

# Machine Learning

## Discriminative Classification & Logistic Regression

Dariusz Hosseini

dariusz.hosseini@ucl.ac.uk  
Department of Computer Science  
University College London

# Lecture Overview

- 1** Lecture Overview
- 2 Classification
- 3 Logistic Regression
- 4 Summary

# Lecture Overview

By the end of this lecture you should:

- 1 Know the problem of **Classification** and understand the different paradigmatic approaches to its solution
- 2 Understand **logistic regression**, its motivation, and its use to solve the **binary classification** problem
- 3 Know that traditional logistic regression can be extended in a number of ways - both in terms of setting (the **Bayesian approach**) and in terms of application (**multinomial classification**)

# Lecture Overview

- 1 Lecture Overview
- 2 Classification**
- 3 Logistic Regression
- 4 Summary

# Setting

- Recall that in classification problems we seek to **learn** a mapping between input features and a discrete output label
- We can then use this mapping to make output **predictions** given novel input data
- In **binary** classification the output set comprises 2 classes, while in **multinomial** classification the output comprises greater than 2 unordered (categorical) labels

# Notation

## ■ Inputs

$$\mathbf{x} = [1, x_1, \dots, x_m]^T \in \mathbb{R}^{m+1}$$

## ■ Binary Outputs

$$y \in \{0, 1\}$$

## ■ Training Data

$$\mathcal{S} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^n$$

# Probabilistic Environment

- We assume:

- $\mathbf{x}$  is the outcome of a random variable  $\mathcal{X}$
- $y$  is the outcome of a random variable  $\mathcal{Y}$
- $(\mathbf{x}, y)$  are drawn i.i.d. from some data generating distribution,  $\mathcal{D}$ , i.e.:

$$(\mathbf{x}, y) \sim \mathcal{D}$$

and:

$$\mathcal{S} \sim \mathcal{D}^n$$

# Learning Problem

## ■ Representation

$$f \in \mathcal{F}$$

## ■ Evaluation

### ■ Loss Measure:

$$\mathcal{E}(f(\mathbf{x}), y) = \mathbb{I}[y \neq f(\mathbf{x})]$$

### ■ Generalisation Loss:

$$L(\mathcal{E}, \mathcal{D}, f) = \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathcal{X})]]$$

Where  $\mathcal{D}$  is characterised by  $p_{\mathcal{X}, y}(\mathbf{x}, y) = p_y(y|\mathbf{x})p_{\mathcal{X}}(\mathbf{x})$  for some pmf,  $p_y(\cdot|\cdot)$ , and some pdf,  $p_{\mathcal{X}}(\cdot)$

## ■ Optimisation

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathcal{X})]]$$



# Optimisation Problem

$$f^* = \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathcal{X})]]$$

## Optimisation Problem

$$\begin{aligned} f^* &= \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathbf{x})]] \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \mathbb{I}[y \neq f(\mathbf{x})] d\mathbf{x} \end{aligned}$$

## Optimisation Problem

$$\begin{aligned} f^* &= \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathbf{x})]] \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \mathbb{I}[y \neq f(\mathbf{x})] d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \left( f(\mathbf{x})(1-y) + (1-f(\mathbf{x}))y \right) d\mathbf{x} \end{aligned}$$

## Optimisation Problem

$$\begin{aligned} f^* &= \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathbf{x})]] \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \mathbb{I}[y \neq f(\mathbf{x})] d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \left( f(\mathbf{x})(1-y) + (1-f(\mathbf{x}))y \right) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int \left( p_y(y=1|\mathbf{x})(1-f(\mathbf{x})) + p_y(y=0|\mathbf{x})f(\mathbf{x}) \right) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

## Optimisation Problem

$$\begin{aligned} f^* &= \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathbf{x})]] \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \mathbb{I}[y \neq f(\mathbf{x})] d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) (f(\mathbf{x})(1-y) + (1-f(\mathbf{x}))y) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int (p_y(y=1|\mathbf{x})(1-f(\mathbf{x})) + p_y(y=0|\mathbf{x})f(\mathbf{x})) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int (p_y(y=1|\mathbf{x})(1-f(\mathbf{x})) + (1-p_y(y=1|\mathbf{x}))f(\mathbf{x})) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

