg}$$

$$\nabla E_n = \left[-\sum_{n=1}^N \left\{ t_n - w^T \phi(x_n) \right\} \phi(x_n)^T \right] + \lambda w \quad \text{with reg}$$

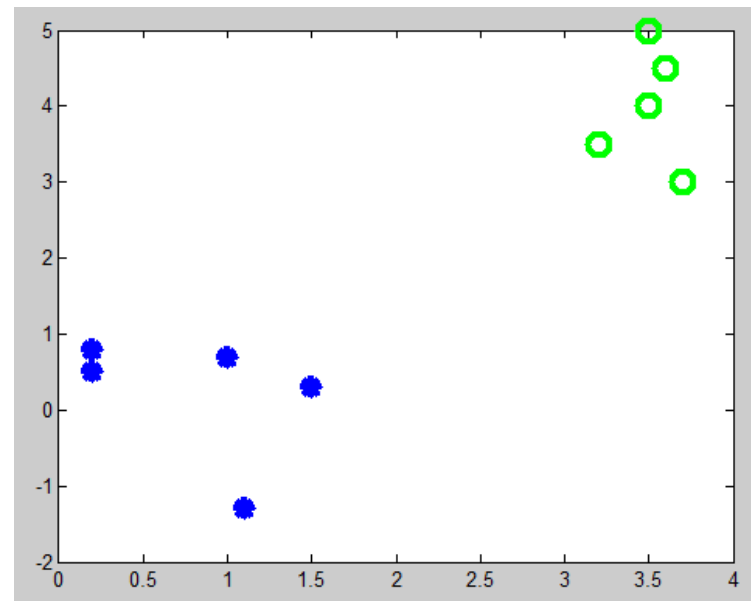
$$\frac{1}{1 + e^{-x}} \quad \frac{1}{1 + e^{-2x}} \quad \frac{1}{1 + e^{-100x}}$$

An Example of 2-class Logistic Regression

- Input Data

C1 =	
3.7000	3.0000
3.2000	3.5000
3.5000	5.0000
3.6000	4.5000
3.5000	4.0000

C2 =	
1.1000	-1.3000
0.2000	0.5000
1.5000	0.3000
0.2000	0.8000
1.0000	0.7000



$\phi_0(\mathbf{x})=1$, dummy feature

Initial Weight Vector, Gradient and Hessian (2-class)

- Weight vector

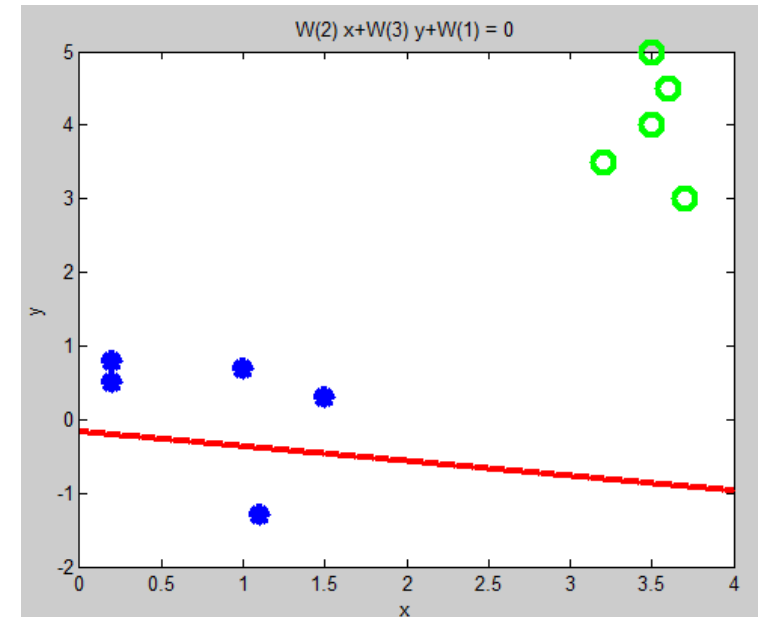
```
W =
    0.1117
    0.1363
    0.6787
```

- Gradient

```
Delta_E =
           0
    6.7500
    9.5000
```

- Hessian

```
H =
    3.5000    5.3750    5.2500
    5.3750   17.4825   17.4950
    5.2500   17.4950   22.4150
```



Final Weight Vector, Gradient and Hessian (2-class)

