



MACHINE LEARNING

LINEAR CLASSIFICATION

AGENDA

01 Introduction

Classification problem and linear regression

02 Binary classification

Logistic regression and Newton's method

03 Generalized Linear Models

Exponential Family, Building GLMs, Softmax

04 Discriminative and generative models

Differences, Gaussian Analysis, Bayes Classifier

05 Binary evaluation metrics

Sensitivity, Specificity, F1 Score, ROC Curve





AI

INTRODUCTION

INTRODUCTION

CLASSIFICATION PROBLEM



The problem is the same as that of **regression**: predicting a set of **output variables** y given a set of **input data** X .

The only **difference** is that the **values** of y take on a set of **discrete values**.

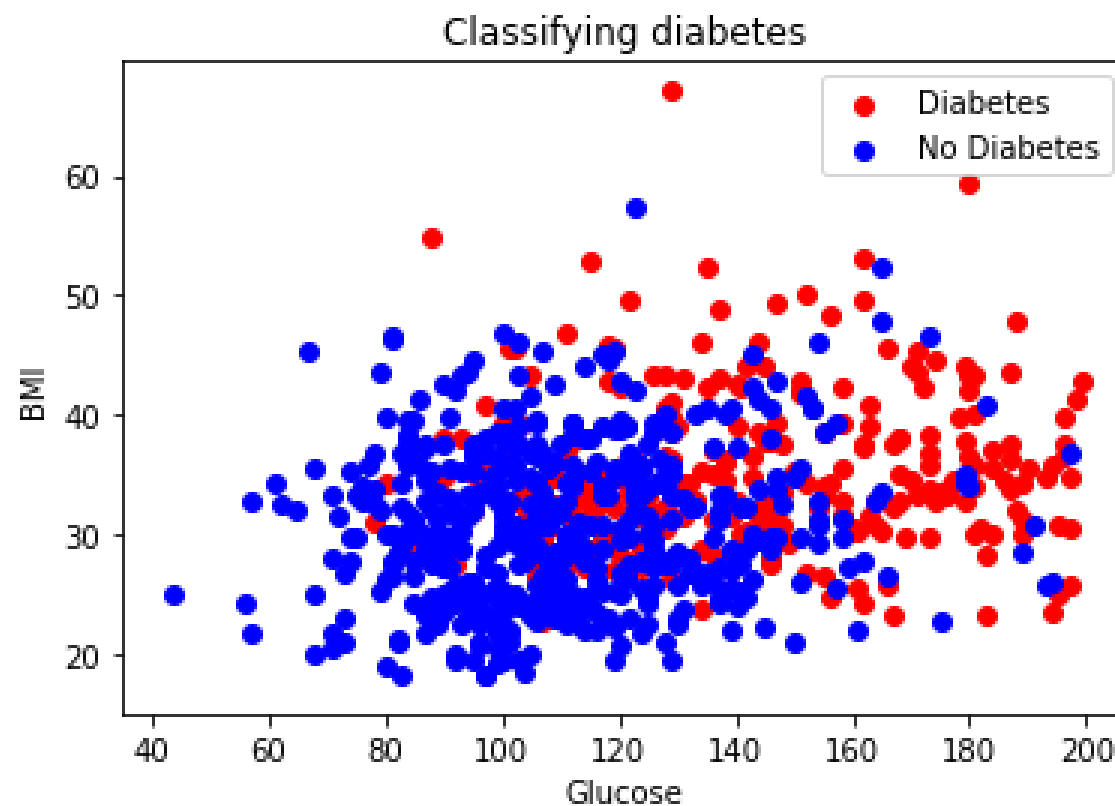
x_1 =Glucose [mg/dl]	x_2 =IMC	y = Presence of diabetes
148	33.6	1
85	26.6	0
183	23.3	1
89	28.1	0
137	43.1	1
⋮	⋮	⋮

INTRODUCTION

CLASSIFICATION PROBLEM



We plot 752 data points (patients) using Python:

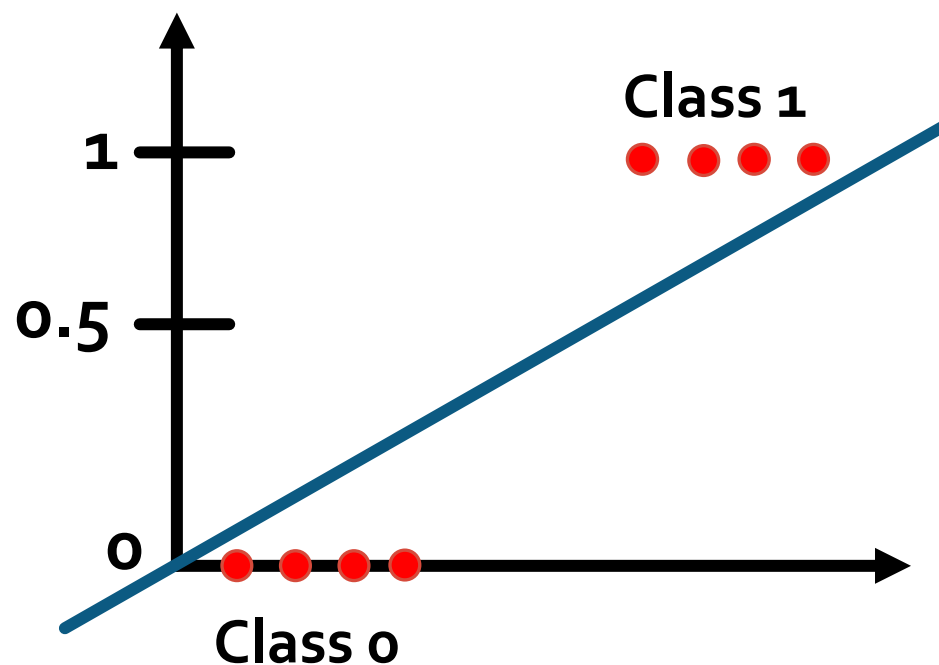


INTRODUCTION

LINEAR REGRESSION



WHY NOT APPLY LINEAR REGRESSION TO CLASSIFY?

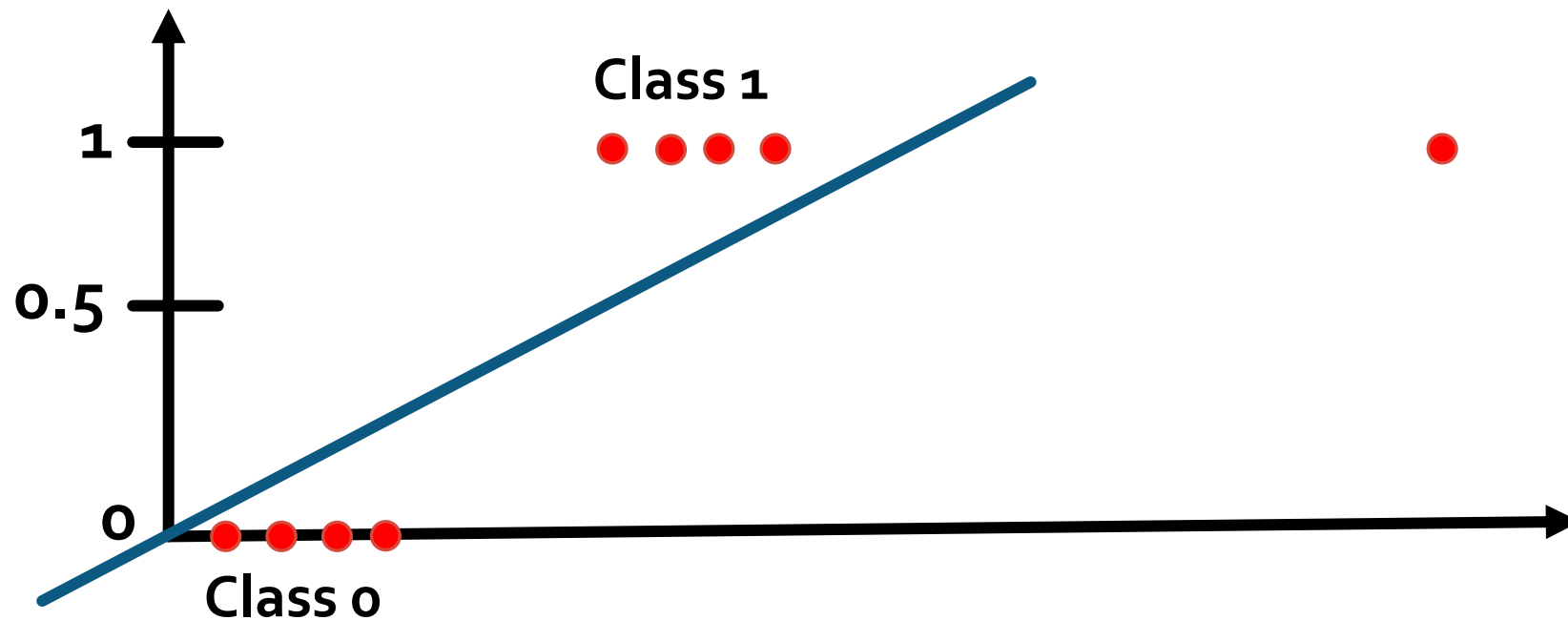


INTRODUCTION

LINEAR REGRESSION



WHY NOT APPLY LINEAR REGRESSION TO CLASSIFY?



The function $w^T X$ does not describe the expected value of $Y \in \{0, 1\}$

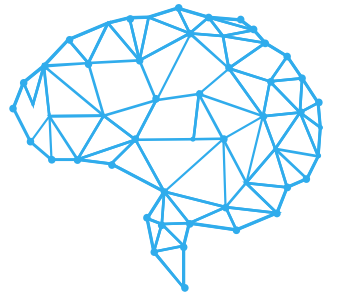


AI

BINARY CLASSIFICATION

LOGISTIC REGRESSION

HYPOTHESIS



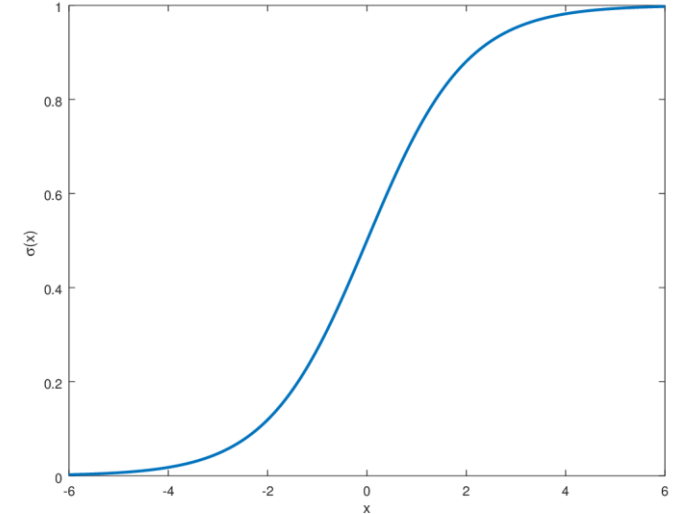
The **hypothesis** that we use for **linear regression** is not adequate to solve our **classification problem**.

