



Supervised Learning: Regression vs Classification

The input vector $\mathbf{x} = [x_1, \dots, x_D]^T$ consists of categorical and continuous features

- **Categorical:** $x_i \in \{1, \dots, K\}$

Often mapped to real values using, e.g., one-hot encoding

- **Continuous:** $x_i \in \mathbb{R}$

The target/label/output y

- **Classification:** y is a category

$$y \in \{1, \dots, C\}$$

- **Regression:** y is a real value

$$y \in \mathbb{R}$$

Our focus today: **Classification**

Discriminative vs Generative Models

Discriminative setting: Model the **conditional distribution** of the output y given the input \mathbf{x} and the model parameters θ

$$p(y | \mathbf{x}, \theta)$$

- Our focus so far: **Discriminative Linear Model (with Gaussian noise)**

$$p(y | \mathbf{x}, \mathbf{w}) = \mathbf{w} \cdot \mathbf{x} + \mathcal{N}(0, \sigma^2)$$

- Other noise models possible, e.g., Laplace
- Non-linearities using basis expansion
- Regularisation to avoid overfitting: Ridge, Lasso
- Validation to choose hyper-parameters
- Optimisation algorithms for model fitting

Generative setting: Model the **full joint distribution** of input \mathbf{x} and output y given the model parameters θ

$$p(\mathbf{x}, y | \theta)$$

Our focus today: **Generative models for classification**

Predicting using Generative Classification Models

Given generative model $p(\mathbf{x}, y | \theta)$ representing the joint distribution of \mathbf{x} and y

For a new input \mathbf{x}_{new} , the conditional distribution for y is

$$p(y = c | \mathbf{x}_{\text{new}}, \theta) = \frac{p(y = c | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c, \theta)}{\sum_{c'=1}^C p(y = c' | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c', \theta)}$$

where $c \in \{1, \dots, C\}$ are the possible categories/classes for y .

- **Numerator** is the joint probability distribution $p(\mathbf{x}_{\text{new}}, y = c | \theta)$
- **Denominator** is the marginal distribution $p(\mathbf{x}_{\text{new}} | \theta)$

To predict the most likely class: $\hat{y} = \arg\max_c p(y = c | \mathbf{x}_{\text{new}}, \theta)$

Defining a Generative Classification Model

In order to fit a generative model, we express the joint distribution as

$$p(\mathbf{x}, y | \theta, \pi) = p(y | \pi) \cdot p(\mathbf{x} | y, \theta)$$

- $p(y | \pi)$: **marginal distribution of outputs y** parameterised by π

We use parameters π_c such that $\sum_c \pi_c = 1$ and model the distribution as

$$p(y = c | \pi) = \pi_c$$

- $p(\mathbf{x} | y, \theta)$: **class-conditional distributions of input \mathbf{x} given the class label y** , parameterised by θ

For each class $c = 1, \dots, C$, we have the model

$$p(\mathbf{x} | y = c, \theta_c)$$

Toy Example: Dataset

Predict voter preference using in US elections

Voted in 2016?	Annual Income	State	Candidate Choice
Y	50K	OK	Beyoncé
N	173K	CA	Beyoncé
Y	80K	NJ	Borat
Y	150K	WA	Beyoncé
N	25K	WV	Kanye West
Y	85K	IL	Beyoncé
:	:	:	:
Y	1050K	NY	Borat
N	35K	CA	Borat
N	100K	NY	?

Toy Example: Classification

Marginal distribution for candidates

$$\begin{aligned} p(y = \text{Beyoncé} \mid \pi) &= \pi_{\text{Beyoncé}} \\ p(y = \text{Borat} \mid \pi) &= \pi_{\text{Borat}} \\ p(y = \text{Kanye West} \mid \pi) &= \pi_{\text{Kanye West}} \end{aligned}$$

Given that a voter supports Borat

$$p(\mathbf{x} \mid y = \text{Borat}, \theta_{\text{Borat}})$$

models the (class-conditional) distribution over \mathbf{x} given $y = \text{Borat}$ and θ_{Borat}

Similarly for $p(\mathbf{x} \mid y = \text{Beyoncé}, \theta_{\text{Beyoncé}})$ and $p(\mathbf{x} \mid y = \text{Kanye West}, \theta_{\text{Kanye West}})$

We need to pick "model" for $p(\mathbf{x} \mid y = c, \theta_c)$

To define the model, we estimate the parameters π_c, θ_c for $c \in \{1, \dots, C\}$