## Optimisation Problem

$$\begin{aligned} f^* &= \operatorname{argmin}_{f \in \mathcal{F}} \mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathbf{x})]] \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) \mathbb{I}[y \neq f(\mathbf{x})] d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \sum_{y \in \{0,1\}} \int p_y(y|\mathbf{x}) p_{\mathbf{x}}(\mathbf{x}) (f(\mathbf{x})(1-y) + (1-f(\mathbf{x}))y) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int (p_y(y=1|\mathbf{x})(1-f(\mathbf{x})) + p_y(y=0|\mathbf{x})f(\mathbf{x})) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int (p_y(y=1|\mathbf{x})(1-f(\mathbf{x})) + (1-p_y(y=1|\mathbf{x}))f(\mathbf{x})) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \\ &= \operatorname{argmin}_{f \in \mathcal{F}} \int (p_y(y=1|\mathbf{x}) + f(\mathbf{x})(1-2p_y(y=1|\mathbf{x}))) p_{\mathbf{x}}(\mathbf{x}) d\mathbf{x} \end{aligned}$$

## Optimisation Problem

- We need to find a function,  $f^*$ , which is optimal for all  $\mathbf{x} \in \text{dom}(\mathcal{X})$ :

$$\begin{aligned} f^*(\mathbf{x}) &= \underset{f(\mathbf{x})}{\operatorname{argmin}} (p_y(y = 1|\mathbf{x}) + f(\mathbf{x})(1 - 2p_y(y = 1|\mathbf{x}))) \\ &= \underset{f(\mathbf{x})}{\operatorname{argmin}} M(\mathbf{x}) \end{aligned}$$

- This is a **discrete optimisation problem**:

- if:  $f(\mathbf{x}) = 1$  then:  $M(\mathbf{x}) = 1 - p_y(y = 1|\mathbf{x})$
- if:  $f(\mathbf{x}) = 0$  then:  $M(\mathbf{x}) = p_y(y = 1|\mathbf{x})$

- This means that if  $f^*(\mathbf{x}) = 1$  optimality implies:

$$\begin{aligned} 1 - p_y(y = 1|\mathbf{x}) &\leq p_y(y = 1|\mathbf{x}) \\ p_y(y = 1|\mathbf{x}) &\geq 0.5 \end{aligned}$$

# Bayes Optimal Classifier

- So the generalisation minimiser for the **Misclassification Loss** can be specified entirely in term of the **posterior distribution**:

$$f^*(\mathbf{x}) = \begin{cases} 1 & \text{if } p_y(y = 1|\mathbf{x}) \geq 0.5 \\ 0 & \text{if } p_y(y = 1|\mathbf{x}) < 0.5 \end{cases}$$

- It is known as the **Bayes Optimal Classifier**
- Of course, different loss functions will lead to different optimal classifiers...



## Alternative Classification Loss Functions

- Misclassification loss can be expressed in terms of a **loss matrix**:

$$\begin{array}{cc} y = 0 & y = 1 \\ \left[ \begin{array}{cc} 0 & 1 \\ 1 & 0 \end{array} \right] & \begin{array}{l} f(\mathbf{x}) = 0 \\ f(\mathbf{x}) = 1 \end{array} \end{array}$$

- Different loss matrices will lead to different classifiers...
- ...In particular they will lead to classification thresholds on the posterior distribution which differ from 0.5
- For example, cancer diagnosis might exhibit the following matrix:

$$\begin{array}{cc} y = 0 & y = 1 \\ \left[ \begin{array}{cc} 0 & 1,000 \\ 1 & 0 \end{array} \right] & \begin{array}{l} f(\mathbf{x}) = 0 \\ f(\mathbf{x}) = 1 \end{array} \end{array}$$