A **hypothesis** that **better fits** the problem is proposed

$$h_w(x) = g(w^T X) = \frac{1}{1 + e^{-w^T X}}$$

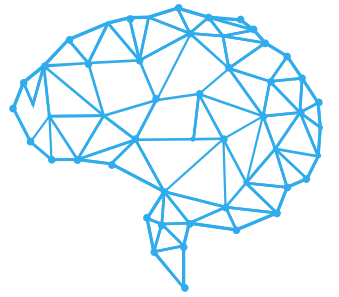
This function is called the **sigmoid** or **logistic function**.

It is **guaranteed** that $g(z) \in \{0, 1\}$.



LOGISTIC REGRESSION

SIGMOID FUNCTION



HOMEWORK

Prove that:

$$g'(z) = g(z)(1 - g(z))$$

LOGISTIC REGRESSION

ASSUMPTIONS



Now, we define the **same problem** of finding the best combination of **weights** that **fits** the **classification problem**.

The following **assumptions** are applied:

1. We assume that the **hypothesis** defines a **probability measure (Bernoulli)**:

$$P(y = 1/x; w) = h_w(x)$$

$$P(y = 0/x; w) = 1 - h_w(x)$$

$$P(y/x; w) = h_w(x)^y (1 - h_w(x))^{1-y}$$

$$L(w) = P(y/x; w)$$

LOGISTIC REGRESSION

ASSUMPTIONS



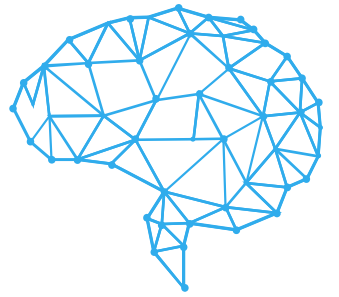
2. Assuming m **samples** are taken, we have that the **likelihood** can be written as:

$$L(\mathbf{w}) = \prod_{i=1}^m P(\mathbf{y}^{(i)} / \mathbf{x}^{(i)}; \mathbf{w})$$

$$L(\mathbf{w}) = \prod_{i=1}^m h_{\mathbf{w}}(\mathbf{x}^{(i)})^{y^{(i)}} (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)}))^{1-y^{(i)}}$$

LOGISTIC REGRESSION

LOGARITHMIC LOSS



The **logarithmic loss** can be written as:

$$l(\mathbf{w}) = \log \prod_{i=1}^m P(\mathbf{y}^{(i)} / \mathbf{x}^{(i)}; \mathbf{w})$$

$$l(\mathbf{w}) = \sum_{i=1}^m \mathbf{y}^{(i)} \log h_{\mathbf{w}}(\mathbf{x}^{(i)}) + (1 - \mathbf{y}^{(i)}) \log (1 - h_{\mathbf{w}}(\mathbf{x}^{(i)}))$$

$$l(\mathbf{w}) = \sum_{i=1}^m \mathbf{y}^{(i)} \log \left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right) + (1 - \mathbf{y}^{(i)}) \log \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{x}}} \right)$$

LOGISTIC REGRESSION

COST FUNCTION



The **cost function** can be expressed as:

$$J(\mathbf{w}) = -l(\mathbf{w})$$

$$J(\mathbf{w}) = -\sum_{i=1}^m y^{(i)} \log \left(\frac{1}{1 + e^{-\mathbf{w}^T \mathbf{X}}} \right) + (1 - y^{(i)}) \log \left(1 - \frac{1}{1 + e^{-\mathbf{w}^T \mathbf{X}}} \right)$$

LOGISTIC REGRESSION

COST FUNCTION



Interpreting the **cost function**:

$$J(w) = - \sum_{i=1}^m y^{(i)} \log h_w(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_w(x^{(i)}))$$

If $y = 0$	If $y = 1$
$J(y = 0, \hat{y}) = -\log(1 - \hat{y})$	$J(y = 1, \hat{y}) = -\log(\hat{y})$
$\lim_{\hat{y} \rightarrow 0} \log(1 - \hat{y}) \rightarrow 0$	$\lim_{\hat{y} \rightarrow 0} \log(\hat{y}) \rightarrow \infty$
$\lim_{\hat{y} \rightarrow 1} \log(1 - \hat{y}) \rightarrow \infty$	$\lim_{\hat{y} \rightarrow 1} \log(\hat{y}) \rightarrow 0$

LOGISTIC REGRESSION

OPTIMIZATION



HOW TO FIND THE BEST WEIGHTS W ?

LOGISTIC REGRESSION

GRADIENT DESCENT

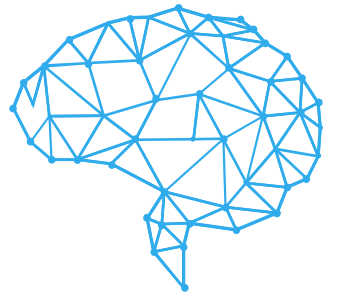


Finding the best combination of **weights** w by **gradient descent** fits the **classification problem**.

$$w := w - \alpha \nabla_w J(w)$$

LOGISTIC REGRESSION

GRADIENT DESCENT



HOMEWORK

Show that the **gradient** with respect to a **single training data** (x, y) and **weight** w_j is given by:

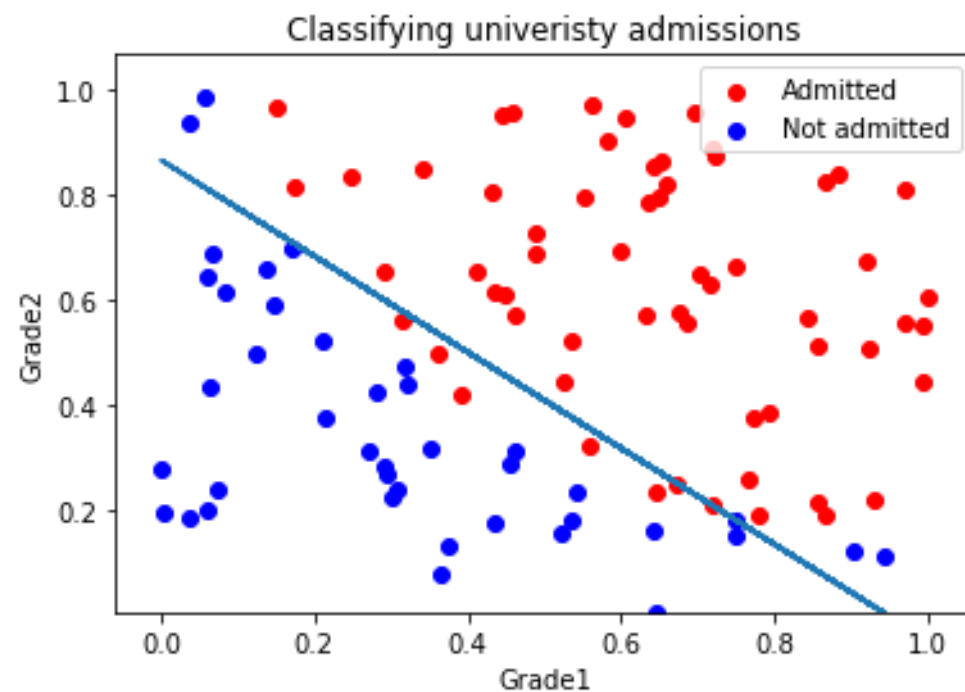
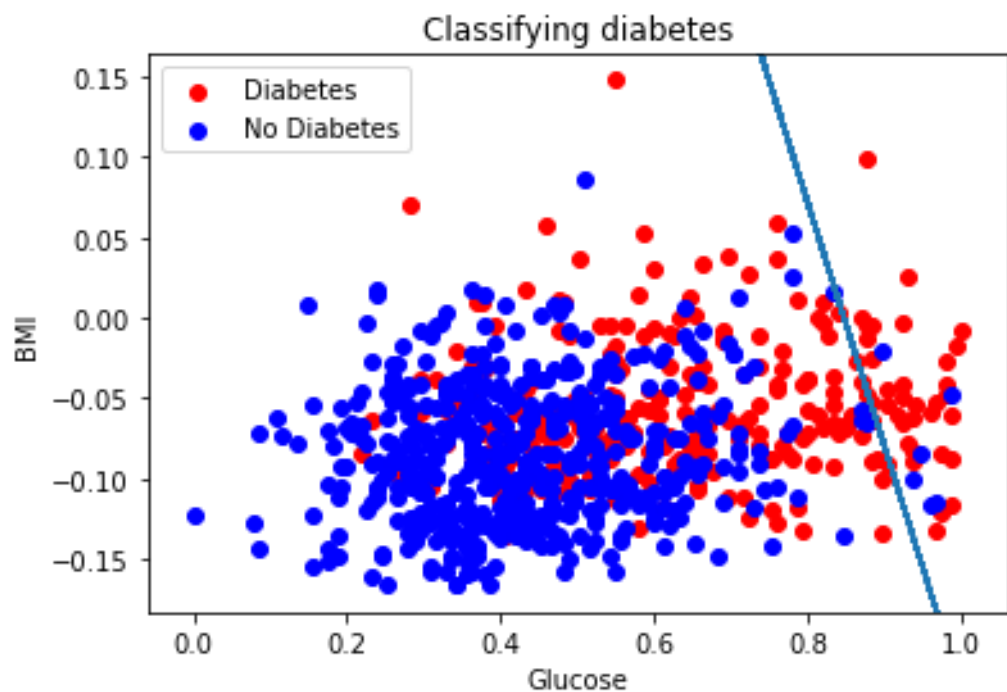
$$\frac{\partial}{\partial w_j} J(w) = (y - h_w(x))x_j$$