- Weight Vector

W =

704.5915
-20.9086
-337.6170

- Gradient

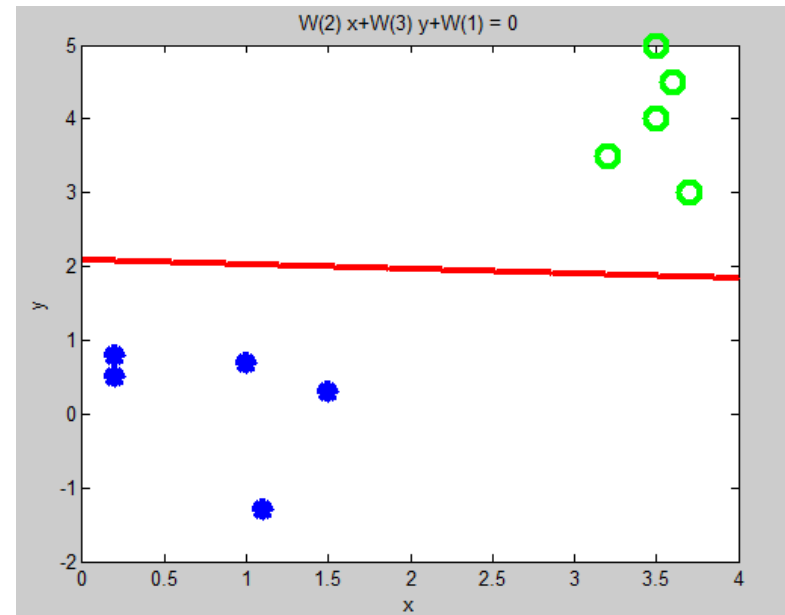
Delta_E =

-12.3917
-1.6321
4.9025

- Hessian

H =

1.0000	0.0000	0.0000
0.0000	1.0000	0.0000
0.0000	0.0000	1.0000



Number of iterations : 10

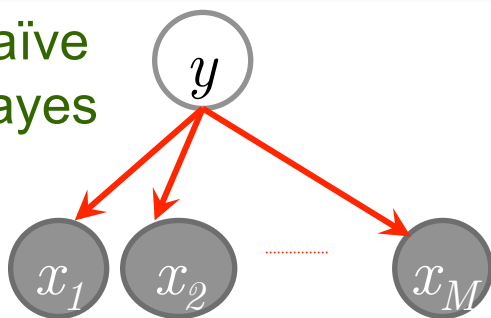
Error (Initial and Final): 15.0642, 1.0000e-009

Generative vs Discriminative Training

Variables $\mathbf{x} = \{x_1, \dots, x_M\}$ and classifier target y

1. Generative: estimate parameters of variables independently

Naïve
Bayes



For classification:
Determine joint:

$$p(y, \mathbf{x}) = p(y) \prod_{i=1}^M p(x_i | y)$$

From joint get required
conditional $p(y | \mathbf{x})$

Simple estimation

independently estimate M sets of parameters

But independence is usually false

We can estimate $M(M+1)/2$ covariance matrix

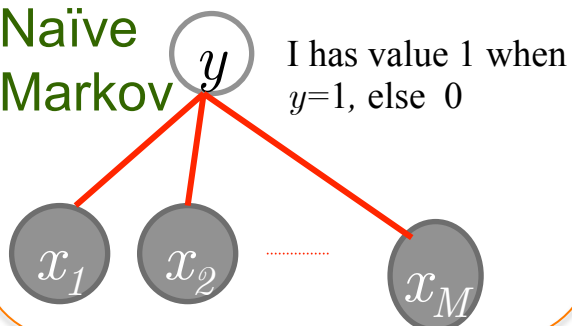
2. Discriminative: estimate joint parameters w_i

Potential Functions (log-linear)

$$\phi_i(x_i, y) = \exp\{w_i x_i \mathbb{I}\{y=1\}\},$$

$$\phi_o(y) = \exp\{w_o \mathbb{I}\{y=1\}\}$$

Naïve
Markov



For classification:

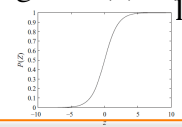
Unnormalized

$$\tilde{P}(y=1 | \mathbf{x}) = \exp\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \tilde{P}(y=0 | \mathbf{x}) = \exp\{0\} = 1$$

Normalized

$$P(y=1 | \mathbf{x}) = \text{sigmoid}\left\{w_0 + \sum_{i=1}^M w_i x_i\right\} \quad \text{where } \text{sigmoid}(z) = \frac{e^z}{1+e^z}$$

Logistic Regression



Jointly optimize M parameters

More complex estimation but correlations
accounted for

Can use much richer features:

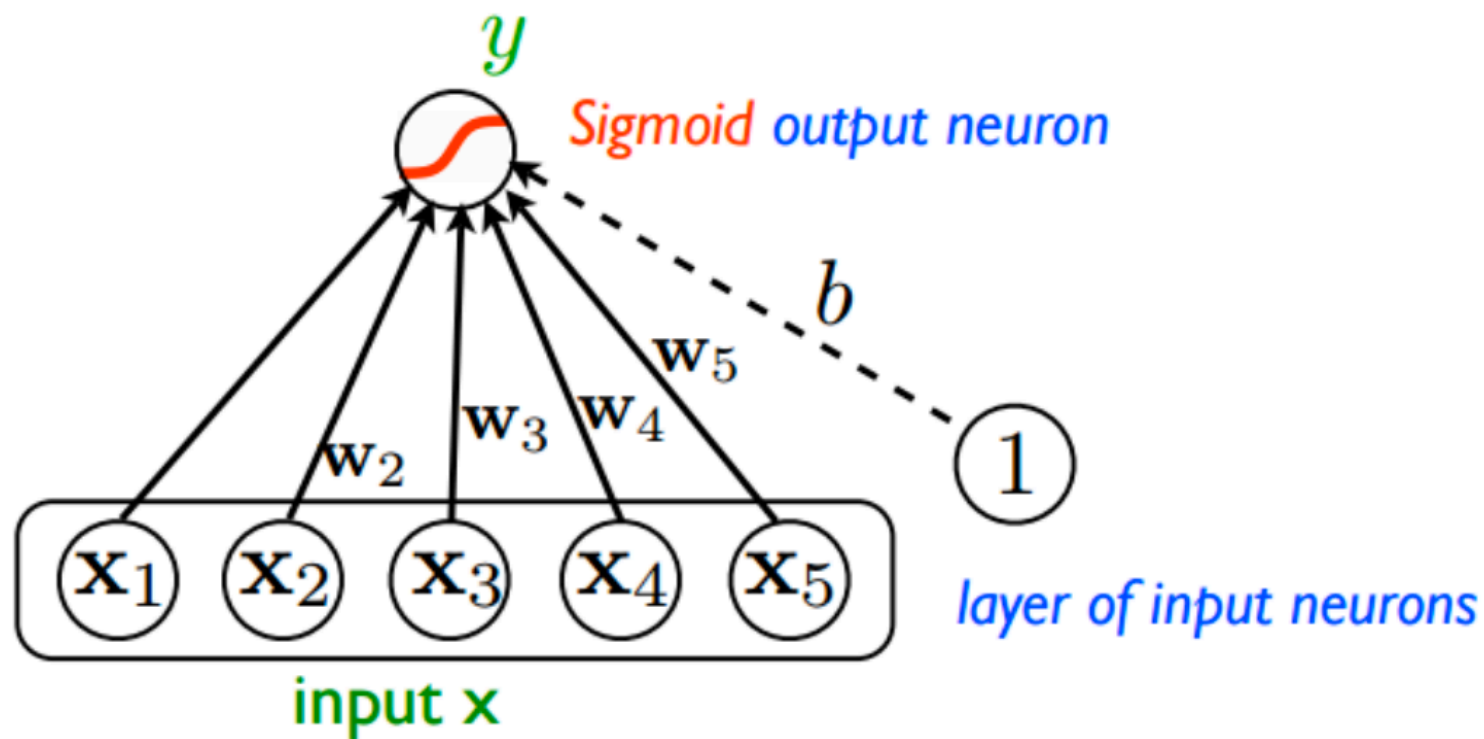
Edges, image patches sharing same pixels

multiclass

$$p(y_i | \phi) = y_i(\phi) = \frac{\exp(a_i)}{\sum_j \exp(a_j)}$$

where $a_j = \mathbf{w}_j^T \phi$

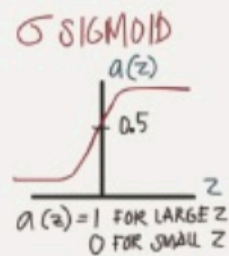
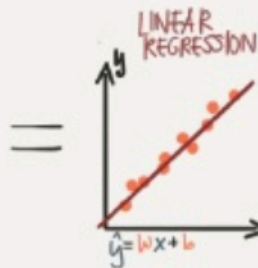
Logistic Regression is a special architecture of a neural network



BINARY CLASSIFICATION

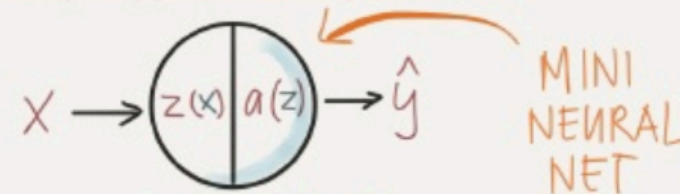


1: CAT
0: NOT CAT
 y

LOGISTIC REGRESSION
AS A NEURAL NETFINDING THE MINIMUM
WITH GRADIENT DESCENT

1. FIND THE DOWNHILL DIRECTION (USING DERIVATIVES)
 2. WALK (UPDATE w & b) AT A α LEARNING RATE
- REPEAT UNTIL YOU REACH BOTTOM (CONVERGE)

PUTTING IT ALL TOGETHER



$$z(x) = wx + b$$

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

1. FORWARD PROPAGATION • CALCULATE \hat{y}
 2. BACKWARD PROPAGATION • GRADIENT DESCENT + UPDATE w & b
- REPEAT UNTIL IT CONVERGES

THE TASK IS TO LEARN w & b BUT HOW?

A: OPTIMIZE HOW GOOD THE GUESS IS BY MINIMIZING THE DIFF BETWEEN GUESS (\hat{y}) AND TRUTH (y)

$$\text{LOSS} = \mathcal{L}(\hat{y}, y)$$

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

COST = LOSS FOR THE ENTIRE DATASET