6

Marginal Distribution of the Target Classes via Maximum Likelihood

Assumptions:

- Dataset $\mathcal{D} = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$, where each vector \mathbf{x}_i has dimension D
- Each data point is drawn independently from the same distribution $p(\mathbf{x}, y)$
- Let N_c be the number of data points with $y_i = c$, so that $\sum_{c=1}^C N_c = N$

The probability for a single data point (\mathbf{x}_i, y_i) is:

$$p(\mathbf{x}_i, y_i \mid \theta, \pi) = p(y_i \mid \pi) \cdot p(\mathbf{x}_i \mid y_i, \theta) = \prod_{c=1}^C \pi_c^{\mathbb{I}(y_i=c)} \cdot p(\mathbf{x}_i \mid y_i, \theta)$$

The likelihood and the log-likelihood of the data are:

$$p(\mathcal{D} \mid \theta, \pi) = \prod_{i=1}^N p(\mathbf{x}_i, y_i \mid \theta, \pi) = \prod_{i=1}^N \left(\prod_{c=1}^C \pi_c^{\mathbb{I}(y_i=c)} \cdot p(\mathbf{x}_i \mid y_i, \theta) \right)$$

$$\log p(\mathcal{D} \mid \theta, \pi) = \sum_{c=1}^C N_c \log \pi_c + \sum_{i=1}^N \log p(\mathbf{x}_i \mid y_i, \theta)$$

The log-likelihood is easily separated into sums involving different parameters!

7

Maximum Log-Likelihood Estimator for Class Probability Distribution

We have the log-likelihood for the NBC

$$\log p(\mathcal{D} \mid \theta, \pi) = \underbrace{\sum_{c=1}^C N_c \log \pi_c}_{\text{dependent on } \pi} + \underbrace{\sum_{i=1}^N \log p(\mathbf{x}_i \mid y_i, \theta)}_{\text{independent of } \pi}$$

Let us obtain estimates for π . We get the constrained optimisation problem:

$$\begin{aligned} &\text{maximise} \quad \sum_{c=1}^C N_c \log \pi_c \\ &\text{subject to:} \quad \sum_{c=1}^C \pi_c = 1 \end{aligned}$$

This problem can be solved using the method of **Lagrange multipliers**

8

Recall: Constrained Optimisation Problems using Lagrange Multipliers

Suppose $f(\mathbf{z})$ is some function that we want to maximise subject to $g(\mathbf{z}) = 0$.

Constrained Objective

$$\underset{\mathbf{z}}{\operatorname{argmax}} f(\mathbf{z}), \quad \text{subject to: } g(\mathbf{z}) = 0$$

Lagrangian (Dual) Form

$$\Lambda(\mathbf{z}, \lambda) = f(\mathbf{z}) - \lambda g(\mathbf{z})$$

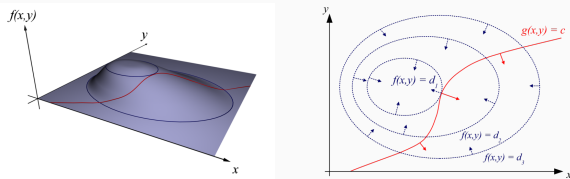
Any optimal solution to the constrained problem is a stationary point of $\Lambda(\mathbf{z}, \lambda)$

9

Recall: Constrained Optimisation Problems using Lagrange Multipliers

Any optimal solution to the constrained problem is a stationary point of

$$\Lambda(\mathbf{z}, \lambda) = f(\mathbf{z}) - \lambda g(\mathbf{z})$$



$$\nabla_{\mathbf{z}} \Lambda(\mathbf{z}, \lambda) = 0 \Rightarrow \nabla_{\mathbf{z}} f = \lambda \nabla_{\mathbf{z}} g$$

$$\frac{\partial \Lambda(\mathbf{z}, \lambda)}{\partial \lambda} = 0 \Rightarrow g(\mathbf{z}) = 0$$

10

Back to Our Constrained Optimisation Problem

Recall that we want to solve:

$$\begin{aligned} &\text{maximise:} \quad \sum_{c=1}^C N_c \log \pi_c \\ &\text{subject to:} \quad \sum_{c=1}^C \pi_c - 1 = 0 \end{aligned}$$

The Lagrangian formulation of our problem:

$$\Lambda(\pi, \lambda) = \sum_{c=1}^C N_c \log \pi_c - \lambda \left(\sum_{c=1}^C \pi_c - 1 \right)$$

