# Early Detection Of Alzheimer's Disease Using Integrated Machine Learning and Probabilistic Distillation

**SIDDHARTH SHYAMSUNDER**

A thesis submitted for the degree of Master of Science in

BIG DATA AND TEXT ANALYTICS

Supervisor: Dr. Ian Daly

School of Computer Science and Electronic Engineering

University of Essex

August 2019

# Acknowledgement

The desire to have the best and get myself closer to perfection has made me try to achieve perfection. Ability and Ambition are not enough to succeed. Success of any project depends solely on support, guidance, encouragement received from the guide and our family and friends well wishes.

Gratitude is often the hardest emotion to express and often one doesn't find the exact words to convey what one feels. I am pleased to express my deepest gratitude to Dr. Prof Ian Daly who opened all the floodgates of his immense knowledge continuous support, encouragement and advice to help bring reality to my dream to make this project successful.

I would also like to thank the Alzheimer's Disease Neuro-Imaging Initiative(ADNI) for providing me with the required data, and expertise in understanding the data, without which going ahead with my project wouldn't have been possible.

And last but not least I would like to thank our Head of Department, for his timely co-operation and help. Thank you for all those who have directly or indirectly contributed towards making my project successful.

# ABSTRACT

In this research project, we have made use of a novel technique to have a computer system aid in the detection of alzheimer's disease. In our techniques we make use of a 2 phase architecture concerning the use of fMRI Image information to create an MRI Biomarker, and then using this MRI Biomarker along with cognitive information to create an aggregate biomarker. To provide a more specific study, we have focused on the Hippocampal sub region of the brain, which is the region of the brain that aids in the formation of new memories and retention of spatial memories. Using principal Component analysis(PCA), we have looked into creating Eigen characteristic vectors, that decrease the dimensionality thus helping ensure increase in the performance of machine learning techniques on the image data. Feature Selection is gathered from a selected region of interest for the image data, and from an Alzheimer's Disease Neuroimaging Initiative(ADNI) Challenge dataset for cognitive information. The PCA output along with the cognitive information are then subjected to a pure ensemble training technique like Random Forest Method, along with a probabilistically distilled ensemble training, integrating Random Forest Metods with other Standard ML Training. The results have proven an accuracy of 97% using our 2 phase model in the detection of AD.

# Table of Contents

# List Of Figures

# List Of Tables

# 1. Introduction

## 1.1. Overview:

Approximately 44 million people around the world suffer from some form of disease related to Alzheimer's. This makes Alzheimer's disease one of the most deadly diseases among the aged population, especially with people over the age of 65. One of the most prominent issues that are being faced by this disease is that the Alzheimer's disease is the root cause for the neurodegeneration of the human brain. And as far as research has dwelled into finding a cure for this disease, many attempts have provided neuroscientists with unsatisfactory results.

Despite, the fact that this disease is incurable at the moment, many attempts have been made to detect an early onset of this disease, which would allow patients to look forward towards better lifestyle measures, that can be heeded. In this report, we attempt to build a system ot rather an aggregate biomarker that attempts to analyze radiological samples in other words, sanples of fMRI samples to first distinguish between samples that follow the standard aging pattern and samples that are in between a normal aging brain and a brain with Alzheimer's disease. This in between stage of a brain between the normal aging pattern and an Alzheimer's brain is called Mild Cognitive Impairment(MCI). In the $2^{nd}$ stage, we will compare MCI Samples with that of the Alzheimer's to have the system learn to classify between an aging brain and an Alzheimer's disease infected brain.

In this section, we will have an overview into what the further sections in the report will contain. Furthermore, we will look at the hypothesis that we have considered while stating our problem statement. We will look further into what our system aims to accomplish along with a brief description of our problem.

In the $2^{nd}$ section, we will look into a detailed background into the domain under consideration, in our case the study of the subtle differences between a normally aging brain and a brain affected by Alzheimer's disease. Furthermore, we will look into a brief description of fMRI image processing along with ways, other experts in the area have approached the application of Machine Learning and Neural Networks to detect the abnormalities in the aging process of the brain. Based on the approaches provided by experts in the area, we will look critically discuss into what can be achieved, and what can be taken forward through our study.

In the $3^{rd}$ Section, we will look into an in detailed view of how we try to achieve the targets provided in the $2^{nd}$ section, not only will we look into how, we planned through the entire process, but we will also look into the software requirements that were taken under consideration, the features that were considered under our experiments, the results that were observed, and the limitations that were faced through our model.

In the $4^{th}$ section, we look into the how we evaluate the results obtained in section 3. We also look into some of the testing parameters and evaluation reports we generate after evaluation.

The last few sections, will look into the conclusions, references and appendices or extra information that need to be learnt in order to understand our system.