Where $x, w_j \in \mathbb{R}^n$ and $y \in \{0, 1\}$

**REMEMBER THE GRADIENT DESCENT IN LINEAR
REGRESSION → SAME RESULT**

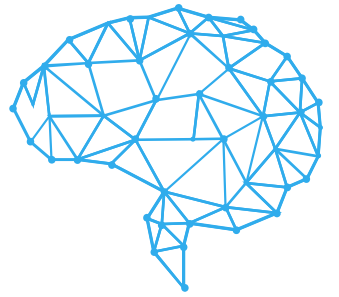
LOGISTIC REGRESSION

REAL EXAMPLE



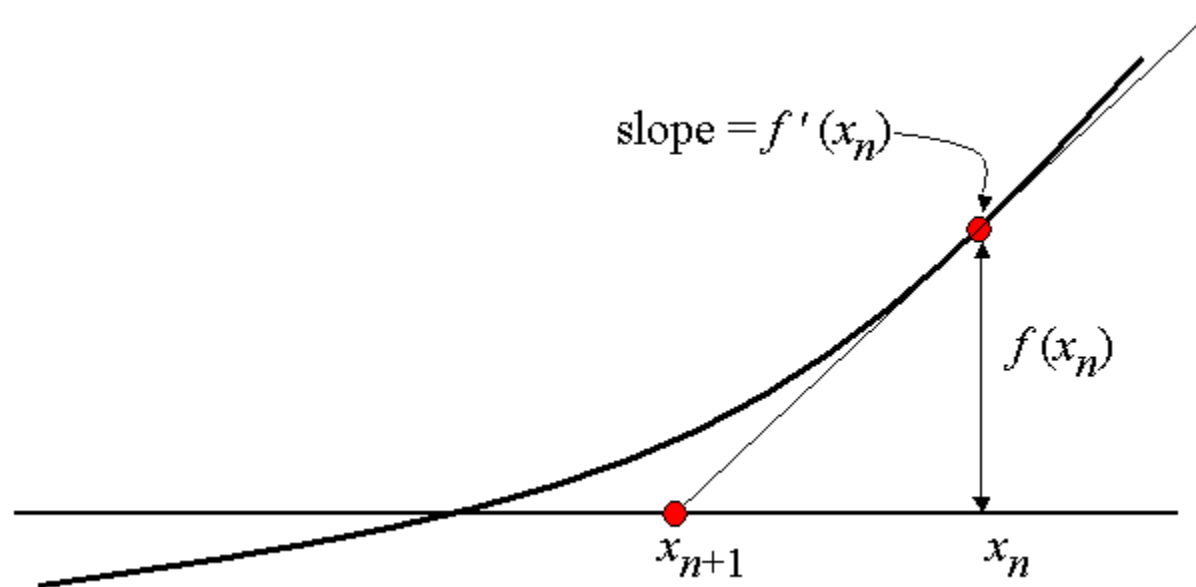
NEWTON'S METHOD

OPTIMIZATION



We **recall** from numerical methods, **Newton's method**, where the value of w is found for which $f(w) = 0$:

$$w := w - \frac{f(w)}{f'(w)}$$



NEWTON'S METHOD

OPTIMIZATION



If we want to find the **minimum** of the **cost function** $J(w)$, this will **correspond** to finding the **points** where $\nabla_w J(w) = 0$, meaning that we can use **Newton's method** (much faster):

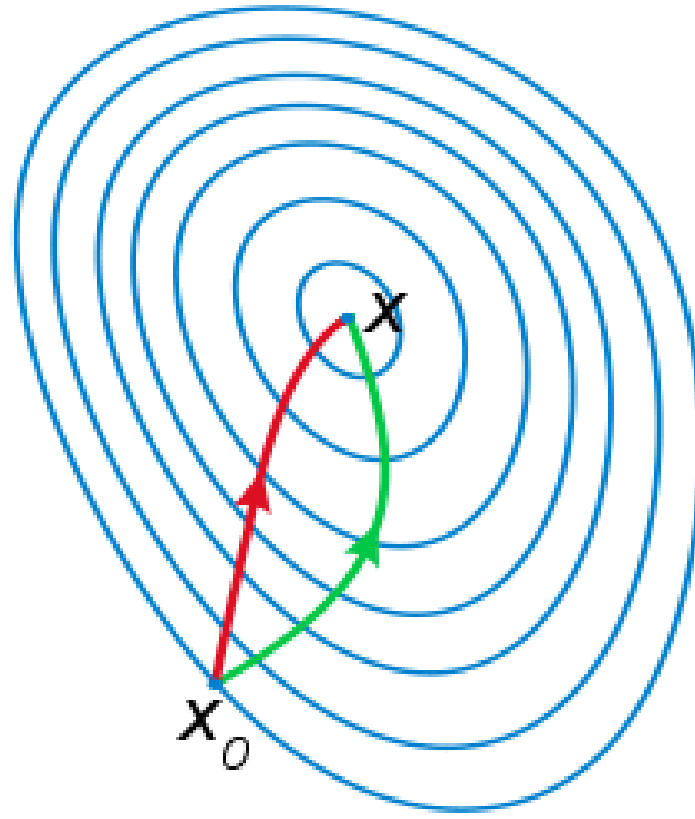
$$w := w - H^{-1} \nabla_w J(w)$$

where H^{-1} is the Hessian matrix, and ∇_w is the gradient.

$$H = \begin{bmatrix} \frac{\partial^2 J(w)}{\partial w_1^2} & \cdots & \frac{\partial^2 J(w)}{\partial w_1 \partial w_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial^2 J(w)}{\partial w_n \partial w_1} & \cdots & \frac{\partial^2 J(w)}{\partial w_n^2} \end{bmatrix}$$

NEWTON'S METHOD

OPTIMIZATION



Gradient descent
(Linear convergence)

Newton's method
(Quadratic convergence)

NEWTON'S METHOD OPTIMIZATION



**WHY NOT USE NEWTON'S METHOD INSTEAD OF GRADIENT
DESCENT?**



AI

GENERALIZED LINEAR MODELS

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY



We **recall** that we made the **following assumptions**:

LEAST SQUARES

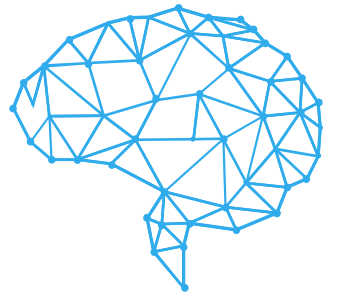
$y \in \mathbb{R} \sim \textit{Gaussian}$ (Linear regression)

LOGISTIC REGRESSION

$y \in \{0, 1\} \sim \textit{Bernoulli}$ (Binary classification)

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY



We say that a **probability distribution** belongs to the **exponential family** if it can be written as follows:

$$p(\mathbf{y}; \boldsymbol{\eta}) = b(\mathbf{y})e^{\left(\boldsymbol{\eta}^T T(\mathbf{y}) - a(\boldsymbol{\eta})\right)}$$

Where:

- $\boldsymbol{\eta}$ is named as the **natural** or **canonical parameter** of the distribution.
- $T(\mathbf{y})$ **sufficient statistic** (statistic that summarizes the complete information of a sample).
- $a(\boldsymbol{\eta})$ the **partition logistics function**.

NOTE: $e^{-a(\boldsymbol{\eta})}$ works as a **normalization constant** to ensure that $p(\mathbf{y}; \boldsymbol{\eta})$ integrates 1.

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY



$$p(\mathbf{y}; \boldsymbol{\eta}) = \mathbf{b}(\mathbf{y})e^{\left(\boldsymbol{\eta}^T \mathbf{T}(\mathbf{y}) - a(\boldsymbol{\eta})\right)}$$

If we **fix** \mathbf{a} , \mathbf{b} and \mathbf{T} , then we can say that we have a **distribution parametrized** only by $\boldsymbol{\eta}$, whereby varying $\boldsymbol{\eta}$ gives us **different distributions**.

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY EXAMPLES



We prove that the **Bernoulli distribution** is part of the **exponential family**:

$$p(y; \phi) = \phi^y (1 - \phi)^{1-y}$$

By **varying ϕ** we obtain **different distributions $p(y; \phi)$** .

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY EXAMPLES



We find:

$$T(y) = y$$

$$b(y) = 1$$

$$\eta = \log\left(\frac{\phi}{1-\phi}\right)$$

$$a(\eta) = -\log(1-\phi)$$

$$p(y; \phi) = e^{\left(\log\left(\frac{\phi}{1-\phi}\right)y + \log(1-\phi)\right)}$$

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY EXAMPLES



We prove that the **Gaussian distribution** (where σ does not matter) is part of the **exponential family** (so $\sigma^2 = 1$):

$$p(y; \mu) = \frac{1}{\sqrt{2\pi}} e^{-\frac{(y-\mu)^2}{2}}$$

HOMework 1

GENERALIZED LINEAR MODELS

CONSTRUCTING THE MODELS



We are going to **make** the **following assumptions**:

1. The output variable $y / X; w \sim FamExp(\eta)$
2. **Objective:** given a matrix of characteristics X , calculate $E(T(y); \eta)$ referred to as a **canonical response function**.

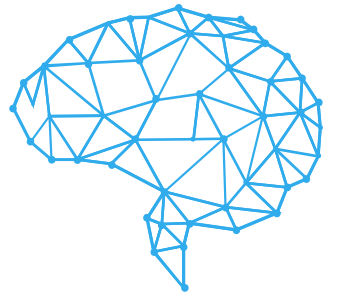
$$h(\eta) = E(T(y); \eta)$$

3. The **relationship** between η , w and X is defined to be **linear** (only if $\eta \in \mathbb{R}$):

$$\eta = w^T X$$

GENERALIZED LINEAR MODELS

CONSTRUCTING THE MODELS



EXAMPLE: Bernoulli

For fixed values X and w the **objective** is to **calculate**:

$$h(w) = E(T(y); \eta)$$

But we know that $T(y) = y$ so the **expected value** of a variable that is **distributed** as **Bernoulli** is:

$$E(y; \eta) = P(y = 1; \eta)$$

Similarly, we know from the exponential family that:

$$\phi = P(y = 1; \eta)$$

GENERALIZED LINEAR MODELS

EXPONENTIAL FAMILY EXAMPLES



It follows from the **definition** of the **exponential family** that the **natural parameter** η is expressed as:

$$\eta = \log \left(\frac{\phi}{1 - \phi} \right)$$

By **clearing** ϕ from the previous equality:

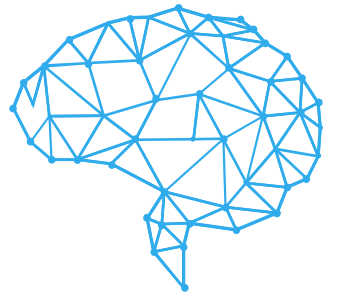
$$\eta = \log \left(\frac{\phi}{1 - \phi} \right)$$

The **following** is **obtained**:

$$\phi = \frac{1}{1 + e^{-\eta}}$$

GENERALIZED LINEAR MODELS

CONSTRUCTING THE MODELS



So the **expected value** is:

$$E(y; \eta) = \phi = \frac{1}{1 + e^{-\eta}}$$

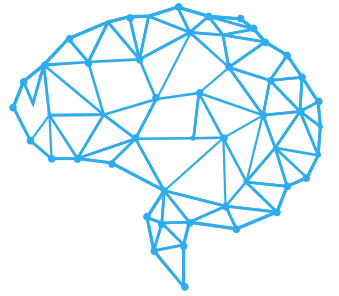
But thanks to the third assumption $\eta = w^T X$

$$E(y; w, X) = \phi = \frac{1}{1 + e^{-w^T X}}$$

**SIGMOID
FUNCTION**

GENERALIZED LINEAR MODELS

CONSTRUCTING THE MODELS



EXAMPLE: Gaussian

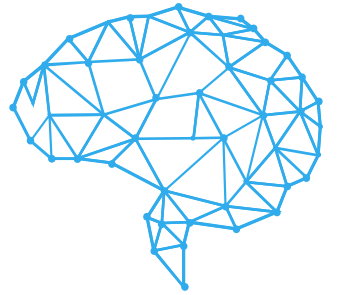
Prove that the canonical response function is equivalent to $\mathbf{w}^T \mathbf{X}$ assuming that $\mathbf{y} \sim \text{Gaussian}$ and that $\boldsymbol{\eta} = \mathbf{W}^T \mathbf{X}$

$$h(\mathbf{w}) = E(T(\mathbf{y}); \boldsymbol{\eta}) = \mathbf{w}^T \mathbf{X}$$

HOMEWORK 2

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



The following problem is defined, where the **output variable** $y \in \{1, \dots, k\}$ is distributed as a **multinomial function** $y \sim \text{Multinomial}$.