We can write the partial derivatives and set them to 0:

$$\begin{aligned} \frac{\partial \Lambda(\pi, \lambda)}{\partial \pi_c} &= \frac{N_c}{\pi_c} - \lambda = 0 & \Rightarrow & \pi_c = \frac{N_c}{\lambda} \\ \frac{\partial \Lambda(\pi, \lambda)}{\partial \lambda} &= \sum_{c=1}^C \pi_c - 1 = 0 & \Rightarrow & \lambda = \sum_{c=1}^C N_c = N \end{aligned}$$

We get the estimates $\pi_c = \frac{N_c}{N}$

11

Summing Up: Marginal Distribution over the Classes

For Generative Classification Models,
the **marginal distribution** $p(y | \pi)$ over the classes is
the **empirical distribution over the classes**.

Remaining challenges:

- Choose suitable model for the **class-conditional distributions** $p(\mathbf{x} | y, \theta)$
- Estimate the parameters θ of the model

We next consider these challenges for several classification models

- **Naïve Bayes**
- **Gaussian (Quadratic/Linear) Discriminant Analysis**

12

Naïve Bayes Classifier (NBC)

Assume that **the features are conditionally independent given the class label c**

The *state, previous voting or annual income* are conditionally independent on being a Beyoncé/Borat/Kanye West supporter

Clearly, this assumption is "naïve" and never satisfied :(

But model fitting becomes very easy :)

Although the generative model is clearly inadequate, it actually works quite well

Goal is predicting class, not modelling the data!

13

Naïve Bayes Classifier (NBC)

Assume that **the features are conditionally independent given the class label c**

Class-conditional distribution for data point (\mathbf{x}_i, y_i) becomes:

$$p(\mathbf{x}_i | y_i = c, \theta) = p(\mathbf{x}_i | y_i = c, \theta_c) = \prod_{j=1}^D p(x_{ij} | \theta_{jc})$$

The joint probability distribution given the parameters θ and π becomes:

$$p(\mathbf{x}_i, y_i | \theta, \pi) = p(y_i | \pi) \cdot p(\mathbf{x}_i | y_i, \theta) = \prod_{c=1}^C \pi_c^{\mathbb{I}(y_i=c)} \cdot \prod_{c=1}^C \prod_{j=1}^D p(x_{ij} | \theta_{jc})^{\mathbb{I}(y_i=c)}$$

Log-likelihood for the entire dataset \mathcal{D} in case of NBC becomes:

$$\log p(\mathcal{D} | \theta, \pi) = \sum_{c=1}^C N_c \log \pi_c + \sum_{c=1}^C \sum_{j=1}^D \sum_{i: y_i=c} \log p(x_{ij} | \theta_{jc})$$

MLE for parameters π_c obtained as before.

14

Naïve Bayes Classifier: Models for Individual Features

Easy to mix and match different models for different features

Real-Valued Features

- x_j is real-valued *e.g.*, annual income
- Gaussian model: $\theta_{jc} = (\mu_{jc}, \sigma_{jc}^2)$
MLE for θ_{jc} given by empirical mean and variance
- Can use other distributions, *e.g.*, Laplace

Categorical Features

- x_j is categorical with categories $1, \dots, K$
- Multinoulli model: $x_j = \ell$ with probability $\theta_{jc, \ell}$ such that $\sum_{\ell=1}^K \theta_{jc, \ell} = 1$
MLE obtained similarly to that for π_c
- The special case when $x_j \in \{0, 1\}$ needs a single parameter $\theta_{jc} \in [0, 1]$
- No need to convert categorical variables to real-valued vectors!

15

Naïve Bayes Classifier: Number of Parameters

Assumptions:

- All the D features are binary, *i.e.*, every $x_j \in \{0, 1\}$, $j \in [D]$
- There are C target classes

With the naïve conditional independence assumption

- Naïve Bayes has at most $C \cdot D$ parameters $(\theta_{jc})_{j \in [D], c \in [C]}$

Without the conditional independence assumption

- We have to assign a probability for each of the 2^D possible feature vectors
- Thus, we have $C \cdot 2^D$ parameters!
- Either overfits if $N \approx C \cdot 2^D$ or computationally prohibitive if $N \gg C \cdot 2^D$

The 'naïve' assumption breaks the *curse of dimensionality* and avoids overfitting!

