# Logistic regression Computer Vision (CSCI 5520G)

Faisal Z. Qureshi

http://vclab.science.ontariotechu.ca



- ► Logistic regression is for binary classification
- $\blacktriangleright$  The target variable y takes on values in  $\{0,1\}$

- Logistic regression is for binary classification
- ▶ The target variable y takes on values in  $\{0,1\}$

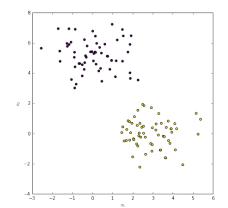
► Data:

$$\mathbf{X} = \left\{ \left( \underbrace{\mathbf{x}^{(i)}}_{\text{sample}}, \underbrace{y^{(i)}}_{\text{label}} \right) \middle| i \in [1, N], \mathbf{x}^{(i)} \in \mathbb{R}^{M}, y^{(i)} \in [0, 1] \right\}$$

- Logistic regression is for binary classification
- ▶ The target variable y takes on values in  $\{0,1\}$

► Data:

$$\mathbf{X} = \left\{ \left( \underbrace{\mathbf{x}^{(i)}}_{\text{sample label}}, \underbrace{y^{(i)}}_{\text{label}} \right) \middle| i \in [1, N], \mathbf{x}^{(i)} \in \mathbb{R}^{M}, y^{(i)} \in [0, 1] \right\}$$





#### Binary classification

The goal of binary classification is to learn  $h_{\theta}(\mathbf{x})$ , which can be used to assign a label  $y \in \{0,1\}$  to the input  $\mathbf{x}$ . Label y takes values in  $\{0,1\}$ , so we can use Bernoulli distribution to specify its probability distribution. Specifically

$$Pr(y = 1) = h_{\theta}(\mathbf{x})$$
$$Pr(y = 0) = 1 - h_{\theta}(\mathbf{x})$$

### Binary classification

The goal of binary classification is to learn  $h_{\theta}(\mathbf{x})$ , which can be used to assign a label  $y \in \{0,1\}$  to the input  $\mathbf{x}$ . Label y takes values in  $\{0,1\}$ , so we can use Bernoulli distribution to specify its probability distribution. Specifically

$$Pr(y = 1) = h_{\theta}(\mathbf{x})$$
$$Pr(y = 0) = 1 - h_{\theta}(\mathbf{x})$$

Or more succinctly

$$Pr(y) = h_{\theta}(\mathbf{x})^{y} \left(1 - h_{\theta}(\mathbf{x})\right)^{1-y}$$

### Binary classification

The goal of binary classification is to learn  $h_{\theta}(\mathbf{x})$ , which can be used to assign a label  $y \in \{0,1\}$  to the input  $\mathbf{x}$ . Label y takes values in  $\{0,1\}$ , so we can use Bernoulli distribution to specify its probability distribution. Specifically

$$Pr(y = 1) = h_{\theta}(\mathbf{x})$$
$$Pr(y = 0) = 1 - h_{\theta}(\mathbf{x})$$

Or more succinctly

$$\Pr(y) = \underbrace{h_{\theta}(\mathbf{x})^y}_{\text{active when } y=1} \underbrace{(1 - h_{\theta}(\mathbf{x}))^{1-y}}_{\text{active when } y=0}$$

7/4

#### Bernoulli distribution

A Bernoulli random variable X takes values in  $\{0,1\}$ 

$$Pr(X|\theta) = \begin{cases} \theta & \text{if } X = 1\\ 1 - \theta & \text{otherwise} \end{cases}$$
$$= \theta^X (1 - \theta)^{1 - X}$$

#### Bernoulli distribution

A Bernoulli random variable X takes values in  $\{0,1\}$ 

$$Pr(X|\theta) = \begin{cases} \theta & \text{if } X = 1\\ 1 - \theta & \text{otherwise} \end{cases}$$
$$= \theta^X (1 - \theta)^{1 - X}$$

#### Example usage

Bernoulli distribution  $\mathrm{Ber}(X|\theta)$  can be used to model coin tosses.