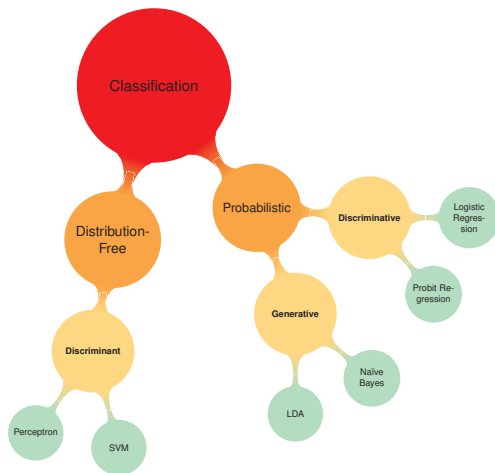
## Probabilistic Classifier

- Characterisation of  $f^*$  in terms of  $p_y(y = 1|\mathbf{x})$  suggests a classification paradigm:
- **Probabilistic Classification**
  - The classification problem reduces to an inference problem in which we must learn the posterior output class probability,  $p_y(y = 1|\mathbf{x})$
  - Here  $p_y(y = 1|\mathbf{x})$  characterises an **inhomogeneous Bernoulli distribution**
  - We must learn a Bernoulli distribution for each  $\mathbf{x} \in \text{dom}(\mathcal{X})$

# Probabilistic Classifier

- We contrast this with an alternative paradigm:
- **Distribution-Free Classification**
  - Here we seek to learn the classification boundary (equivalently  $f^*$ ) directly, without resorting to probabilistic inference
  - An example is the PAC approach where we seek to approximate  $\mathbb{E}_{\mathcal{D}} [\mathbb{I}[y \neq f(\mathcal{X})]]$  without reference to any explicit pdf and then to optimise this new quantity in order to learn  $f^*$

# Classification Approaches



## Generative Classification

- We note **Bayes' Theorem**:

$$\begin{aligned} p_y(y|\mathbf{x}) &= \frac{p_{\mathbf{x}}(\mathbf{x}|y)p_y(y)}{p_{\mathbf{x}}(\mathbf{x})} \\ &= \frac{p_{\mathbf{x}}(\mathbf{x}|y)p_y(y)}{\sum_{y \in \{0,1\}} p_{\mathbf{x}}(\mathbf{x}|y)p_y(y)} \end{aligned}$$

- Learn  $p_y(y|\mathbf{x})$  indirectly by inferring  $p_{\mathbf{x}}(\mathbf{x}|y)$  and  $p_y(y)$  for each class separately
- Most demanding approach in terms of number of parameters to learn
- Allows us to learn  $p_{\mathbf{x}}(\mathbf{x})$  which can be useful in **novelty detection**

# Discriminative Classification

- Attempts to learn  $p_y(y|\mathbf{x})$  directly
- Less demanding in terms of number of parameters to learn

## Discriminant Classification

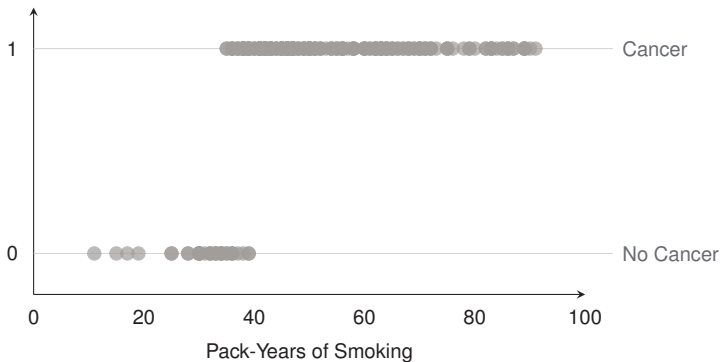
- As mentioned, this does not seek to infer the pdf at all.
  - Instead it seeks to learn  $f^*$  directly
- Least demanding in terms of number of parameters to learn
- Inflexible - requires us to run algorithm afresh for changes of the loss function
  - c.f. probabilistic approaches, which can be used to update  $f^*$  trivially
- Does not deal well with **class imbalance**