16

NBC: Prediction for Examples With Missing Data

Recall the prediction rule in a generative model:

$$p(y = c | \mathbf{x}_{\text{new}}, \theta) = \frac{p(y = c | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c, \theta)}{\sum_{c'=1}^C p(y = c' | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c', \theta)}$$

$$\stackrel{\text{NBC}}{=} \frac{\pi_c \cdot \prod_{j=1}^D p(x_j | y = c, \theta_{cj})}{\sum_{c'=1}^C \pi_{c'} \cdot \prod_{j=1}^D p(x_j | y = c', \theta_{c'j})}$$

Assume our data point is $\mathbf{x}_{\text{new}} = (?, x_2, \dots, x_D)$, *e.g.*, $(?, 100K, NY)$

Since x_1 is **missing**, we skip it:

$$p(y = c | \mathbf{x}_{\text{new}}, \theta) = \frac{\pi_c \cdot \prod_{j=2}^D p(x_j | y = c, \theta_{cj})}{\sum_{c'=1}^C \pi_{c'} \cdot \prod_{j=2}^D p(x_j | y = c', \theta_{c'j})}$$

This can be done for other generative models, but marginalisation requires summation/integration

17

NBC: Training With Missing Data

For Naïve Bayes Classifiers, training with missing entries is easy

Voted in 2016?	Annual Income	State	Candidate Choice
Y	50K	OK	Beyoncé
N	173K	CA	Beyoncé
Y	80K	NJ	Borat
Y	150K	WA	Beyoncé
N	25K	WV	Kanye West
Y	85K	IL	Beyoncé
⋮	⋮	⋮	⋮
Y	1050K	NY	Borat
N	35K	CA	Borat
?	100K	NY	?

Assume for Beyoncé voters, 103 had voted in 2016, 54 had not, and 25 didn't answer

We can set $\theta = \frac{103}{103+54}$ as the probability that a voter had voted in 2016, conditioned on being a Beyoncé supporter

18

Gaussian Discriminant Analysis

Recall the form of the joint distribution in a generative model

$$p(\mathbf{x}, y | \theta, \pi) = p(y | \pi) \cdot p(\mathbf{x} | y, \theta)$$

For classes, we use parameters π_c such that $\sum_c \pi_c = 1$

$$p(y = c | \pi) = \pi_c$$

We model the **class-conditional density** for class $c \in \{1, \dots, C\}$, as a **multivariate normal distribution** with mean μ_c and covariance matrix Σ_c

$$p(\mathbf{x} | y = c, \theta_c) = \mathcal{N}(\mathbf{x} | \mu_c, \Sigma_c)$$

Note: Name 'discriminant' unfortunate as this is a *generative* model.

It refers to the shape of the boundaries that discriminate between the classes.

More suitable than NBC in some cases:

Predict zebra or giraffe based on animal height and weight

19

MLE for Gaussian Discriminant Analysis

Given data $\mathcal{D} = ((\mathbf{x}_1, y_1), \dots, (\mathbf{x}_N, y_N))$, the log-likelihood is:

$$\log p(\mathcal{D} | \theta) = \sum_{c=1}^C N_c \log \pi_c + \sum_{c=1}^C \left(\sum_{i: y_i=c} \log \mathcal{N}(\mathbf{x}_i | \mu_c, \Sigma_c) \right)$$

As in the case of Naïve Bayes, we have $\pi_c = \frac{N_c}{N}$

For parameters μ_c and Σ_c , we have (Section 4.1 from Murphy)

$$\hat{\mu}_c = \frac{1}{N_c} \sum_{i: y_i=c} \mathbf{x}_i$$

$$\hat{\Sigma}_c = \frac{1}{N_c} \sum_{i: y_i=c} (\mathbf{x}_i - \hat{\mu}_c)(\mathbf{x}_i - \hat{\mu}_c)^\top$$

20

Quadratic Discriminant Analysis (QDA)

The prediction rule for this model is:

$$p(y = c | \mathbf{x}_{\text{new}}, \theta) = \frac{p(y = c | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c, \theta)}{\sum_{c'=1}^C p(y = c' | \theta) \cdot p(\mathbf{x}_{\text{new}} | y = c', \theta)}$$

When the densities $p(\mathbf{x} | y = c, \theta_c)$ are multivariate normal:

$$p(y = c | \mathbf{x}, \theta) = \frac{\pi_c |2\pi \Sigma_c|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^\top \Sigma_c^{-1}(\mathbf{x} - \mu_c)\right)}{\sum_{c'=1}^C \pi_{c'} |2\pi \Sigma_{c'}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c'})^\top \Sigma_{c'}^{-1}(\mathbf{x} - \mu_{c'})\right)}$$

The denominator is the same for all classes $c' \in \{1, \dots, C\}$.