### Likelihood for binary classification

Under the assumption that data is independant and identically distributed (i.e., i.i.d.) the likelihood for the entire data is

$$\Pr(y|\mathbf{X},\theta) = \prod_{i=1}^{N} h_{\theta}(\mathbf{x}^{(i)})^{y^{(i)}} \left(1 - h_{\theta}(\mathbf{x}^{(i)})\right)^{1 - y^{(i)}}$$

### Likelihood for binary classification

Under the assumption that data is independent and identically distributed (i.e., i.i.d.) the likelihood for the entire data is

$$\Pr(y|\mathbf{X}, \theta) = \prod_{i=1}^{N} h_{\theta}(\mathbf{x}^{(i)})^{y^{(i)}} \left(1 - h_{\theta}(\mathbf{x}^{(i)})\right)^{1 - y^{(i)}}$$

What form should  $h_{\theta}(.)$  take?

# Aside: Mean (Expectation)

- ▶ The mean is the "average" or "center of mass" of data.
- **Sample mean** (finite data):

$$\bar{x} = \frac{1}{n} \sum_{i=1}^{n} x_i$$

▶ **Probabilistic definition** (random variable *X*):

$$\mu = \mathbb{E}[X] = \begin{cases} \sum_{x} x \, P(X = x), & \text{if } X \text{ is discrete} \\ \int_{-\infty}^{\infty} x \, p(x) \, dx, & \text{if } X \text{ is continuous} \end{cases}$$

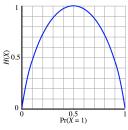
▶ Interpretation: Weighted average of possible values, weighted by their probabilities.

### Aside: Entropy

- Average level of information in a random variable.
- ▶ Given a discrete random variable X, which takes values in the alphabet  $\mathcal{X}$  and is distributed according to  $p: \mathcal{X} \to [0,1]$ :

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x)$$

- ► Choice of base for log varies with applications
  - Base 2 gives the unit of bits or shannons
  - ► Base *e* gives units of nats
  - Base 10 gives units of dits, bans, or hartley

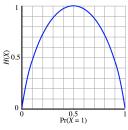


### Aside: Entropy

- Average level of information in a random variable.
- ▶ Given a discrete random variable X, which takes values in the alphabet  $\mathcal{X}$  and is distributed according to  $p: \mathcal{X} \to [0,1]$ :

$$H(X) = -\sum_{x \in \mathcal{X}} p(x) \log p(x) = \mathbb{E}_{x \sim p(x)} [-\log p(x)]$$

- Choice of base for log varies with applications
  - ► Base 2 gives the unit of bits or shannons
  - Base e gives units of nats
  - Base 10 gives units of dits, bans, or hartley

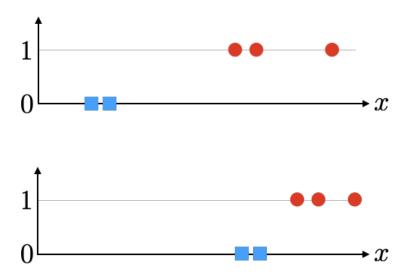


## Aside: Cross entropy

▶ Cross-entropy beween two distributions p and q is a measure of the average number of bits needed to identify an event from a set  $\mathcal{X}$  with true distribution p when the coding scheme used for the set is optimized for an estimated probability distribution q

$$H(p,q) = -\sum_{x \in \mathcal{X}} p(x) \log q(x) = -\mathbb{E}_{x \sim p(x)}[\log q(x)]$$

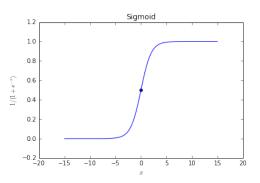
# Lets consider a simple 1D case for binary classification



# Sigmoid function

 $\operatorname{sigm}(x)$  refers to a  $\operatorname{\textit{sigmoid}}$  function, also known as the  $\operatorname{\textit{logistic}}$  or  $\operatorname{\textit{logit}}$  function.