- Parameters: $\phi_1, \phi_2, \dots, \phi_{k-1}$
- $P(y = i) = \phi_i$
- $\phi_k = 1 - (\phi_1 + \phi_2 + \dots + \phi_{k-1})$

The **sufficient statistic** $T(y)$ is defined as a **vector**:

$$T(1) = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} T(2) = \begin{bmatrix} 0 \\ 1 \\ \vdots \\ 0 \end{bmatrix} \dots T(k-1) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 1 \end{bmatrix} T(k) = \begin{bmatrix} 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \in \mathbb{R}^{k-1}$$

Indicator function:

$$1\{True\} = 1 \quad 1\{False\} = 0$$

$$T(y)_i = 1\{y = i\}$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



The **multinomial distribution** is **expressed** in the form of the **exponential family**:

$$P(y) = \phi_1^{1\{y=1\}} \phi_2^{1\{y=2\}} \dots \phi_k^{1\{y=k\}}$$

$$P(y) = \phi_1^{T(y)_1} \phi_2^{T(y)_2} \dots \phi_k^{T(y)_{k-1}} \phi_k^{1 - \sum_{j=1}^{k-1} T(y)_j}$$

$$P(y) = e^{[T(y)_1 \log(\phi_1) + T(y)_2 \log(\phi_2) + \dots + (1 - \sum_{j=1}^{k-1} T(y)_j) \log(\phi_k)]}$$

$$P(y) = e^{[T(y)_1 \log\left(\frac{\phi_1}{\phi_k}\right) + T(y)_2 \log\left(\frac{\phi_2}{\phi_k}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_k}\right) + \log(\phi_k)]}$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



It is compared to the exponential family form:

$$P(y) = e^{\left[T(y)_1 \log\left(\frac{\phi_1}{\phi_k}\right) + T(y)_2 \log\left(\frac{\phi_2}{\phi_k}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_k}\right) + \log(\phi_k) \right]}$$

$$p(y; \eta) = b(y) e^{\left(\eta^T T(y) - a(\eta) \right)}$$

Therefore:

$$b(y) = 1$$

$$a(\eta) = -\log(\phi_k)$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



$$P(y) = e^{\left[T(y)_1 \log\left(\frac{\phi_1}{\phi_k}\right) + T(y)_2 \log\left(\frac{\phi_2}{\phi_k}\right) + \dots + T(y)_{k-1} \log\left(\frac{\phi_{k-1}}{\phi_k}\right) + \log(\phi_k) \right]}$$

$$p(y; \eta) = b(y) e^{\left(\eta^T T(y) - a(\eta) \right)}$$

Finally:

$$\eta = \begin{bmatrix} \log\left(\frac{\phi_1}{\phi_k}\right) \\ \vdots \\ \log\left(\frac{\phi_{k-1}}{\phi_k}\right) \end{bmatrix}$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



The **canonical response function** would be:

$$\eta_i = \log \left(\frac{\phi_i}{\phi_k} \right)$$

$$\eta_k = \log \left(\frac{\phi_k}{\phi_k} \right) = 0$$

$$\phi_k e^{\eta_i} = \phi_i$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



The **sum of probabilities** must be **equal** to 1:

$$\sum_{i=1}^k \phi_k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1$$

$$\phi_k \sum_{i=1}^k e^{\eta_i} = \sum_{i=1}^k \phi_i = 1$$

$$\phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}}$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



Substituting $\phi_k = \frac{1}{\sum_{i=1}^k e^{\eta_i}}$ in the canonical response function

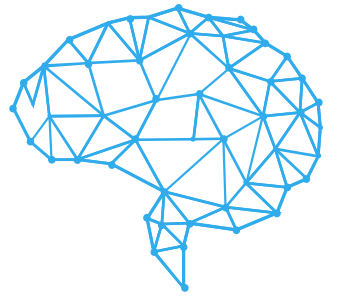
$$\phi_k e^{\eta_i} = \phi_i$$

$$\phi_i = \frac{e^{\eta_i}}{\sum_{j=1}^k e^{\eta_j}}$$

SOFTMAX
FUNCTION

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



Using the **third assumption** that $\eta_i = w_i^T X$:

$$\phi_i = \frac{e^{w_i^T X}}{\sum_{j=1}^k e^{w_j^T X}}$$

Knowing that the objective is to calculate $h(w) = E(T(y); \eta)$:

$$E\left(\begin{bmatrix} T(y)_1 \\ T(y)_2 \\ \vdots \\ T(y)_{k-1} \end{bmatrix}\right) = E\left(\begin{bmatrix} \mathbf{1}\{y = 1\} \\ \mathbf{1}\{y = 2\} \\ \vdots \\ \mathbf{1}\{y = k-1\} \end{bmatrix}\right) = \begin{bmatrix} \phi_1 \\ \phi_2 \\ \vdots \\ \phi_{k-1} \end{bmatrix}$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



The hypothesis $h(\mathbf{w})$ is **expressed** as a **vector**:

$$h(\mathbf{w}) = \begin{bmatrix} \frac{e^{w_1^T X}}{\sum_{j=1}^k e^{w_j^T X}} \\ \frac{e^{w_2^T X}}{\sum_{j=1}^k e^{w_j^T X}} \\ \vdots \\ \frac{e^{w_{k-1}^T X}}{\sum_{j=1}^k e^{w_j^T X}} \end{bmatrix} = \begin{bmatrix} P(y = 1) \\ P(y = 2) \\ \vdots \\ P(y = k - 1) \end{bmatrix}$$

which has as **components** the **probability** of **each class** i : $P(y = i)$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



We calculate the **logarithmic loss** and **cost function** for a **single training data**:

$$l(\mathbf{w}) = \sum_{l=1}^k \mathbf{1}(y = l) \log(p(y/x)) = \sum_{i=1}^k \mathbf{1}(y = l) \log \left(\frac{e^{w_l^T X}}{\sum_{j=1}^k e^{w_j^T X}} \right)$$

$$J(\mathbf{w}) = -l(\mathbf{w})$$

GENERALIZED LINEAR MODELS

SOFTMAX REGRESSION



We calculate the **logarithmic loss** and **cost function** for m training data:

$$l(\mathbf{w}) = \sum_{i=1}^m \sum_{l=1}^k \mathbf{1}(y^{(i)} = l) \log \left(\frac{e^{\mathbf{w}_l^T \mathbf{x}^{(i)}}}{\sum_{j=1}^k e^{\mathbf{w}_j^T \mathbf{x}^{(i)}}} \right)$$

$$J(\mathbf{w}) = -l(\mathbf{w})$$

THE NEWTON OR GRADIENT DESCENT
METHOD MAY BE USED



AI

**DISCRIMINATIVE AND
GENERATIVE MODELS**

DISCRIMINATIVE AND GENERATIVE MODELS DIFFERENCES



Discriminative models:

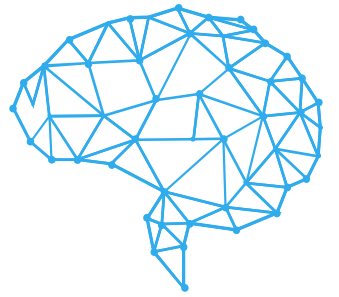
- The models that have been seen so far are called **discriminative models**, where learning about the **probability distribution** $p(y / x)$ is done **directly**.
- **Examples:** logistic regression and least squares.

Generative models:

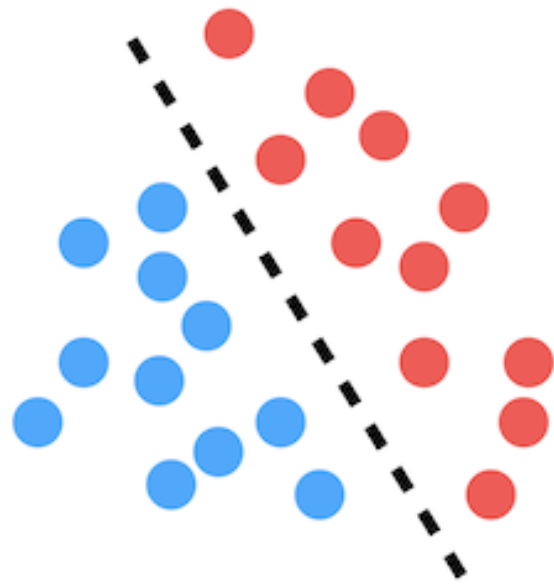
- Models that try to learn $p(x / y)$ and $p(y)$. For the case of **binary classification** the algorithm would learn **three different distributions** $p(x / y = 0)$, $p(x / y = 1)$ and $p(y)$.
- Thus, a new **data** would be **classified** by **comparing** it with both **distributions** $p(x / y = 0)$ and $p(x / y = 1)$.

DISCRIMINATIVE AND GENERATIVE MODELS

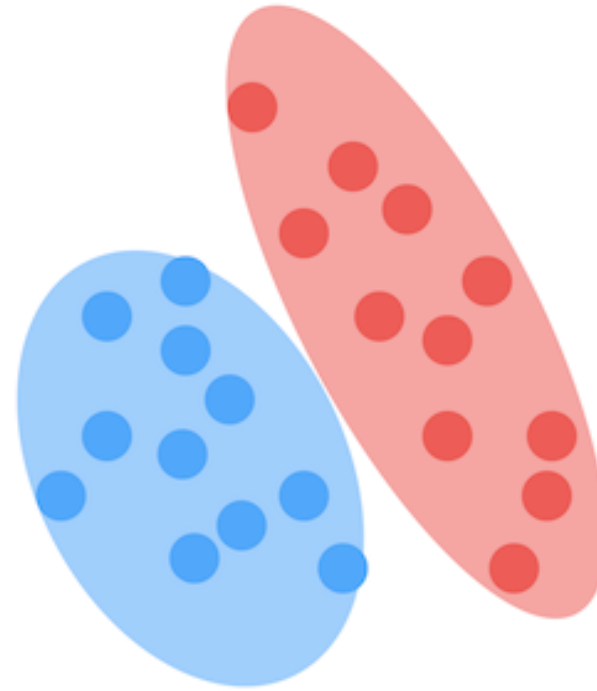
D I F F E R E N C E S



Discriminative



Generative



GENERATIVE MODELS

DESCRIPTION



A **generative model** models the **characteristics** that are **conditioned** by the **response variable** $p(x / y)$.

