

Probabilistic Discriminative Models: Fixed Basis Functions

Karan Nathwani

Topics in Linear Classification using Probabilistic Discriminative Models

- Distinction between 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

Probabilistic Generative Models

- **Two-class classification:** Posterior of class C_1 : $p(C_1|\mathbf{x})$ can be written as a logistic sigmoid operating on linear function of \mathbf{x} , i.e., $\sigma(\mathbf{w}^T\mathbf{x} + w_0)$, for wide choice of forms for $p(\mathbf{x}|C_k)$
 - For Gaussians with same covariance matrix $\boxed{\mathbf{w} = \Sigma^{-1}(\mu_1 - \mu_2)} \quad \boxed{w_0 = -\frac{1}{2}\mu_1^T \Sigma^{-1} \mu_1 + \frac{1}{2}\mu_2^T \Sigma^{-1} \mu_2 + \ln \frac{p(C_1)}{p(C_2)}}$
 - For Discrete binary features x_i
 - which are linear functions of features
- **Multiclass case:**
 - Posterior probability of class C_k i.e.,
 - $p(C_k|\mathbf{x})$ is given by a softmax transformation of a linear function of \mathbf{x}
- **MLE used to get parameters of $p(\mathbf{x}|C_k)$ and $p(C_k)$**
 - This indirect approach to find parameters of a generalized linear model is called generative
 - Since we can use such a model to generate synthetic data from marginal $p(\mathbf{x})$

$$a_k(\mathbf{x}) = \ln(p(\mathbf{x} | C_k)p(C_k))$$

$$= \sum_{i=1}^D \{x_i \ln \mu_{ki} + (1 - x_i) \ln(1 - \mu_{ki})\} + \ln p(C_k)$$

$$p(C_k | \mathbf{x}) = \frac{p(\mathbf{x} | C_k)p(C_k)}{\sum_j p(\mathbf{x} | C_j)p(C_j)}$$

$$= \frac{\exp(a_k)}{\sum_j \exp(a_j)}$$

Normalized
Exponential
Or Softmax

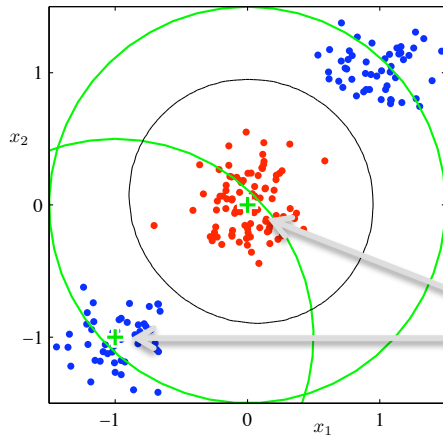
Probabilistic Discriminative Models

- An alternative approach
 - Use the functional form of the generalized linear model explicitly
 - Determine its parameters using maximum likelihood
 - There is an efficient algorithm to find such solutions known as Iterative reweighted least squares (IRLS)
- Advantages
 - Fewer adaptive parameters
 - Improved predictive performance when $p(\mathbf{x}|C_k)$ assumptions are poor approximations

Fixed Basis Functions

Although we use linear classification models
Linear-separability in feature space does not imply
linear-separability in input space

Original Input Space (x_1, x_2)

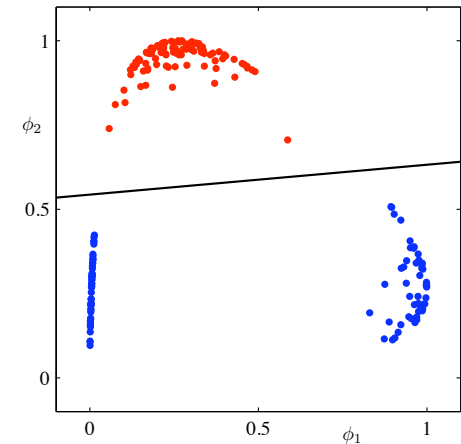


not linearly separable

Nonlinear transformation
of inputs using vector of
basis functions $\phi(x)$

Basis functions are
Gaussian with centers
Shown by crosses and
Green contours

Feature Space (ϕ_1, ϕ_2)



linearly separable

Basis functions with increased dimensionality is often used

Limitation of Fixed Basis Functions

- Nonlinear transformations cannot remove overlap between classes
 - They can even increase the overlap!
 - Still fixed nonlinear basis functions play an important role
- One solution: basis functions that adapt to data
 - SVMs use basis functions centered on the data points and select a fixed subset of them