# Lecture Overview

- 1 Lecture Overview
- 2 Classification
- 3 Logistic Regression**
- 4 Summary



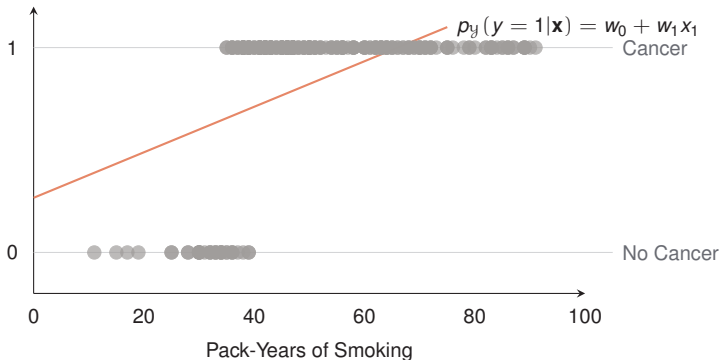
# Throat Cancer Prediction



## Discriminative Classification

- Let us focus on learning the inhomogeneous Bernoulli distribution  $p_y(y = 1|\mathbf{x})$  directly
- Can we model  $p_y(y = 1|\mathbf{x})$  as a linear function of  $\mathbf{x}$ , i.e.  $\mathbf{w} \cdot \mathbf{x}$
- In other words can we just use Linear Regression?

# Throat Cancer Prediction: Linear Regression

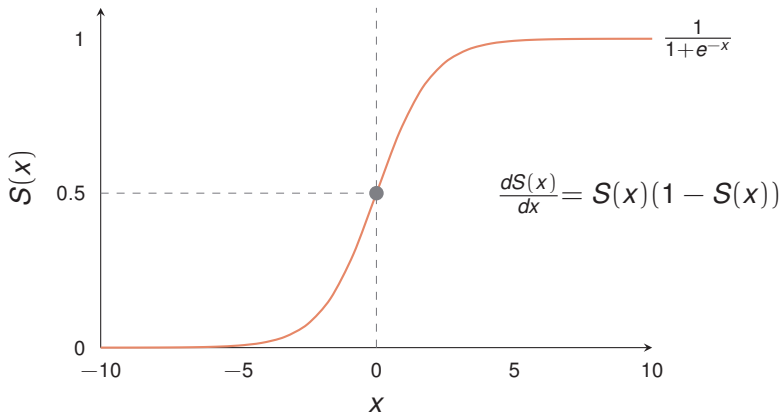


## Discriminative Classification: Linear Regression

- Can we just use Linear Regression?
- No!
  - $p_y(y|\mathbf{x})$  must lie in the range  $[0, 1]$ , while linear functions are unbounded
  - We need to learn a function that will **squash** the output of our model into  $[0, 1]$
  - One such function is the **logistic sigmoid**,  $S$ :

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

# Logistic Sigmoid



## Logistic Regression Model

- So we could attempt to model  $p_y(y = 1|\mathbf{x})$  as:

$$p_y(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}$$

- In other words our representation,  $\mathcal{F}$ , becomes:

$$\mathcal{F} = \left\{ f_{\mathbf{w}}(\mathbf{x}) = \mathbb{I}[p_y(y = 1|\mathbf{x}) \geq 0.5] \mid p_y(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}, \mathbf{w} \in \mathbb{R}^{m+1} \right\}$$

- This is the **Logistic Regression Model**

## Odds Ratio

- Re-arranging:

$$\begin{aligned}p_y(y = 1|\mathbf{x}) &= \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} \\e^{-\mathbf{w} \cdot \mathbf{x}} &= \frac{1}{p_y(y = 1|\mathbf{x})} - 1 \\ \mathbf{w} \cdot \mathbf{x} &= \ln \left( \frac{p_y(y = 1|\mathbf{x})}{1 - p_y(y = 1|\mathbf{x})} \right) \\ &= \text{logit}(p_y(y = 1|\mathbf{x}))\end{aligned}$$

- Here  $\frac{p_y(y=1|\mathbf{x})}{1-p_y(y=1|\mathbf{x})}$  is the **odds ratio**