Using the **Bayes' Theorem**, $p(y / x)$ can be calculated.

$$p(y/x) = \frac{p(x/y)p(y)}{p(x)}$$

We have that the **denominator** can be **calculated** as follows (**binary classification**):

$$p(x) = p(x/y = 1)p(y = 1) + p(x/y = 0)p(y = 0)$$

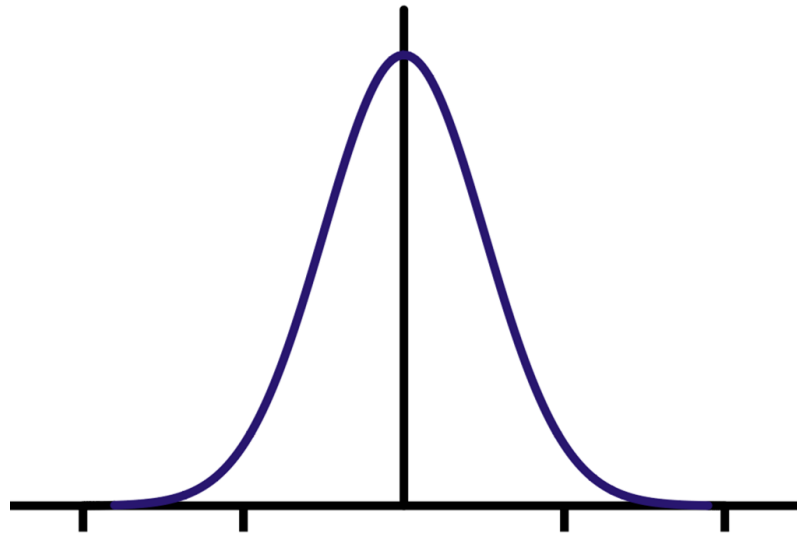
GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



In a **Gaussian discriminative analysis** we assume that:

$$p(x/y) \sim \text{Gaussian}$$



GENERATIVE MODELS

MULTIVARIATE GAUSSIANS



A **Gaussian distribution** of d dimensions is **parametrized** by:

- **Mean vector:** $\mu \in \mathbb{R}^d$
- **Covariance matrix:** $\Sigma \in \mathbb{R}^{d \times d}$

where Σ is **symmetric** and **positive semi-definite**. Its density can be calculated as:

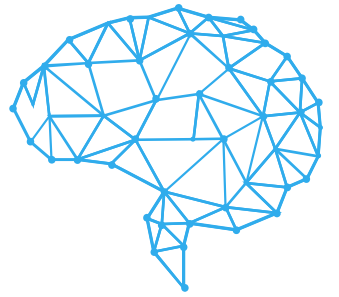
$$p(x; \mu, \Sigma) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} e^{\left(-\frac{1}{2}(x-\mu)^T \Sigma^{-1} (x-\mu)\right)}$$

Properties:

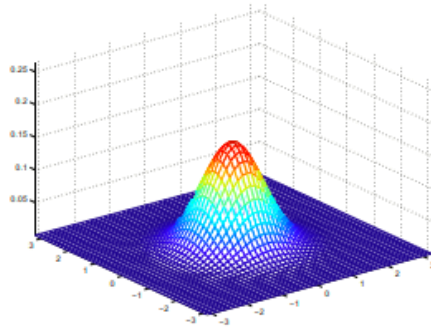
$$E[X] = \int_x x p(x; \mu, \Sigma) dx = \mu$$
$$\text{Cov}(X) = E[(X - E[X])(X - E[X])^T] = \Sigma$$

GENERATIVE MODELS

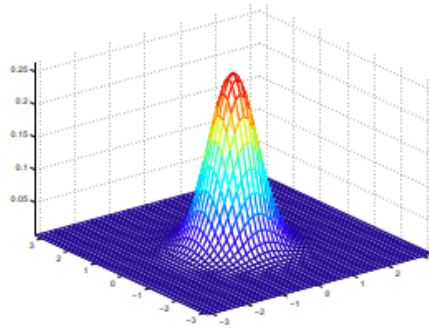
MULTIVARIATE GAUSSIANS



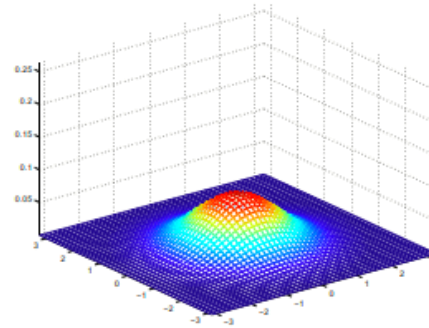
Example: Varying Σ .



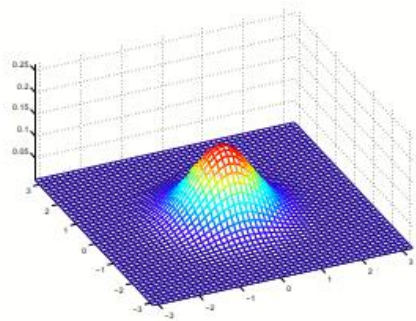
$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



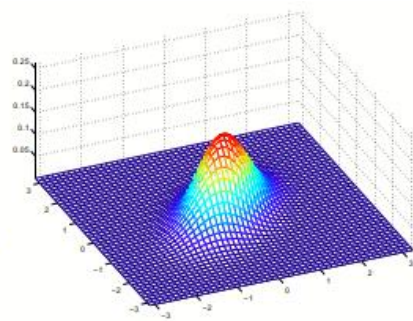
$$\Sigma = \begin{bmatrix} 0.6 & 0 \\ 0 & 0.6 \end{bmatrix}$$



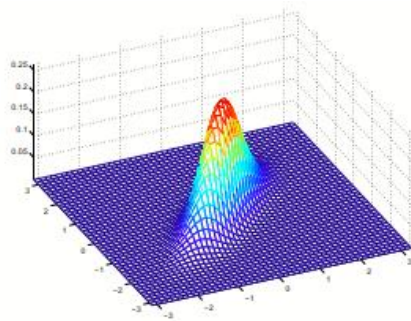
$$\Sigma = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$$



$$\Sigma = \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}$$



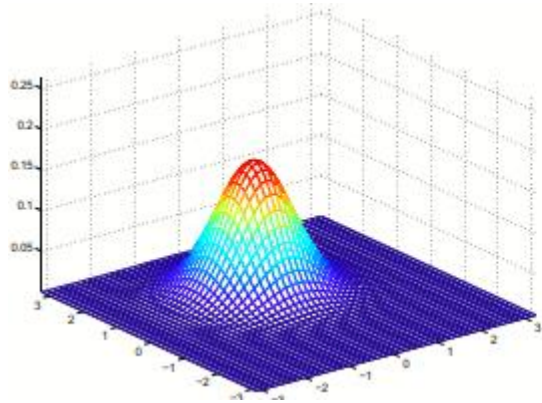
$$\Sigma = \begin{bmatrix} 1 & 0.8 \\ 0.8 & 1 \end{bmatrix}$$

GENERATIVE MODELS

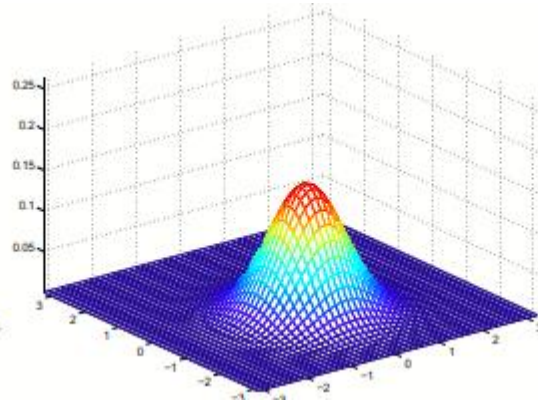
MULTIVARIATE GAUSSIANS



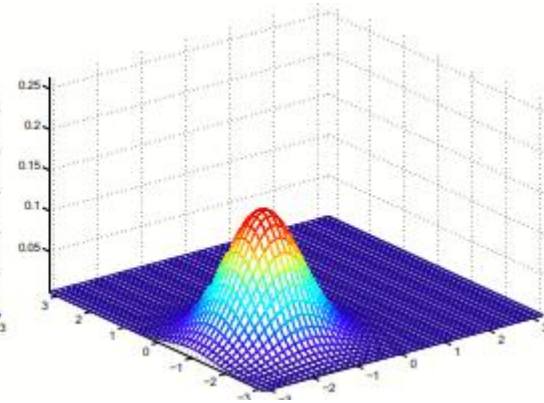
Example: Varying μ .



$$\mu = \begin{bmatrix} 1 \\ 0 \end{bmatrix}$$



$$\mu = \begin{bmatrix} -0.5 \\ 0 \end{bmatrix}$$



$$\mu = \begin{bmatrix} -1 \\ -1.5 \end{bmatrix}$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Gaussian discriminative analysis solves the **binary classification problem** where:

$$\begin{aligned}x &\in \mathbb{R}^n \\y &\in \{0, 1\}\end{aligned}$$

The **assumptions** of the **model** are:

$$y \sim \text{Bernoulli}(\phi)$$

$$x / y = 0 \sim N(\mu_0, \Sigma)$$

$$x / y = 1 \sim N(\mu_1, \Sigma)$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



The **distributions** are given by:

$$p(y) = \phi^y (1 - \phi)^{1-y}$$

$$p(x|y=0) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_0)^T \Sigma^{-1} (x - \mu_0) \right)$$

$$p(x|y=1) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp \left(-\frac{1}{2} (x - \mu_1)^T \Sigma^{-1} (x - \mu_1) \right)$$

The **parameters** of the **distributions** would be: $\phi, \mu_0, \mu_1, \Sigma$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



The **joint likelihood** of the data would be given by:

$$l(\phi, \mu_k, \Sigma) = \prod_{i=1}^m p(x^{(i)}, y^{(i)}; \phi, \mu_k, \Sigma)$$

$$l(\phi, \mu_k, \Sigma) = \prod_{i=1}^m p(x^{(i)} / y^{(i)}; \mu_k, \Sigma) p(y^{(i)}; \phi)$$

$$\log l(\phi, \mu_k, \Sigma) = \log \prod_{i=1}^m p(x^{(i)} / y^{(i)}; \mu_k, \Sigma) p(y^{(i)}; \phi)$$

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \log(p(x^{(i)} / y^{(i)}; \mu_k, \Sigma) + \log p(y^{(i)}; \phi)$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



The **joint likelihood** of the data would be given by:

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \log(p(x^{(i)} / y^{(i)}; \mu_k, \Sigma) + \log p(y^{(i)}; \phi)$$

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \log\left(\frac{1}{(2\pi)^{\frac{d}{2}} |\Sigma|^{\frac{1}{2}}}\right) + \log \left[e^{\left(-\frac{1}{2}(x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k)\right)} \right] + \log(\Phi^{y^{(i)}} (1 - \Phi)^{1-y^{(i)}})$$

$$\log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \log\left(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}\right) - \frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) + y^{(i)} \log(\phi) + (1 - y^{(i)}) \log((1 - \phi))$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Maximizing with respect to ϕ :

$$\nabla_{\phi} \log l(\phi, \mu_k, \Sigma) = \nabla_{\phi} \sum_{i=1}^m \mathbf{y}^{(i)} \log \left(\frac{\phi}{1 - \phi} \right) + \log((1 - \phi)) = 0$$

$$\nabla_{\phi} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \left(\frac{\mathbf{y}^{(i)}}{\phi(1 - \phi)} - \frac{1}{1 - \phi} \right) = 0$$

$$\sum_{i=1}^m \frac{\mathbf{y}^{(i)}}{\phi(1 - \phi)} = \frac{m}{1 - \phi}$$

$$\nabla_{\phi} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \frac{\mathbf{y}^{(i)}}{m} = \sum_{i=1}^m \frac{1(\mathbf{y}^{(i)} = \mathbf{k})}{m} = \phi = \frac{\text{Number of data in class } k}{\text{Total data}}$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Maximizing with respect to μ_k :