Logistic Regression

Karan Nathwani

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 $\Phi = [\phi_1, \dots, \phi_M]^T$

$$p(C_1|\Phi) = y(\Phi) = \sigma(\mathbf{w}^T \Phi)$$

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

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

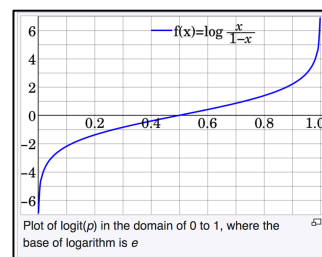
– Known as logistic regression²⁰ in statistics

- Although a model for classification rather than for regression

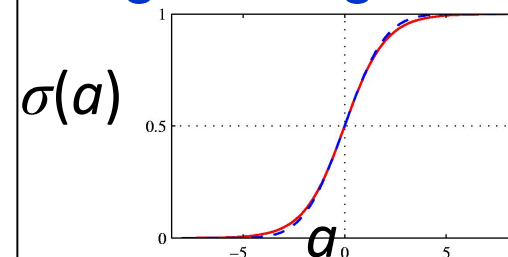
- Logit function:

– 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

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

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

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

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

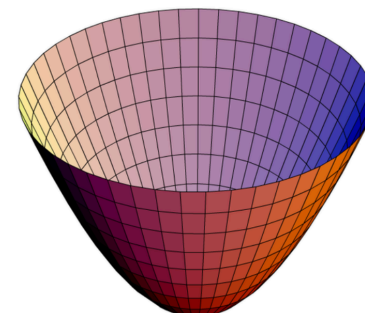
$$\nabla E(\mathbf{w}) = 0$$

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



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

Gradient of the error function

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

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

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

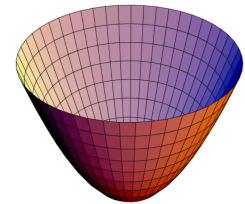
Simple Sequential Algorithm

- 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

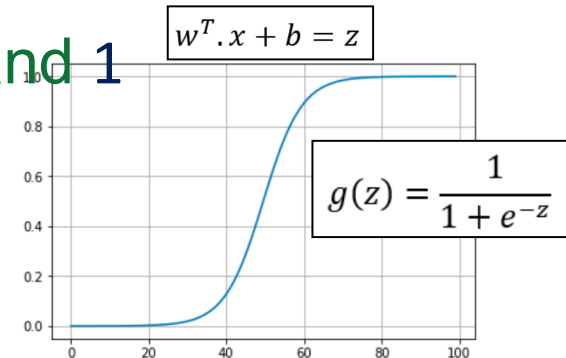
Samples are presented one at a time in
which each each of the weight vectors is updated

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

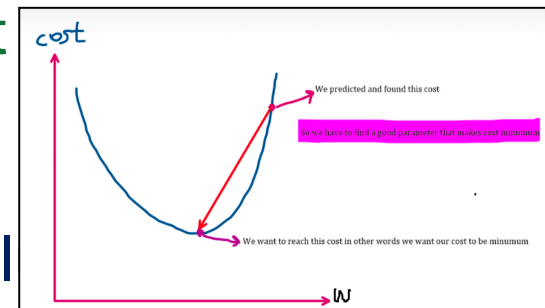
Loss function is the loss for a training example z

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 .

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

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

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

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

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

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

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

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

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

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

db

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

$$\frac{\partial L(p, y)}{\partial w} = dw = \frac{1}{m} \sum_{k=1}^m x dz^k$$

dw

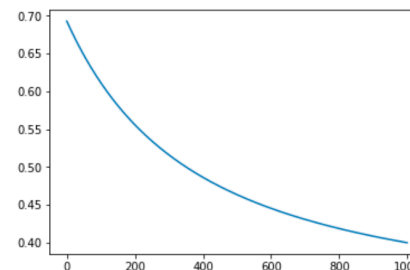
Logistic Regression Code in Python

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):