## Linear Discriminant

- If  $f_{\mathbf{w}}(\mathbf{x}) = 1$  then  $p_y(y = 1|\mathbf{x}) \geq 0.5$ ,  $\frac{p_y(y=1|\mathbf{x})}{1-p_y(y=1|\mathbf{x})} \geq 1$ ,  $\mathbf{w} \cdot \mathbf{x} \geq 0$
- If  $f_{\mathbf{w}}(\mathbf{x}) = 0$  then  $p_y(y = 1|\mathbf{x}) < 0.5$ ,  $\frac{p_y(y=1|\mathbf{x})}{1-p_y(y=1|\mathbf{x})} < 1$ ,  $\mathbf{w} \cdot \mathbf{x} < 0$
- And  $\mathbf{w} \cdot \mathbf{x} = 0$  defines a **linear discriminant** separating hyperplane



# Evaluation

- How should we learn the parameters  $\mathbf{w}$ ?
- This is equivalent to asking how we should learn the distribution  $p_y(y|\mathbf{x})$
- From our earlier work on Probability and Point Estimation recall that we may adopt a number of different approaches
- The most common one is to assume a **frequentist** setting and use **maximum likelihood estimation**

## Evaluation

- Let us construct our log-likelihood function given Bernoulli outcomes:

$$\begin{aligned}\ln(L(\mathbf{w})) &= \ln \left( \prod_{i=1}^n p_y(y^{(i)} | \mathbf{x}^{(i)}) \right) \\ &= \sum_{i=1}^n \ln \left( p_y(y^{(i)} | \mathbf{x}^{(i)}) \right) \\ &= \sum_{i=1}^n y^{(i)} \ln \left( p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) \right) \\ &\quad + (1 - y^{(i)}) \ln \left( p_y(y^{(i)} = 0 | \mathbf{x}^{(i)}) \right)\end{aligned}$$

- This expression is known as the **Cross-Entropy** loss function

## Evaluation

- Previous two slides are true in general, we now specialise to the logistic regression model by substituting the following:

$$p_y(y = 1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} \quad p_y(y = 0|\mathbf{x}) = \frac{e^{-\mathbf{w} \cdot \mathbf{x}}}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}} \\ = \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}$$

- Into the log-likelihood function:

$$\ln(L(\mathbf{w})) = \sum_{i=1}^n \ln(p_y(y^{(i)} = 0|\mathbf{x}^{(i)})) + y^{(i)} \ln \frac{p_y(y^{(i)} = 1|\mathbf{x}^{(i)})}{p_y(y^{(i)} = 0|\mathbf{x}^{(i)})} \\ = \sum_{i=1}^n y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} - \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}})$$

# Optimisation

- Thus we seek  $\mathbf{w}_{\text{MLE}}$ , such that:

$$\mathbf{w}_{\text{MLE}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}) - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$$

# Motivation

- What leads us to adopt logistic regression as our model?
- After all, plenty of other 'squashing' functions exist, apart from the sigmoid
- For example the **cumulative distribution function** of the **standard normal distribution**

## Motivation 1: Additive Logistic Noise Latent Variable Model

- Let us assume that our data is generated by some hidden (**latent**) random variable,  $y^*$ , such that:

$$y^* = \mathbf{w} \cdot \mathbf{x} + \varepsilon$$

Where  $\varepsilon$  is an outcome associated with a random variable  $\epsilon$ , distributed according to a Logistic Distribution,  $\varepsilon \sim \text{Logistic}(0, 1)$ ; and where  $y^*$  are the outcomes associated with  $y^*$

## Logistic Distribution

- Let  $\epsilon$  be a continuous random variable, taking values  $\epsilon \in \mathbb{R}$ , with a **Logistic** distribution:

$$\epsilon \sim \text{Logistic}(\mu, s) \quad \text{where:} \quad \mu \in \mathbb{R}, s > 0$$

- This has a characteristic pdf,  $f_\epsilon$ :

$$f_\epsilon(\epsilon; \mu, s) = \frac{e^{-\frac{\epsilon - \mu}{s}}}{s \left(1 + e^{-\frac{\epsilon - \mu}{s}}\right)^2}$$

$$\mathbb{E}_{\mathcal{D}}[\epsilon] = \mu$$

$$\mathbb{P}(\epsilon < z) = \frac{1}{1 + e^{-\frac{z - \mu}{s}}}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