$$sigm(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1}$$



For logistic regression, we set  $h_{\theta}(\mathbf{x}) = \operatorname{sigm}(\mathbf{x}^T \theta)$ . So

$$\Pr(y|\mathbf{X}, \theta) = \prod_{i=1}^{N} \left[ \frac{1}{1 + e^{-\mathbf{x}^{(i)}^{T} \theta}} \right]^{y^{(i)}} \left[ 1 - \frac{1}{1 + e^{-\mathbf{x}^{(i)}^{T} \theta}} \right]^{1 - y^{(i)}}$$

where

$$\mathbf{x}^T \theta = \theta_0 + \sum_{j=1}^M \theta_j \mathbf{x}_j$$

.

## Sigmoid function

$$\Pr(y|x,\theta) = \left[\frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}\right]^y \left[1 - \frac{1}{1 + e^{-(\theta_0 + \theta_1 x)}}\right]^{1 - y}$$

- lacktriangledown  $heta=( heta_0, heta_1)$  are model parameters.
- $\triangleright$   $\theta_0$  controls the shift.
- $\theta_1$  controls the scale (how steep is the slope of the sigmoid function).





#### Likelihood

$$L(\theta) = \Pr(y|\mathbf{X}, \theta)$$

#### Negative log-likelihood

$$l(\theta) = -\log L(\theta)$$

$$= -\sum_{i=1}^{N} y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))$$

We prefer to work in the log domain for mathematical convenience. Plus there are numerical advantages of working in the log domain.

#### Likelihood

$$L(\theta) = \Pr(y|\mathbf{X}, \theta)$$

#### Negative log-likelihood

$$\begin{split} l(\theta) &= -\log L(\theta) \\ &= -\sum_{i=1}^N y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)})) \text{Cross entropy for sa} \\ &= \sum_{i=1}^N \underbrace{-\left(y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) + (1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)}))\right)}_{\text{Cross entropy for sample } i} \end{split}$$

For binary classifiers, we need to minimize the negative log-likelihood, or minimize the cross-entropy between ground truth and predicted distributions.

#### Goal

Our goal is to find parameters  $\theta$  that maximize the likelihood (or minimize the negative log-likelihood).

$$\theta^* = \operatorname*{arg\,min}_{\theta} l(\theta)$$

# Derivative of sigmoid

$$\begin{split} \frac{d}{dx} \mathrm{sigm}(x) &= \frac{d}{dx} \frac{1}{1 + e^{-x}} \\ &= \frac{-(-1)e^{-x}}{(1 + e^{-x})^2} \\ &= \left(\frac{e^{-x}}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) \\ &= \left(\frac{1 - 1 + e^{-x}}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) \\ &= \left(1 - \frac{1}{1 + e^{-x}}\right) \left(\frac{1}{1 + e^{-x}}\right) \\ &= (1 - \mathrm{sigm}(x)) \, \mathrm{sigm}(x) \end{split}$$

# Gradient of a sigmoid w.r.t. $\theta$

We know that

$$\frac{d}{dx}\operatorname{sigm}(x) = (1 - \operatorname{sigm}(x))\operatorname{sigm}(x)$$

It follows

$$\frac{d}{d\theta} \operatorname{sigm}(\mathbf{x}^T \theta) = \left(1 - \operatorname{sigm}(\mathbf{x}^T \theta)\right) \operatorname{sigm}(\mathbf{x}^T \theta) \mathbf{x}$$

$$l^{(i)}(\theta) = -y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)}) - (1 - y^{(i)}) \log(1 - h_{\theta}(\mathbf{x}^{(i)}))$$

$$l^{(i)}(\theta) = -y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)})$$
$$-(1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)}))$$
$$= -y^{(i)} \log \operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta)$$
$$-(1-y^{(i)}) \log(1-\operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta))$$

$$l^{(i)}(\theta) = -y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)})$$
$$-(1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)}))$$
$$= -y^{(i)} \log \operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta)$$
$$-(1-y^{(i)}) \log(1-\operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta))$$

$$l^{(i)}(\theta) = -y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)})$$
$$-(1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)}))$$
$$= -y^{(i)} \log \operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta)$$
$$-(1-y^{(i)}) \log(1-\operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta))$$