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \nabla_{\mu_k} \sum_{i=1}^m -\frac{1}{2} (\mathbf{x}^{(i)} - \mu_k)^T \Sigma^{-1} (\mathbf{x}^{(i)} - \mu_k) = 0$$

Applying the property $\nabla_x x^T A x = 2Ax$:

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \Sigma^{-1} (\mathbf{x}^{(i)} - \mu_k) = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \Sigma^{-1} \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu_k) = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \mathbf{x}^{(i)} - \sum_{i=1}^m \mu_k = 0$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



We are **only interested** in the $\mathbf{x}^{(i)}$ that belong to **class k**:

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \mathbf{x}^{(i)} - \sum_{i=1}^m \mu_k = 0$$

$$\nabla_{\mu_k} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k}) \mathbf{x}^{(i)} - \sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k}) \mu_k = 0$$

$$\mu_k = \frac{\sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k}) \mathbf{x}^{(i)}}{\sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k})} = \frac{\text{Sum of } \mathbf{x} \text{ that belong to class } k}{\text{Class } k \text{ data number}}$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Maximizing with respect to Σ^{-1} :

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^m \log\left(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}\right) - \frac{1}{2} (x^{(i)} - \mu_k)^T \Sigma^{-1} (x^{(i)} - \mu_k) = 0$$

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^m \log\left(\frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}}\right) - \frac{1}{2} (x^{(i)} - \mu_k)^T (x^{(i)} - \mu_k) \Sigma^{-T} = 0$$

Because Σ is symmetrical to $\Sigma^{-1} = \Sigma^{-T}$:

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \nabla_{\Sigma^{-1}} \sum_{i=1}^m -\frac{d}{2} \log(2\pi) + \frac{1}{2} \log(|\Sigma^{-1}|) - \frac{1}{2} (x^{(i)} - \mu_k)^T (x^{(i)} - \mu_k) \Sigma^{-1} = 0$$

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m \nabla_{\Sigma^{-1}} \frac{1}{2} \log(|\Sigma^{-1}|) - \nabla_{\Sigma^{-1}} \frac{1}{2} (x^{(i)} - \mu_k)^T (x^{(i)} - \mu_k) \Sigma^{-1} = 0$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Applying the property $\nabla_x b^T x = b$:

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \frac{1}{2} \sum_{i=1}^m \nabla_{\Sigma^{-1}} [\log(|\Sigma^{-1}|)] - (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T = 0$$

Applying the property $\nabla_x \log(|X|) = X^{-T}$:

$$\nabla_{\Sigma^{-1}} \log l(\phi, \mu_k, \Sigma) = \sum_{i=1}^m (\Sigma^{-1})^{-T} - (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T = 0$$

$$m \Sigma^T - \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T = 0$$

$$\Sigma^T = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Because Σ is symmetrical to $\Sigma^T = \Sigma$:

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)} - \boldsymbol{\mu}_k)^T$$

VARIANCE OF K-CLASS DATA

Remembering that it is an **outer product**: $\Sigma \in \mathbb{R}^{n \times n}$

$$\Sigma = \begin{bmatrix} \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}_1 - \boldsymbol{\mu}_k)^2 & \dots & \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}_1 - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)}_n - \boldsymbol{\mu}_k) \\ \vdots & \ddots & \vdots \\ \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}_n - \boldsymbol{\mu}_k)(\mathbf{x}^{(i)}_1 - \boldsymbol{\mu}_k) & \dots & \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)}_n - \boldsymbol{\mu}_k)^2 \end{bmatrix}$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



Once we have **found** the **parameters**:

$$\phi = \sum_{i=1}^m \frac{1(\mathbf{y}^{(i)} = \mathbf{k})}{m}$$

$$\mu_k = \frac{\sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k}) \mathbf{x}^{(i)}}{\sum_{i=1}^m 1(\mathbf{y}^{(i)} = \mathbf{k})}$$

$$\Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu_k)(\mathbf{x}^{(i)} - \mu_k)^T$$

We can **start making predictions**.

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



To make a prediction we calculate the **posterior probability** $P(y = 1 / x)$ with the **new data** x :

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

We assign a threshold: if $P(y = 1 / x) \geq 0.5$ then it belongs to class 1, otherwise it belongs to class 0. That is, we are **maximizing the probability**:

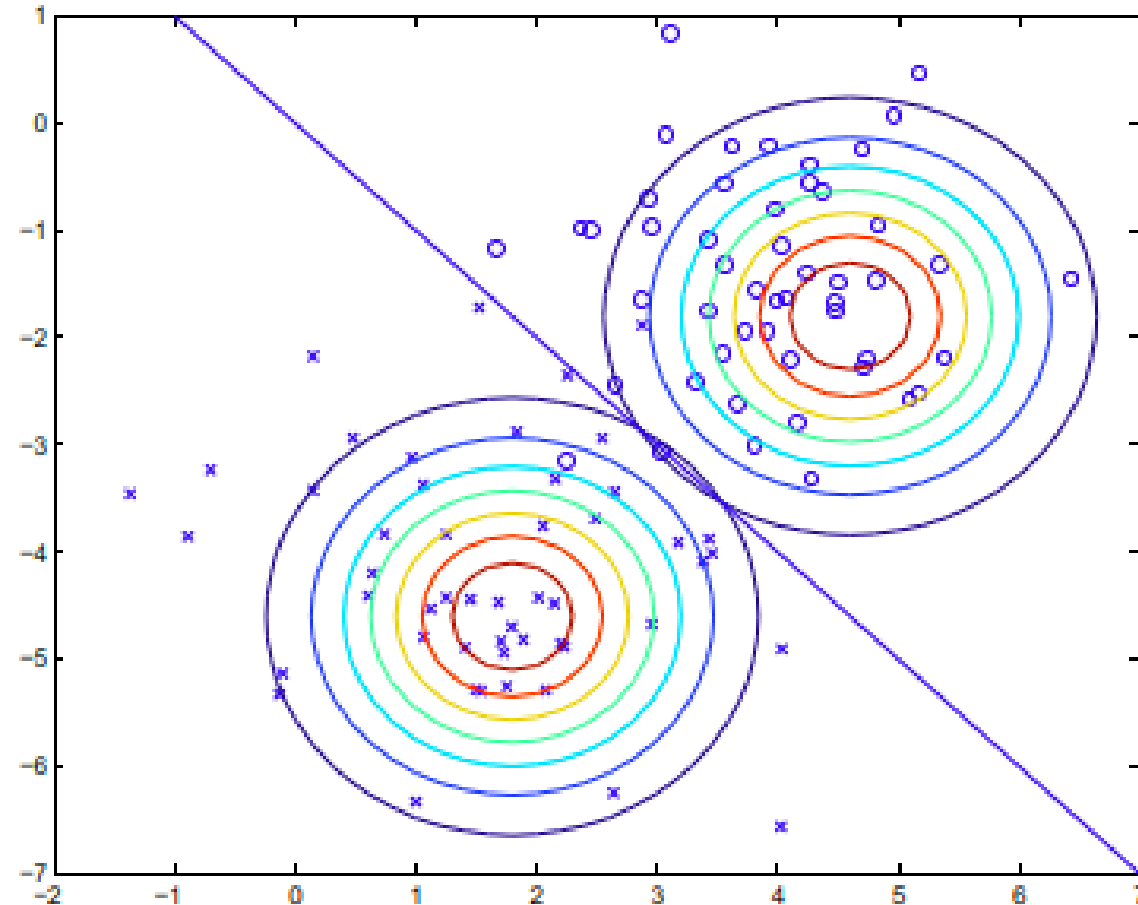
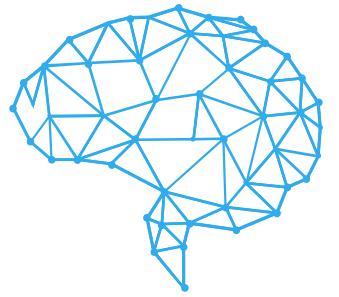
$$\operatorname{argmax}_y p(y/x) = \operatorname{argmax}_y \frac{p(x/y)p(y)}{p(x)}$$

If $p(y = 1) = p(y = 0) = 0.5$ we have:

$$\operatorname{argmax}_y p(y/x) = \operatorname{argmax}_y p(x/y)$$

GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



GENERATIVE MODELS

GAUSSIAN DISCRIMINATIVE ANALYSIS



SUMMARY:

1. Compute the parameters for each Gaussian (in binary classification it would be two curves) from the data.

$$\phi, \mu, \Sigma$$

2. Calculate the probabilities

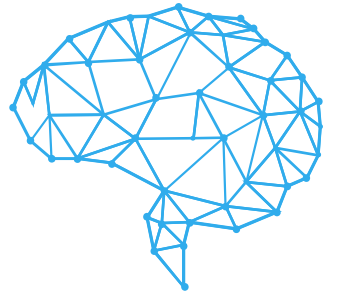
$$P(x/y = k) \quad P(y) \quad P(x)$$

3. Use Bayes' Theorem to make a new prediction.

$$p(y = k/x) = \frac{p(x/y = k)p(y = k)}{p(x)}$$

GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



If we analyze the probability $p(y = 1 / x)$, in **Bayes' Theorem**, we can realize that we obtain the **same logistic regression curve**.