- We link a latent variable outcome,  $y^*$ , to an output,  $y$ , by assuming the following classification model:

$$y^{(i)} = \begin{cases} 1 & \text{if } y^{*(i)} \geq 0 \\ 0 & \text{if } y^{*(i)} < 0 \end{cases}$$



## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) = \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)})$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \end{aligned}$$

# Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \end{aligned}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \mathbb{P}(\epsilon^{(i)} < -\mathbf{w} \cdot \mathbf{x}^{(i)}) \end{aligned}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \mathbb{P}(\epsilon^{(i)} < -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \end{aligned}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \mathbb{P}(\epsilon^{(i)} < -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \\ &= \frac{e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \end{aligned}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \mathbb{P}(\epsilon^{(i)} < -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \\ &= \frac{e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \\ &= \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}^{(i)}}} \end{aligned}$$

## Motivation 1: Additive Logistic Noise Latent Variable Model

■ Then:

$$\begin{aligned} p_y(y^{(i)} = 1 | \mathbf{x}^{(i)}) &= \mathbb{P}(y^{*(i)} \geq 0 | \mathbf{x}^{(i)}) \\ &= \mathbb{P}(\mathbf{w} \cdot \mathbf{x}^{(i)} + \epsilon^{(i)} \geq 0) \\ &= \mathbb{P}(\epsilon^{(i)} \geq -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \mathbb{P}(\epsilon^{(i)} < -\mathbf{w} \cdot \mathbf{x}^{(i)}) \\ &= 1 - \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \\ &= \frac{e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}} \\ &= \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}^{(i)}}} \end{aligned}$$

■ ...Which is the logistic regression model



## Motivation 2

- Logistic regression leads to an optimisation problem which itself has directly attractive properties, it is:
  - **Smooth**
  - **Convex**
- We will see that we have turned a **non-convex problem** - trying to optimise  $\mathbb{E}_{\mathcal{D}} [\mathbb{I}[f(\mathcal{X}) \neq \mathcal{Y}]]$  - into a **convex problem**
- We do this by shifting our focus to probabilistic inference, and selecting a tractable form for  $p_{\mathcal{X}, \mathcal{Y}}$
- In practice it works well as a classifier (although it does not always get the probability prediction right)

## Recap

### ■ Representation:

$$\mathcal{F} = \left\{ f_{\mathbf{w}}(\mathbf{x}) = \mathbb{I}[p_y(y=1|\mathbf{x}) \geq 0.5] \mid p_y(y=1|\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{w} \cdot \mathbf{x}}}, \mathbf{w} \in \mathbb{R}^{m+1} \right\}$$

### ■ Evaluation:

$$\sum_{i=1}^n y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} - \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}) \quad \text{i.e. the log-likelihood}$$

### ■ Optimisation:

$$\mathbf{w}_{\text{MLE}} = \underset{\mathbf{w}}{\operatorname{argmin}} \sum_{i=1}^n \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}) - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)}$$

equivalent to a maximisation of the log-likelihood

# Optimisation

- We seek:

$$\mathbf{w}_{\text{MLE}} = \underset{\mathbf{w}}{\operatorname{argmin}} \left( \sum_{i=1}^n \ln(1 + e^{\mathbf{w} \cdot \mathbf{x}^{(i)}}) - y^{(i)} \mathbf{w} \cdot \mathbf{x}^{(i)} \right)$$

- But this has no analytic solution...
- ...However, we can apply a numerical technique to optimise, such as **gradient descent**
- But we should really check for **convexity**