21

Quadratic Discriminant Analysis: Decision Boundaries are Quadratic Curves

The boundary between classes c_1 and c_2 is given by the data points \mathbf{x} with

$$p(y = c_1 | \mathbf{x}, \theta) = p(y = c_2 | \mathbf{x}, \theta)$$

This means:

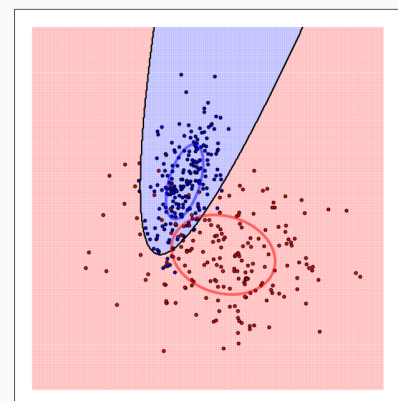
$$\frac{\pi_{c_1} |2\pi \Sigma_{c_1}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c_1})^\top \Sigma_{c_1}^{-1}(\mathbf{x} - \mu_{c_1})\right)}{\pi_{c_2} |2\pi \Sigma_{c_2}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c_2})^\top \Sigma_{c_2}^{-1}(\mathbf{x} - \mu_{c_2})\right)} = 1$$

By taking the log:

$$\frac{1}{2} \left((\mathbf{x} - \mu_{c_2})^\top \Sigma_{c_2}^{-1}(\mathbf{x} - \mu_{c_2}) - (\mathbf{x} - \mu_{c_1})^\top \Sigma_{c_1}^{-1}(\mathbf{x} - \mu_{c_1}) \right) = \log \left(\frac{\pi_{c_2} |2\pi \Sigma_{c_2}|^{-\frac{1}{2}}}{\pi_{c_1} |2\pi \Sigma_{c_1}|^{-\frac{1}{2}}} \right)$$

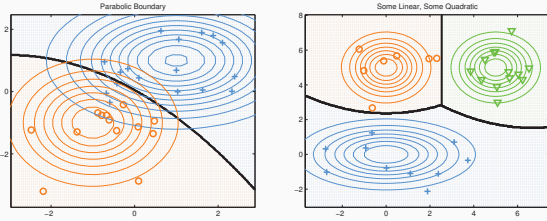
22

Quadratic Discriminant Analysis: Decision Boundary for Two Classes



23

Quadratic Discriminant Analysis: Quadratic and Linear Decision Boundaries



Not every point \mathbf{x} satisfying

$$p(y = c_1 | \mathbf{x}, \theta) = p(y = c_2 | \mathbf{x}, \theta)$$

lies on the decision boundary between the classes c_1 and c_2 !

Point \mathbf{x} belongs to class c_3 if $p(y = c_3 | \mathbf{x}, \theta) > p(y = c_i | \mathbf{x}, \theta), i \in \{1, 2\}$

Boundaries are given by piecewise quadratic curves

24

Linear Discriminant Analysis (LDA)

Special case: The covariance matrices are shared or tied across different classes

$$p(y = c | \mathbf{x}, \theta) = \frac{\pi_c |2\pi \Sigma_c|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma_c^{-1} (\mathbf{x} - \mu_c)\right)}{\sum_{c'=1}^C \pi_{c'} |2\pi \Sigma_{c'}|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma_{c'}^{-1} (\mathbf{x} - \mu_{c'})\right)}$$

becomes

$$\begin{aligned} p(y = c | \mathbf{x}, \theta) &= \frac{\pi_c |2\pi \Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma^{-1} (\mathbf{x} - \mu_c)\right)}{\sum_{c'=1}^C \pi_{c'} |2\pi \Sigma|^{-\frac{1}{2}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma^{-1} (\mathbf{x} - \mu_{c'})\right)} \\ &= \frac{\pi_c \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma^{-1} (\mathbf{x} - \mu_c)\right)}{\sum_{c'=1}^C \pi_{c'} \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_{c'})^T \Sigma^{-1} (\mathbf{x} - \mu_{c'})\right)} \end{aligned}$$