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x/y = 1)p(y = 1) + p(x/y = 0)p(y = 0)} \left(\frac{\frac{1}{p(x/y = 1)p(y = 1)}}{\frac{1}{p(x/y = 1)p(y = 1)}} \right)$$

$$p(y = 1/x) = \frac{1}{1 + \frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)}}$$

GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\frac{1}{2}(x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) + \frac{1}{2}(x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)\right)} \left(\frac{1 - \phi}{\phi}\right)$$

$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\frac{1}{2}(x^{(i)} - \mu_0)^T \Sigma^{-1} (x^{(i)} - \mu_0) + \frac{1}{2}(x^{(i)} - \mu_1)^T \Sigma^{-1} (x^{(i)} - \mu_1)\right)} \left(e^{\log\left(\frac{1 - \phi}{\phi}\right)}\right)$$

$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\frac{1}{2}\Sigma^{-1} \left[(x^{(i)} - \mu_0)^T (x^{(i)} - \mu_0) - (x^{(i)} - \mu_1)^T (x^{(i)} - \mu_1)\right] + \log\left(\frac{1 - \phi}{\phi}\right)\right)}$$

$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\frac{1}{2}\Sigma^{-1} \left[\sum_{j=1}^n (x^{(i)}_j - \mu_0)^2 - (x^{(i)}_j - \mu_1)^2\right] + \log\left(\frac{1 - \phi}{\phi}\right)\right)}$$

GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\frac{1}{2}\Sigma^{-1}\left[\sum_{j=1}^n -2\mu_0 x^{(i)}_j + \mu_0^2 + 2\mu_1 x^{(i)}_j - \mu_1^2\right] + \log\left(\frac{1-\phi}{\phi}\right)\right)}$$

Because $x^{(i)}_0 = 1$

$$\frac{p(x/y = 0)p(y = 0)}{p(x/y = 1)p(y = 1)} = e^{\left(-\left[\sum_{j=1}^n \frac{1}{2}\Sigma^{-1}(\mu_0 - \mu_1)(x^{(i)}_j)\right] + \left[\frac{1}{2}\Sigma^{-1}(\mu_0^2 - \mu_1^2) + \log\left(\frac{1-\phi}{\phi}\right)\right](x^{(i)}_0)\right)}$$

Comparing $w^T X$:

$$w^T = \begin{bmatrix} -\frac{1}{2}\Sigma^{-1}(\mu_0^2 - \mu_1^2) - \log\left(\frac{1-\phi}{\phi}\right) \\ \frac{1}{2}\Sigma^{-1}(\mu_0 - \mu_1) \\ \vdots \\ \frac{1}{2}\Sigma^{-1}(\mu_0 - \mu_1) \end{bmatrix}$$

GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



Therefore:

$$p(y = 1/x) = \frac{1}{1 + e^{-\left(\left[\sum_{j=1}^n \frac{1}{2} \Sigma^{-1}(\mu_0 - \mu_1)(x^{(i)}_j)\right] + \left[-\frac{1}{2} \Sigma^{-1}(\mu_0^2 - \mu_1^2) + \log\left(\frac{1-\phi}{\phi}\right)\right](x^{(i)}_0)\right)}} = \boxed{\frac{1}{1 + e^{-w^T X}}}$$

Where w^T :

$$w^T = \begin{bmatrix} -\frac{1}{2} \Sigma^{-1}(\mu_0^2 - \mu_1^2) - \log\left(\frac{1-\phi}{\phi}\right) \\ \frac{1}{2} \Sigma^{-1}(\mu_0 - \mu_1) \\ \vdots \\ \frac{1}{2} \Sigma^{-1}(\mu_0 - \mu_1) \end{bmatrix}$$

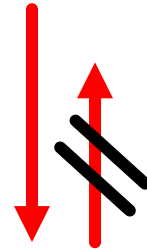
GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



In a general fashion:

Likelihood function
 $x/y \sim \text{ExponentialFamily}(\eta)$



Posterior Distribution
 $P(y/x) = \text{sigmoid function}$

THAT'S WHY LOGISTIC REGRESSION IS USED → WORKS FOR MANY TYPES
OF ASSUMPTIONS

GENERATIVE MODELS

COMPARISON DISCRIMINATIVE MODELS



A **generative model** will have **better performance** than a discriminative one, if the assumption that we made of the shape of the distribution $P(x / y)$ holds for the real data (more information is provided to the algorithm).

Otherwise, if the data does not behave as we assumed, the **discriminative model** will have **better results**, because even though our assumptions were not as accurate, **logistic regression** works for many different assumptions.

GENERATIVE	DISCRIMINATIVE
Better performance when the $P(x / y)$ distribution is known	Better performance when $P(x / y)$ is unknown (robust to incorrect assumptions)
Asymptotic efficiency	There may be a better model
Needs little data	Needs a lot of data



DISCRIMINATIVE AND GENERATIVE MODELS

Naive Bayes Classifier

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



BERNOULLI MULTIVARIATE EVENT MODEL:

For the case of the **Gaussian discriminative analysis**, it was the case that $x^{(i)}_j \in \mathbb{R}$, that is, they were **continuous values**.

For cases, such as **text classification** ("desired mail" and "spam mail") the **x vectors** can have a very **high dimensionality** (the number of words is immense).

$$x \in \{0, 1\}^{5000}$$

2^{5000} **parameters** would be needed.

$$x = \begin{bmatrix} 1 \\ 0 \\ 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{bmatrix} \begin{matrix} a \\ \text{aardvark} \\ \text{aardwolf} \\ \vdots \\ \text{buy} \\ \vdots \\ \text{zygmurgy} \end{matrix}$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



SOLUTION: naive Bayes assumption (very strong)

We suppose that $x^{(i)}_j$ are **conditionally independent** given y .

That is, if a text is known to be “**spam**” $y = 1$, the fact that the word $x_{2087} = \text{“buy”}$ appears in the text does not affect beliefs about the appearance of any other word, such as $x_{39831} = \text{“price”}$.

$$p(x_1, \dots, x_j, \dots, x_n/y) = p(x_1)p(x_2/x_1, y)p(x_3/x_1, x_2, y) \dots = p(x_1/y)p(x_2/y) \dots = \prod_{j=1}^n p(x_j/y)$$

EVEN WHEN THE ASSUMPTION IS NOT TRUE, THE ALGORITHM PERFORMS GOOD IN MANY APPLICATIONS

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



As the **Bayes classifier** is a **generative model**, it is of interest to model the distributions $p(x_j/y)$ and $p(y)$. In particular, we want to **maximize the joint likelihood**:

$$l(\phi_y, \phi_{j/y=0} \phi_{j/y=1}) = \prod_{i=1}^m p(x^{(i)}, y^{(i)})$$

Where the **assumptions** are:

$$p(x/y) = p(x_1, \dots, x_j, \dots, x_n/y) = \prod_{j=1}^n p(x_j/y)$$

$$x_j / y = 0 \sim \text{Bernoulli}(\phi_{j/y=0})$$

$$x_j / y = 1 \sim \text{Bernoulli}(\phi_{j/y=1})$$

$$y \sim \text{Bernoulli}(\phi_j)$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



The **result** of the **maximum likelihood** estimation gives:

$$\phi_{j/y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)}, y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} = \frac{\text{\textit{\#times the word } j \text{ is repeated in "spam" }}}{\text{\textit{total "spam" mail}}}$$

$$\phi_{j/y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)}, y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} = \frac{\text{\textit{\#times the word } j \text{ is repeated in "desired" }}}{\text{\textit{total "desired" mail}}}$$

$$\phi_y = \frac{\sum_{i=1}^m 1\{y^{(i)} = 1\}}{m} = \frac{\text{\textit{total spam mail}}}{\text{\textit{total mail}}}$$

If the **characteristics** $x_j^{(i)}$ take **more values**, they can be **modeled as multinomials**.

If the **characteristics** $x_j^{(i)}$ take **continuous values**, they are **discretized in intervals**.

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



To make a **new prediction** we use **Bayes' theorem**:

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)}$$

Where:

$$p(x/y = 1) = \prod_{j=1}^n p(x_j/y = 1) = \prod_{j=1}^n (\phi_j/y = 1)$$

$$p(y) = \phi_y = \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\}}{m}$$

$$p(x) = \prod_{j=1}^n p(x_j/y = 1) p(y = 1) + \prod_{j=1}^n p(x_j/y = 0) p(y = 0)$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



What happens if a word appears in an email that the algorithm has never seen in another email?

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



The **probability** that it will assign, of seeing the word in **either of the two emails** will be **0**.

$$\phi_{j/y=1} = \frac{\sum_{i=1}^m 1\{x_j^{(i)}, y^{(i)} = 1\}}{\sum_{i=1}^m 1\{y^{(i)} = 1\}} = \frac{\text{\textit{\#times the word } j \text{ is repeated in "spam" }}}{\text{\textit{total "spam" mail}}} = \frac{0}{m_{neg}}$$

$$\phi_{j/y=0} = \frac{\sum_{i=1}^m 1\{x_j^{(i)}, y^{(i)} = 0\}}{\sum_{i=1}^m 1\{y^{(i)} = 0\}} = \frac{\text{\textit{\#times the word } j \text{ is repeated in "desired" }}}{\text{\textit{total "desired" mail}}} = \frac{0}{m_{pos}}$$

When making the **prediction** we will have an **incongruity**:

$$p(y = 1/x) = \frac{p(x/y = 1)p(y = 1)}{p(x)} = \frac{0}{0}$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER

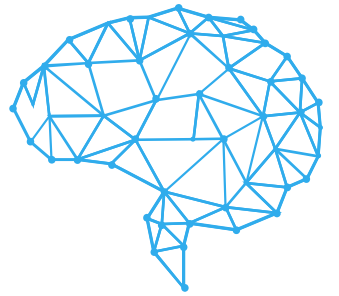


**STATISTICALLY IT IS NOT ADEQUATE TO SAY THAT THE
PROBABILITY OF AN EVENT IS ZERO JUST BECAUSE YOU HAVE
NOT SEEN IT IN YOUR DATA**

NOTE: LAPLACE AND THE SUN

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



LAPLACE SMOOTHING:

The **solution** is to **change our estimate** (general multinomial case where $y = \{1, 2, \dots, k\}$):

$$p(y = 1) = \frac{(\# \text{"1"s} + 1)}{(\# \text{"0"s} + 1) + (\# \text{"1"s} + 1)}$$

$$p(y = j) = \phi_y = \frac{1 + \sum_{i=1}^m \mathbf{1}\{y^{(i)} = j\}}{m + k}$$

When calculating the **other estimators** we have:

$$\phi_{j/y=1} = \frac{1 + \sum_{i=1}^m \mathbf{1}\{x_j^{(i)}, y^{(i)} = 1\}}{2 + \sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\}}$$

$$\phi_{j/y=0} = \frac{1 + \sum_{i=1}^m \mathbf{1}\{x_j^{(i)}, y^{(i)} = 0\}}{2 + \sum_{i=1}^m \mathbf{1}\{y^{(i)} = 0\}}$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



MULTINOMIAL EVENT MODEL:

Let's see the case when $x_j = \{1, 2, \dots, k\}$.

For example, now the value of x_j represents the **position** of the **word** in the **dictionary**, and the index j indicates the **position** of the **word** in the **mail**.

So n now represents the **length of the email** (and varies according to each email).