# Convexity

- The sum of convex functions is also a convex function (think about this)
- So if we can prove that  $-y\mathbf{w} \cdot \mathbf{x}$  and  $\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}})$  are convex, then we have proved convexity of our objective
- Consider  $-y\mathbf{w} \cdot \mathbf{x}$ :

$$\nabla_{\mathbf{w}}(-y\mathbf{w} \cdot \mathbf{x}) = -y\mathbf{x}$$

- Taking second derivatives:

$$\nabla_{\mathbf{w}}^2(-y\mathbf{w} \cdot \mathbf{x}) = \mathbf{0}$$

$$\implies \mathcal{H} = \mathbf{0} \quad \text{which demonstrates convexity}$$

# Convexity

- Consider  $\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}})$ :

$$\nabla_{\mathbf{w}}(\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}})) = \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \mathbf{x}$$

- So:

$$\frac{\partial(\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i} = \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} x_i$$

# Convexity

■ Thus:

$$\frac{\partial^2 (\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} = \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial \left( \frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i$$

# Convexity

■ Thus:

$$\begin{aligned}\frac{\partial^2(\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\ &= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i\end{aligned}$$

# Convexity

■ Thus:

$$\begin{aligned}\frac{\partial^2(\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\&= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\&= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i\end{aligned}$$



# Convexity

■ Thus:

$$\begin{aligned}
 \frac{\partial^2 (\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} x_j e^{\mathbf{w} \cdot \mathbf{x}} x_i
 \end{aligned}$$

# Convexity

■ Thus:

$$\begin{aligned}
 \frac{\partial^2 (\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} x_j e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \left(1 - \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right) x_i x_j
 \end{aligned}$$

# Convexity

■ Thus:

$$\begin{aligned}
 \frac{\partial^2 (\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} x_j e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \left(1 - \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right) x_i x_j \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right) x_i x_j
 \end{aligned}$$

# Convexity

■ Thus:

$$\begin{aligned}
 \frac{\partial^2 (\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}}))}{\partial w_i \partial w_j} &= \frac{\partial(e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} + \frac{\partial\left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right)}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \frac{\partial(1 + e^{\mathbf{w} \cdot \mathbf{x}})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} \frac{\partial(\mathbf{w} \cdot \mathbf{x})}{\partial w_j} e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= e^{\mathbf{w} \cdot \mathbf{x}} x_j \frac{x_i}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} - \frac{1}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} e^{\mathbf{w} \cdot \mathbf{x}} x_j e^{\mathbf{w} \cdot \mathbf{x}} x_i \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \left(1 - \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right) x_i x_j \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{1 + e^{\mathbf{w} \cdot \mathbf{x}}} \left(\frac{1}{1 + e^{\mathbf{w} \cdot \mathbf{x}}}\right) x_i x_j \\
 &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} x_i x_j
 \end{aligned}$$

# Convexity

- Now, note that:

$$\mathbf{xx}^T = \begin{bmatrix} x_0x_0 & x_0x_1 & \cdot & \cdot & x_0x_m \\ x_1x_0 & x_1x_1 & \cdot & \cdot & x_1x_m \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ \cdot & \cdot & \cdot & \cdot & \cdot \\ x_mx_0 & x_mx_1 & \cdot & \cdot & x_mx_m \end{bmatrix}$$

$$\implies [\mathbf{xx}^T]_{ij} = x_i x_j$$

- This allows us to express the Hessian for  $\ln(1 + e^{\mathbf{w} \cdot \mathbf{x}})$  as:

$$\mathcal{H} = \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \mathbf{xx}^T$$

# Convexity

■ Thus:

$$\begin{aligned}\mathbf{a}^T \mathcal{H} \mathbf{a} &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \mathbf{a}^T \mathbf{x} \mathbf{x}^T \mathbf{a} \\ &= \frac{e^{\mathbf{w} \cdot \mathbf{x}}}{(1 + e^{\mathbf{w} \cdot \mathbf{x}})^2} \|\mathbf{a}^T \mathbf{x}\|_2^2 \geq 0 \\ \implies \mathcal{H} \succeq 0\end{aligned}$$

■ So our objective is convex...