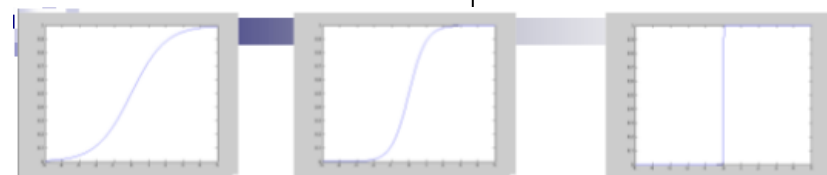
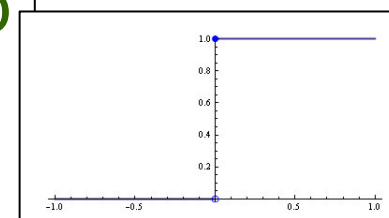
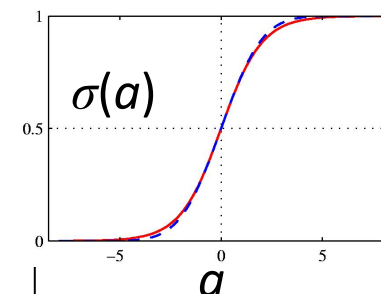
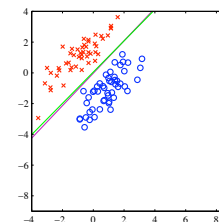
$$\begin{aligned}z &= w^T \cdot x + b \\p &= s(z) \\dz &= p - y \\dw &= \frac{1}{m} X dz^T \\db &= \frac{1}{m} \sum_{k=1}^m dz^k \\b &:= b - \alpha db \\w &:= w - \alpha dw\end{aligned}$$



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

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 $a = \mathbf{w}^T \boldsymbol{\phi} = 0$
 - Logistic sigmoid becomes infinitely steep
 - A Heavyside step function $||\mathbf{w}||$ goes to ∞
 - Solution
 - Penalizing wts
 - Recall in linear regression



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

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

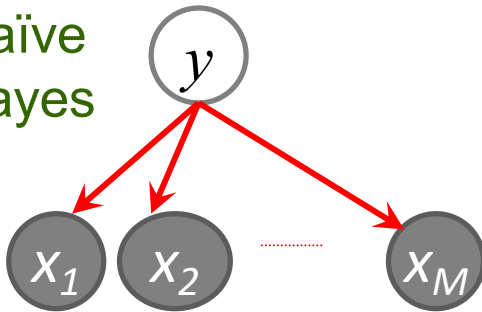
$\frac{1}{1 + e^{-x}}$	$\frac{1}{1 + e^{-2x}}$	$\frac{1}{1 + e^{-100x}}$
------------------------	-------------------------	---------------------------

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

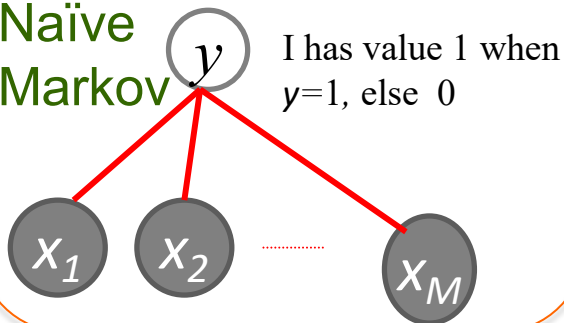
2. Discriminative: estimate joint parameters w_i

Potential Functions (log-linear)

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

$$\phi_0(y) = \exp\{w_0 \mathbb{I}\{y=1\}\}$$

Naïve
Markov



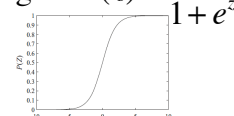
For classification:
Unnormalized

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

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:

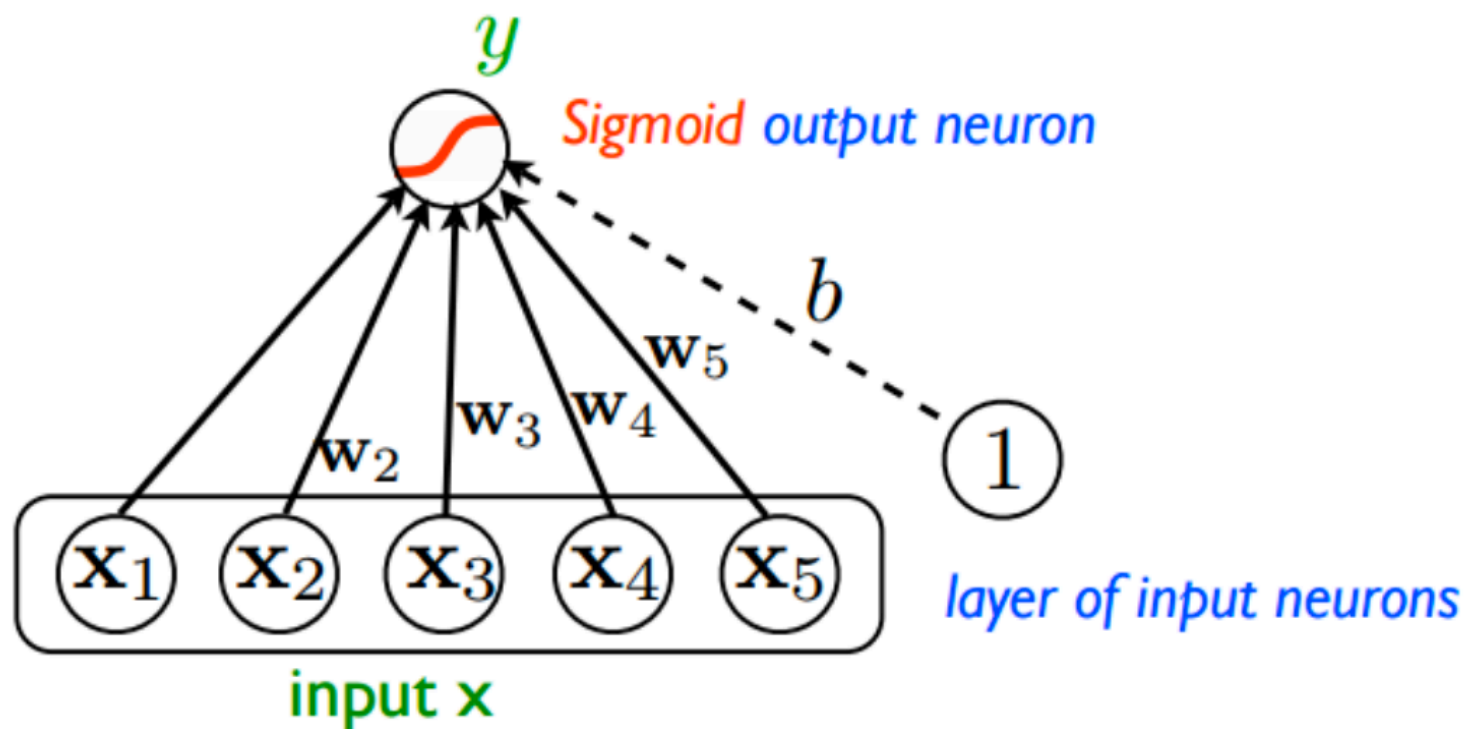
Edges, image patches sharing same pixels

multiclass

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

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

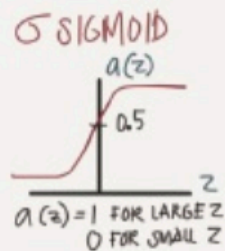
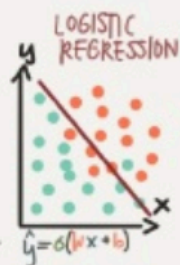
Logistic Regression is a special architecture of a neural network



BINARY CLASSIFICATION

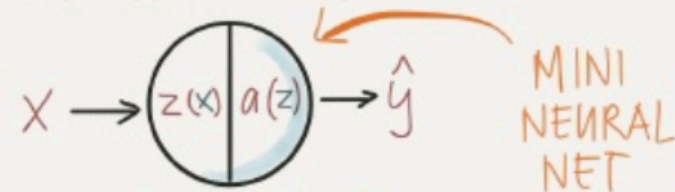


1: CAT
0: NOT CAT
 y

LOGISTIC REGRESSION
AS A NEURAL NETFINDING THE MINIMUM
WITH GRADIENT DESCENT

1. FIND THE DOWNHILL DIRECTION (USING DERIVATIVES)
 2. WALK (UPDATE $w \in b$) 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 $w \in b$
- REPEAT UNTIL IT CONVERGES

THE TASK IS TO LEARN $w \in b$ 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