The **assumptions** would be:

$$p(x/y) = p(x_1, \dots, x_j, \dots, x_n/y) = \prod_{j=1}^n p(x_j/y)$$

$$x_j / y = 0 \sim \text{Multinomial}(\phi_{j/y=0})$$

$$x_j / y = 1 \sim \text{Multinomial}(\phi_{j/y=1})$$

$$y \sim \text{Bernoulli}(\phi_j)$$

GENERATIVE MODELS

NAIVE BAYES CLASSIFIER



The **parameters**, when calculating the **maximum likelihood** would be (with **Laplace smoothing**):

$$\phi_{k/y=1} = \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{1}\{x_j^{(i)} = k, y^{(i)} = 1\} + 1}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\} d_i + n_i + k} = \frac{\text{\textit{\# of times that the word } k \text{ appears in mails } ND}}{\text{\textit{total words in mails } ND}}$$

$$\phi_{k/y=0} = \frac{\sum_{i=1}^m \sum_{j=1}^{n_i} \mathbf{1}\{x_j^{(i)} = k, y^{(i)} = 0\} + 1}{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = 0\} d_i + n_i + k} = \frac{\text{\textit{\# of times that the word } k \text{ appears in mails } D}}{\text{\textit{total words in mails } D}}$$

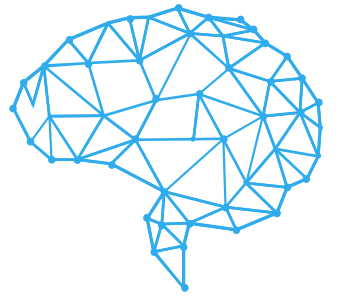
$$\phi_y = \frac{\sum_{i=1}^m \mathbf{1}\{y^{(i)} = 1\} + 1}{m}$$



BINARY EVALUATION METRICS

BINARY EVALUATION METRICS

DEFINITIONS

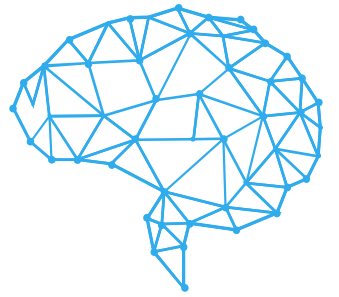


False positive (error type I):

- **Statistics:** A true null hypothesis is incorrectly rejected.
- **Machine learning:** the model **predicts** that the **class** of a training data $p(y / x) = 1$ is **positive**, when the **reality** is that the **data belonged** to the **negative class** $y = 0$.

BINARY EVALUATION METRICS

DEFINITIONS

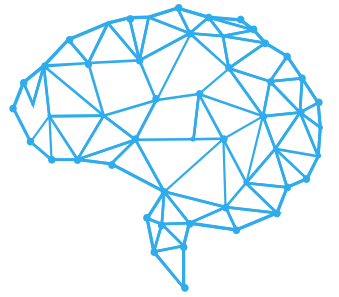


False negative (error type II):

- **Statistics:** The **false null hypothesis** is **incorrectly accepted**.
- **Machine learning:** the model **predicts** that the **class** of a **training data** $p(y / x) = 0$ is **negative**, when the **reality** is that the **data belonged** to the **positive class** $y = 1$.

BINARY EVALUATION METRICS

DEFINITIONS

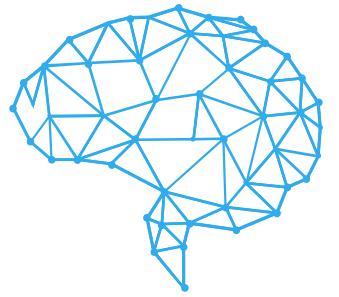


True positive:

- **Statistics:** The **false null hypothesis** is **correctly accepted**.
- **Machine learning:** the model **predicts** that the **class** of a **training data** $p(y / x) = 1$ is **positive**, when the **reality** is that the **data belonged** to the **positive class** $y = 1$.

BINARY EVALUATION METRICS

DEFINITIONS



True negative:

- **Statistics:** The **true null hypothesis** is **correctly accepted**.
- **Machine learning:** the model **predicts** that the **class** of a **training data** $p(y / x) = 0$ is **negative**, when the **reality** is that the **data belonged** to the **negative class** $y = 0$.

BINARY EVALUATION METRICS

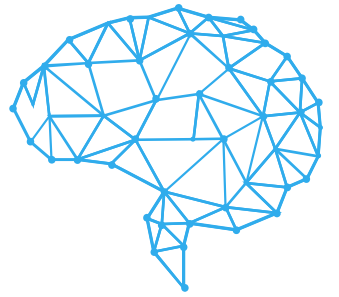
C O N F U S I O N M A T R I X



		PREDICTIVE VALUES	
		POSITIVE (1)	NEGATIVE (0)
ACTUAL VALUES	POSITIVE (1)	TP	FN
	NEGATIVE (0)	FP	TN

BINARY EVALUATION METRICS

S E N S I T I V I T Y



Sensitivity (“true positive rate”, “recall”): measures the **performance** of the model to **predict** the **positive class** $y = 1$.

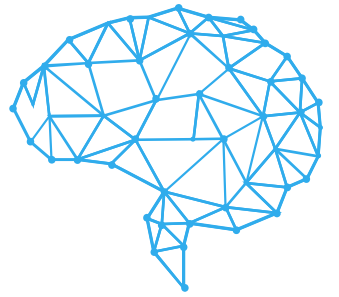
It is **difficult** for a **sensitive algorithm** to make a **mistake** in **predicting** the **positive class**.

Still, **high sensitivity** can be **accompanied** by many **false positives**.

$$Sens = \frac{TP}{TP + FN}$$

BINARY EVALUATION METRICS

S P E C I F I C I T Y



Specificity (“true negative rate”): measures the **performance** of the model to **predict** the **negative class** $y = 0$.

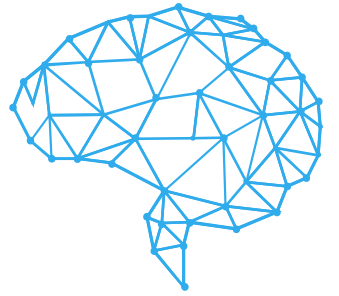
It is **difficult** for a **specitive algorithm** to make a **mistake** in **predicting** the **negative class**.

Still, **high specitivity** can be **accompanied** by many **false negatives**.

$$Spec = \frac{TN}{TN + FP}$$

BINARY EVALUATION METRICS

A C C U R A C Y

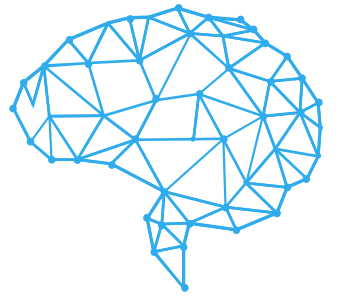


Accuracy (not to be confused with "positive predictive value" or "precision") measures the performance of the model in a general way. How much the predictions deviate from the actual values.

$$Acc = \frac{TP + TN}{TP + TN + FN + FP}$$

BINARY EVALUATION METRICS

POSITIVE PREDICTIVE VALUE



The "**positive predictive value**" or "**precision**" measures the **fraction of positives** that were **correctly predicted**.

$$PPV = \frac{TP}{TP + FP}$$

BINARY EVALUATION METRICS

F 1 S C O R E



The **F1 score** is the **harmonic mean** (average for rates) **between** the **recall** and the **PPV**:

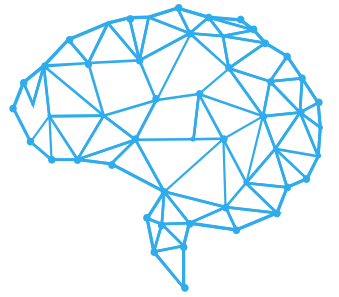
$$Sens = \frac{TP}{TP + FN}$$

$$PPV = \frac{TP}{TP + FP}$$

$$F1 = \frac{2}{Sens^{-1} + PPV^{-1}}$$

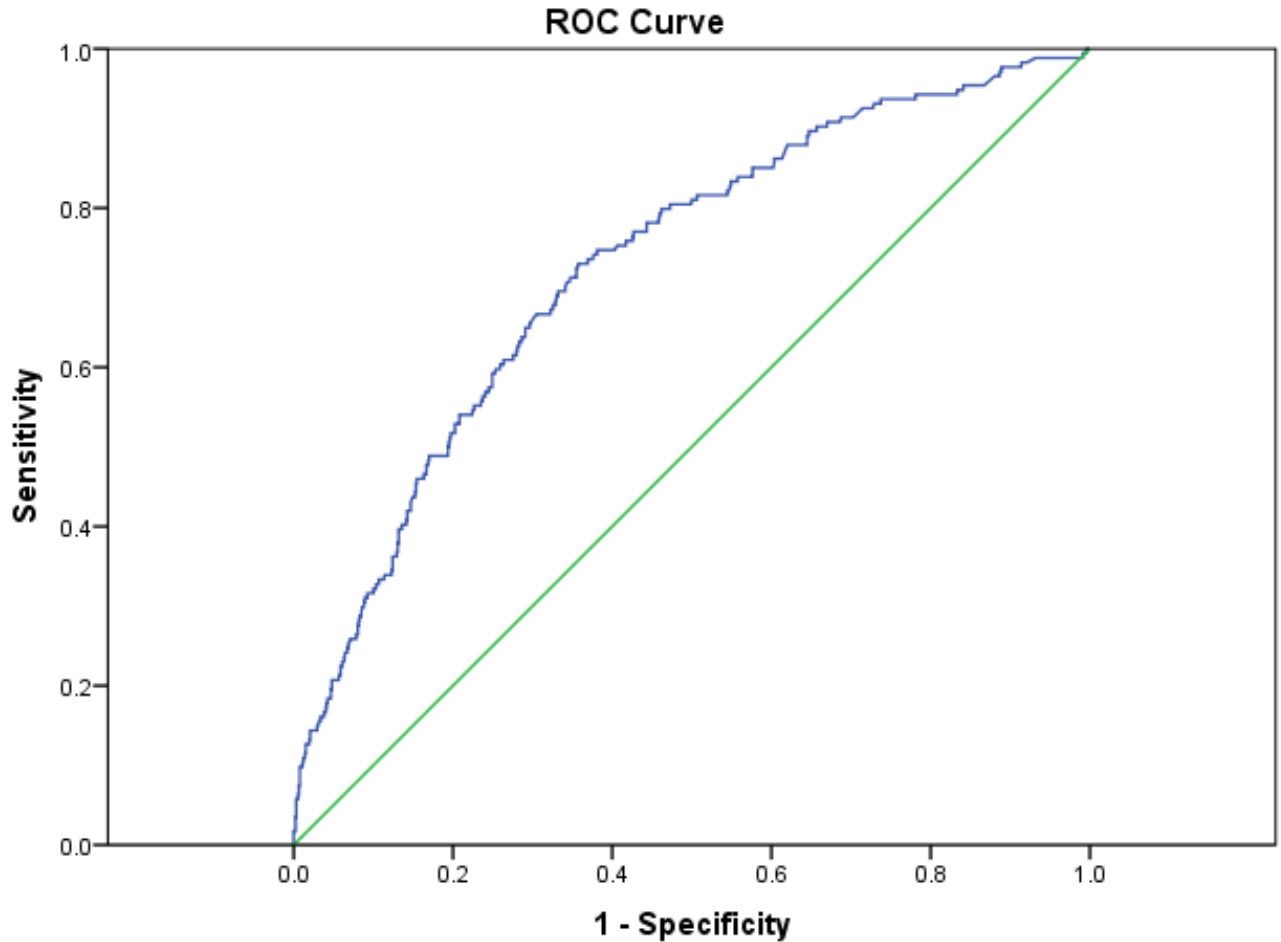
BINARY EVALUATION METRICS

R O C C U R V E



FPR vs TPR

The threshold is varied



Diagonal segments are produced by ties.

BINARY EVALUATION METRICS

R O C C U R V E