■ ...And gradient descent will converge to a **global** (but not necessarily **unique**) optimal solution...if one exists

## Perfectly Separable Data

- Consider what happens if our training data is **linearly separable**
- And consider the objective function which we are seeking to maximise, for some  $\mathbf{w} = c\tilde{\mathbf{w}}$ , where  $c > 0$ , which separates the training data:

$$\sum_{i=1}^n \ln(1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}) - y^{(i)} c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}$$

## Perfectly Separable Data: Overfitting

- Now let's take the derivative of this with respect to  $c$ :

$$\sum_{i=1}^n \frac{e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}}{1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}} \tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} - y^{(i)} \tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} = \sum_{i=1}^n \tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} \left( \frac{e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}}{1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}} - y^{(i)} \right)$$

- But remember that for this  $\tilde{\mathbf{w}}$  the data is well classified, so:

- If  $y^{(i)} = 1$  then  $\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} > 0$  and  $\frac{e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}}{1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}} - y^{(i)} < 0$
- If  $y^{(i)} = 0$  then  $\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} < 0$  and  $\frac{e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}}{1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}} - y^{(i)} > 0$
- In other words  $\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} \left( \frac{e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}}{1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}} - y^{(i)} \right) < 0$  for all  $i$

- So if the data are linearly separable then there exists a  $\tilde{\mathbf{w}}$  such that the objective has a negative derivative with respect to  $c$



## Perfectly Separable Data: Overfitting

- And in this situation gradient descent will cause  $c$  to grow without bound
- This drives the sigmoid function to the **Heaviside function** which is infinitely steep at its inflection point
- So this is an example of the non-existence of a (finite) solution to the Logistic Regression optimisation problem
- We can regard it as a case of **overfitting** - our model becomes overly confident about the data and uses very large weights to achieve a 'perfect fit'.

## Perfectly Separable Data: Regularisation

- We can circumvent this problem via **regularisation**. For example an  $\ell_2$  norm regulariser will lead to the following objective in this situation:

$$\sum_{i=1}^n \ln(1 + e^{c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)}}) - y^{(i)} c\tilde{\mathbf{w}} \cdot \mathbf{x}^{(i)} + \lambda c^2 \|\tilde{\mathbf{w}}\|_2^2$$

- And the derivative of this with respect to  $c$  includes a strictly positive component,  $2\lambda c \|\tilde{\mathbf{w}}\|_2^2$ , which allows us to discern a (unique) stationary point with respect to  $c$ .

# Regularisation

- Of course, this doesn't necessarily mean that other situations will not lead to problems in optimisation...
- ...But actually, it turns out that, apart from the non-regularised, separable case, the optimisation of the Logistic Regression objective does in fact lead to unique solutions [Silvapulle ('81)]
- So the addition of a ridge regulariser serves to ensure that gradient descent will lead to a globally unique solution for any input data

# Multinomial Logistic Regression

- We may consider **multinomial logistic regression** in an equivalent formulation
  - Assume we have  $k$  classes such that:  $y \in \{1, \dots, k\}$
- The Bayes Optimal Classifier becomes:

$$f^*(\mathbf{x}) = \operatorname{argmax}_{y \in \{1, \dots, k\}} p_y(y|\mathbf{x})$$

## Multinomial Logistic Regression

- The model is now defined using the **softmax function**, and we seek to learn an **inhomogeneous multinomial distribution**:

$$p_y(y = j|\mathbf{x}) = \frac{e^{\mathbf{w}_j \cdot \mathbf{x}}}{\sum_{j=1}^k e^{\mathbf{w}_j \cdot \mathbf{x}}}$$

Which is characterised by  $\{\mathbf{w}_j\}_{j=1}^k$

- The maximum likelihood solution will now make use of the **multinomial cross entropy**, and will again result in a convex optimisation problem

# Lecture Overview

- 1 Lecture Overview
- 2 Classification
- 3 Logistic Regression
- 4 Summary**

# Summary

- 1 The **Classification** task can be tackled using probabilistic and distribution-free approaches
- 2 Probabilistic classification can be split into **Discriminative** and **Generative** approaches
- 3 **Logistic regression** is a well-motivated approach to discriminative classification which leads to a smooth, convex, optimisation problem

In the next lecture we will move on to consider Generative Classification in more detail, and in particular the **Naïve Bayes** algorithm