25

Linear Discriminant Analysis: Further Simplifications

Simplify the numerator:

$$\begin{aligned} &\pi_c \exp\left(-\frac{1}{2}(\mathbf{x} - \mu_c)^T \Sigma^{-1} (\mathbf{x} - \mu_c)\right) \\ &= \exp\left(\underbrace{\mu_c^T \Sigma^{-1} \mathbf{x}}_{\beta_c^T} - \underbrace{\frac{1}{2} \mu_c^T \Sigma^{-1} \mu_c}_{\gamma_c} + \log \pi_c\right) \cdot \exp\left(-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}\right) \\ &= \exp\left(\beta_c^T \mathbf{x} + \gamma_c\right) \cdot \exp\left(-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}\right) \end{aligned}$$

After expansion, the term $\exp\left(-\frac{1}{2} \mathbf{x}^T \Sigma^{-1} \mathbf{x}\right)$ also appears at denominator.

$$p(y = c | \mathbf{x}, \theta) = \frac{\exp\left(\beta_c^T \mathbf{x} + \gamma_c\right)}{\sum_{c'} \exp\left(\beta_{c'}^T \mathbf{x} + \gamma_{c'}\right)} = \text{softmax}(\boldsymbol{\eta})_c$$

where, $\boldsymbol{\eta} = [\beta_1^T \mathbf{x} + \gamma_1, \dots, \beta_C^T \mathbf{x} + \gamma_C]$.

26

Softmax Function

$$\text{softmax}([z_1, \dots, z_n])_i = \frac{e^{z_i}}{\sum_{j \in [n]} e^{z_j}}$$

Normalises the vector values into a probability distribution consisting of probabilities proportional to their exponentials.

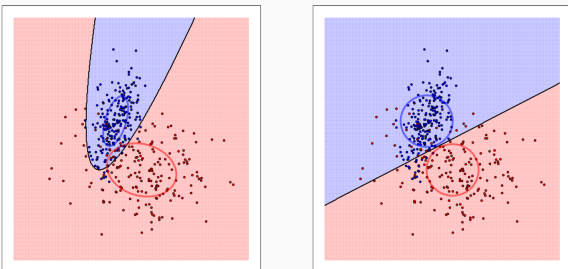
- Smooth approximation of **argmax** and not max!
- Used in statistical mechanics as the **Boltzmann distribution**
- Large inputs \rightarrow large probabilities
- Translation invariant: Multiplying all values by the same quantity does not change the ranking

$$\begin{aligned} \text{softmax}([1, 2, 3]) &\approx [0.090, 0.245, 0.665] \\ \text{softmax}([10, 20, 30]) &\approx [2 \times 10^{-9}, 4 \times 10^{-5}, 1] \end{aligned}$$

- Not scale invariant: De-emphasises the max value for values in $[0, 1]$
- Often used as the last activation function of a neural network to normalise the network output to a probability distribution over predicted output classes

27

QDA and LDA



28

Two-Class LDA

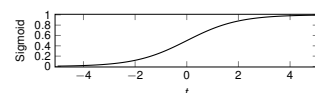
When we have only 2 classes, say 0 and 1:

$$\begin{aligned} p(y = 1 | \mathbf{x}, \theta) &= \frac{\exp\left(\beta_1^T \mathbf{x} + \gamma_1\right)}{\exp\left(\beta_1^T \mathbf{x} + \gamma_1\right) + \exp\left(\beta_0^T \mathbf{x} + \gamma_0\right)} \\ &= \frac{1}{1 + \exp\left(-((\beta_1 - \beta_0)^T \mathbf{x} + (\gamma_1 - \gamma_0))\right)} \\ &= \text{sigmoid}\left((\beta_1 - \beta_0)^T \mathbf{x} + (\gamma_1 - \gamma_0)\right) \end{aligned}$$

Sigmoid Function

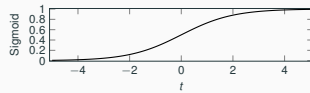
The sigmoid function is defined as:

$$\text{sigmoid}(t) = \frac{1}{1 + e^{-t}}$$



29

The Sigmoid Functions



- In general: Functions with S-shaped curves
- Examples: logistic (previous slide), hyperbolic tangent
- Logistic function: special case of softmax for 1D axis in 2D space
- Properties:
 - differentiable
 - non-negative derivative at each point
 - monotonic
 - convex for $t < 0$ and concave for $t > 0$

30

How to Prevent Overfitting

- The number of parameters in QDA is roughly $C \cdot D^2$
- In high-dimensions this can lead to overfitting
- Use diagonal covariance matrices – this corresponds to Naïve Bayes!
- Use weight tying/parameter sharing: LDA vs QDA
- Use a discriminative classifier (+ regularisation if needed)

31