Negative log likelihood contribution by sample i

$$l^{(i)}(\theta) = -y^{(i)} \log h_{\theta}(\mathbf{x}^{(i)})$$
$$-(1-y^{(i)}) \log(1-h_{\theta}(\mathbf{x}^{(i)}))$$
$$= -y^{(i)} \log \operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta)$$
$$-(1-y^{(i)}) \log(1-\operatorname{sigm}(\mathbf{x}^{(i)^{T}}\theta))$$

Gradient of  $l^{(i)}(\theta)$ :

$$\nabla_{\theta} l^{(i)} = ?$$

- ▶ Replacing  $\operatorname{sigm}(\mathbf{x}^{(i)^T})$  with s
- ightharpoonup Replacing  $y^{(i)}$  with y
- lacktriangle Replacing  $\mathbf{x}^{(i)}$  with  $\mathbf{x}$

$$\nabla_{\theta} l^{(i)} = \nabla_{\theta} \left[ -y \log s - (1 - y) \log(1 - s) \right]$$

- ▶ Replacing  $\operatorname{sigm}(\mathbf{x}^{(i)^T})$  with s
- ightharpoonup Replacing  $y^{(i)}$  with y
- ightharpoonup Replacing  $\mathbf{x}^{(i)}$  with  $\mathbf{x}$

$$\nabla_{\theta} l^{(i)} = \nabla_{\theta} \left[ -y \log s - (1 - y) \log(1 - s) \right]$$
$$= -y \frac{s(1 - s)\mathbf{x}}{s} - (1 - y) \frac{s(1 - s)\mathbf{x}}{1 - s}$$

- ▶ Replacing  $\operatorname{sigm}(\mathbf{x}^{(i)^T})$  with s
- ightharpoonup Replacing  $y^{(i)}$  with y
- ightharpoonup Replacing  $\mathbf{x}^{(i)}$  with  $\mathbf{x}$

$$\nabla_{\theta} l^{(i)} = \nabla_{\theta} \left[ -y \log s - (1 - y) \log(1 - s) \right]$$
$$= -y \frac{s(1 - s)\mathbf{x}}{s} - (1 - y) \frac{s(1 - s)\mathbf{x}}{1 - s}$$
$$= -y \mathbf{x} + y s \mathbf{x} - s \mathbf{x} - y s \mathbf{x}$$

- ▶ Replacing  $\operatorname{sigm}(\mathbf{x}^{(i)^T})$  with s
- ightharpoonup Replacing  $y^{(i)}$  with y
- ightharpoonup Replacing  $\mathbf{x}^{(i)}$  with  $\mathbf{x}$

$$\nabla_{\theta} l^{(i)} = \nabla_{\theta} \left[ -y \log s - (1 - y) \log(1 - s) \right]$$

$$= -y \frac{s(1 - s)\mathbf{x}}{s} - (1 - y) \frac{s(1 - s)\mathbf{x}}{1 - s}$$

$$= -y \mathbf{x} + y s \mathbf{x} - s \mathbf{x} - y s \mathbf{x}$$

$$= -y \mathbf{x} - s \mathbf{x}$$

#### Notation change

- ▶ Replacing  $\operatorname{sigm}(\mathbf{x}^{(i)^T})$  with s
- ightharpoonup Replacing  $y^{(i)}$  with y
- ightharpoonup Replacing  $\mathbf{x}^{(i)}$  with  $\mathbf{x}$

$$\nabla_{\theta} l^{(i)} = \nabla_{\theta} \left[ -y \log s - (1 - y) \log(1 - s) \right]$$

$$= -y \frac{s(1 - s)\mathbf{x}}{s} - (1 - y) \frac{s(1 - s)\mathbf{x}}{1 - s}$$

$$= -y \mathbf{x} + y s \mathbf{x} - s \mathbf{x} - y s \mathbf{x}$$

$$= -y \mathbf{x} - s \mathbf{x}$$

$$= -\mathbf{x}(y - s)$$

Therefore (after fixing the notation),

$$\nabla_{\theta} l^{(i)} = -\mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$$

