# Task – 2

## Sub Task – 1

- Prakhar Tripathi

A detailed mathematical explanation of the Gaussian Discriminative Analysis model describing the principles, assumptions, and equations involved and explaining how the model learns from the data and makes predictions.

# Contents

# Generative Learning Algorithms

In general models train on p(y| x; $\theta$), the conditional distribution of y over x. For instance, logistic regression modelled p(y| x; $\theta$), as $h_\theta(x) = g(\theta^T x)$ where g is the sigmoid function.

Lets take an example where we have to distinguish between car(y=0) and aeroplane(y=1) based on some features of vehicles. If we are given a training set, an algorithm like Logistic Regression tries to find a straight line that is a decision boundary that separates the regions of cars and aeroplanes on the graph.

But let's see a different approach where we can build a model of what cars look like. Then looking at aeroplane, we can build a separate model of what aeroplanes look like. Finally, to classify a new vehicle, we can match the new vehicle against the cars model, and match it against the elephant model, to see whether the new vehicle looks more like the aeroplanes or more like the cars we had seen in the training set.

Algorithms that try to learn p(y| x) directly (such as logistic regression), or algorithms are called discriminative learning algorithms. Here, well talk about algorithms that instead try to model p(x| y) and p(y). These algorithms are called generative learning algorithms. For instance, if y indicates whether an example is a car (0) or an aeroplane (1), then p(x|y = 0) models the distribution of car features, and p(x|y = 1) models the distribution of aeroplanes features.

After modelling p(y) (called the class priors) and p(x| y), our algorithm can then use Bayes rule to derive the posterior distribution on y given x:

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

Here, the denominator is given by p(x) = p(x| y = 1)p(y = 1) + p(x| y = 0)p(y = 0)

Actually, if were calculating p(y| x) in order to make a prediction, then we don't actually need to calculate the denominator, since

$$arg \max_{y} p(y|x) = arg \max_{y} \frac{p(x|y)p(y)}{p(x)}$$
$$= arg \max_{y} p(x|y)p(y)$$

# Gaussian Discriminative Analysis

In this model, well assume that p(x| y) is distributed according to a multivariate normal distribution. Let's talk briefly about the properties of multivariate normal distributions before moving on to the GDA model itself.

## The multivariate Normal Distribution

The multivariate normal distribution in d-dimensions, also called the multi variate Gaussian distribution, is parameterized by a mean vector $\mu \in \mathbb{R}^d$ and a covariance matrix $\Sigma \in \mathbb{R}^{d \times d}$ where $\Sigma \geq 0$ is symmetric and positive semi-definite. Also written $\aleph(\mu, \Sigma)$, its density is given by

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

In the above equation, $|\Sigma|$ denotes the determinant of $\Sigma$.

## The Gaussian Discriminant Analysis Model

When we have a classification problem in which the input features x are continuous-valued random variables, we can then use the Gaussian Discriminant Analysis (GDA) model, which models p(x| y) using a multivariate normal distribution. The model is:

$$y \sim Bernoulli(\phi)$$
$$x|y = 0 \sim \aleph(\mu_0, \Sigma)$$
$$x|y = 1 \sim \aleph(\mu_1, \Sigma)$$

Writing out the distributions this is

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

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

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

The log-likelihood of the data is given by

$$l(\phi, \mu_0, \mu_1, \Sigma) = \log \prod_{i=1}^{n} p(x^{(i)}, y^{(i)}; \phi, \mu_0, \mu_1, \Sigma)$$

$$= \log \prod_{i=1}^{n} p(x^{(i)}|y^{(i)}; \phi, \mu_0, \mu_1, \Sigma) \, p(y^{(i)}; \phi)$$

By maximizing with respect to the parameters, we find the maximum like likelihood estimate of the parameters:

$$\phi = \frac{1}{n} \sum_{i=1}^{n} 1\{y^{(i)} = 1\}$$

$$\mu_0 = \frac{\sum_{i=1}^{n} 1(\{y^{(i)} = 0\}x^{(i)}}{1\{y^{(i)} = 0\}}$$

$$\mu_1 = \frac{\sum_{i=1}^{n} 1(\{y^{(i)} = 1\}x^{(i)}}{1\{y^{(i)} = 1\}}$$

$$\Sigma = \frac{1}{n} \sum_{i=1}^{n} (x^{(i)} - \mu_{y^{(i)}})(x^{(i)} - \mu_{y^{(i)}})^T$$

# Prediction and Output

The two models can be built by these values of mu and variances. Then each data's probability is calculated in both classes. Then the one with the bigger value is allotted the class.