Gradient of  $l(\theta)$  for *i*th example

$$\nabla_{\theta} l^{(i)} = -\mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$$

Stochastic gradient descent rule

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} l^{(i)}$$

Gradient of  $l(\theta)$  for *i*th example

$$\nabla_{\theta} l^{(i)} = -\mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$$

Stochastic gradient descent rule

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} l^{(i)}$$
  
=  $\theta^{(k)} + \eta \mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$ 

Gradient of  $l(\theta)$  for *i*th example

$$\nabla_{\theta} l^{(i)} = -\mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$$

Stochastic gradient descent rule

$$\theta^{(k+1)} = \theta^{(k)} - \eta \nabla_{\theta} l^{(i)}$$

$$= \theta^{(k)} + \eta \mathbf{x}^{(i)} (y^{(i)} - h_{\theta}(\mathbf{x}^{(i)}))$$

$$= \theta^{(k)} + \eta \mathbf{x}^{(i)} (y^{(i)} - \operatorname{sigm}(\mathbf{x}^{(i)^{T}} \theta)),$$

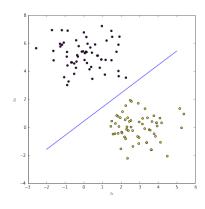
where  $\eta$  is the learning rate and k refers the the gradient descent iteration (step).

# Logistic regression for binary classification

Given a point  $\mathbf{x}^{(*)}$ , classify using the following rule

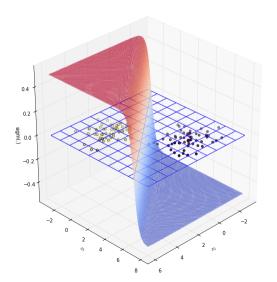
$$y^{(*)} = \begin{cases} 1 & \text{if } \Pr(y|\mathbf{x}^{(*)}, \theta) \ge 0.5\\ 0 & \text{otherwise} \end{cases}$$

The decision boundary is  $\mathbf{x}^T \theta = 0$ . Recall that this is where the sigmoid function is 0.5.



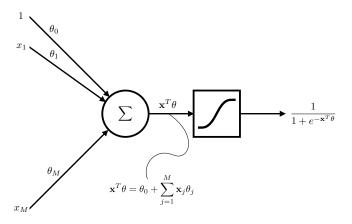
# Logistic regression for binary classification

- ► The decision boundary is  $\mathbf{x}^T \theta = 0$ 
  - $\blacktriangleright$  This is where sigm function is 0.5



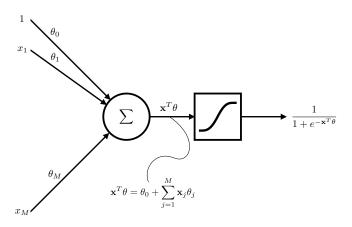
### Network view of logisitc regression

By changing the activation function to sigmoid and using the cross-entropy loss instead the least-squares loss that we use for linear regression, we are able to perform binary classification.



## Network view of logisitc regression

By changing the activation function to sigmoid and using the cross-entropy loss instead the least-squares loss that we use for linear regression, we are able to perform binary classification.



#### Artificial neuron

## Summary

- ▶ We looked at logisitc regression, a binary classifier.
- ► Bernoulli distribution

#### Summary

- We looked at logisite regression, a binary classifier.
- Bernoulli distribution
- Linear regression and logistic regression topics provide an excellent opportunity to study and understand the concepts underpinning neural networks

# Copyright and License

©Faisal Z. Qureshi



This work is licensed under a Creative Commons Attribution-NonCommercial 4.0 International License.