## 1.2. Research Hypothesis

When the idea to have the system automatically detect if a patient suffers from Alzheimer's disease was considered, it was important to frame a precise, testifiable question that gives a clear idea as to what we wish to achieve. Moreover the question that we have phrased, should be simplistic enough to give the general idea, of what our results should aim to depict. Looking at all these parameters under consideration, we decided to ask ourselves, "Is it possible for a computer system, given a set of functional brain scans, to automatically detect

2

with as little human intervention as possible, if a particular patient would likely suffer from Alzheimer's disease in the long run or not?".

Considering the above mentioned parameters, we will be performing a number of experiments right from processing an fMRI Image, along with creating an intermediate MRI Biomarker along with an aggregate biomarker through making use of Machine Learning Tools to aid with our decision making. A detailed information on the fMRI image processing as well as the machine learning algorithms that we intend to use will be discussed in the further sections to prove if our above question is testifiable and precise.

## 1.3. Aim

As can be seen from the earlier sections, we aim to design a system that can aid in the detection of alzheimer's disease on patients by making use of an aggregate Biomarker which will take in the properties of a patient's MRI Scans along with some other cognitive information such as age, Mini mental State examination report and so on to identify, if a patient's brain is aging in the normal process else if there is an abnormality in the aging process. By the end of our experiments we wish to obtain 2 labeled outputs whether a patient will have Alzheimer's disease or not, and furthermore, we also aim to achieve a high level of accuracy in the detection of Alzheimer's disease.

Our aim is to consider fMRI samples, taken from a set of 50+ patient, adjust by removing the standard aging parameters, and then preprocess the images, and study the major functional areras of the brain, especially the role of regions like the Hippocampus, in the formative stage of the Alzheimer's disease. Through comparisons of different volumetric features in the hippocampus, we will be able to gather if a patient, is likely to be affected by a stable Mild Cognitive Impairment(sMCI) in other words, if the abnormality in the aging process in the brain would traverse abnormally but at a rate, wherein Alzheimer's disease is not likely to

3

progress or in any case if the patient is likely to be affected by progressive Mild Cognitive Impairment(pMCI) wherein the aging process is so abnormal that a patient is likely to get affected by Alzheimer's disease.

Furthermore in a $2^{nd}$ stage process, we plan to consider critical cognitive information, in other words details like the results of an MMSE examination and so on that might play a key role in further deciding if an sMCI patient is likely to get affected by Alzheimer's disease.

# 2. Background And Literature

In this section, a detailed background and the domain related knowledge has been provided. Under the Literature Survey section, more of the technical aspects in relating Machine Learning, Deep Learning and the Neuro imaging domain which have been carried out by previous experts in the area have been provided. The last section, will have us, look at our aim discued in the previous chapter, and evaluate and discuss past work, and how, we intend to implement past work to evaluate better results.

## 2.1. Background:

A brain is a complex and integral part of the human body, as it works as the heart of the central nervous system. The brain works through means of billions of nerve cells called neurons, each communication information to each other through means of synapses. An abnormality in the proper functioning of the brain, including death of neuron cells or damage to synapses, can make a terrible difference to the standard functioning of the whole human body[1]. Such abnormality leads to one of the most leading causes of death in the aged population across the world called Alzheimer's disease. It is a neurodegenerative progressive and irreversible disease[2]-[8], that affects the entire brain in destroying neuron cells, in a highly abnormal way, so that the entire functionality of the brain begins to get affected, especially during the later stages. Some of the major causes of the disease is in the formation of intracellular neurofibrillary tangles and extracellular β-Amyloid plaques[1]–[9], which creates an abnormal flow towards the aging process of a brain by blocking the synapses, which play the essential role of helping neuron cells communicate with one another. The below figure gives a more clear view on how neurofibrillary tangles and β-Amyloid plaques play a role in affecting the standard aging process of the brain.

*Figure 1: The effect of how neurofibrillary tangles and β-Amyloid plaques cause AD.*

The discontinuity in the aging process of the brain can be calculated through a middle stage or a prodromal stage called Mild Cognitive Impairment(MCI), which lies between the normally aging or a Healthy Controlled(HC) brain and a completely affected brain with AD[1]-[4], [7]–[12], [13]–[22], [23]–[26]. This middle MCI stage can be further subdivided into the early or stable Mild Cognitive Impairment, a stage that lies between the normal aging process of the brain and which gives a low risk ratio of having the HC Brain progress into a fully affected Alzheimer's disease brain, and the late or progressive MCI(pMCI), which is a stage that lies between the HC brain, and which is more likely to progress into Alzheimer's disease. As of today, there is no practical cure to Alzheimer's disease, however, efforts and study have been taken, into detecting the onset of this disease, or an abnormality of the standard aging process from an early stage, using information from radiological images like fMRI [1]-[4], [6]–[14], [16]–[21], [24],[26], [27], PET Scans[1], [2], [9]–[13], [15], [20], [23], [24], [27] and Cerebral Spinal Fluid(CSF)[1], [8], [10], [12]-[14], [18], [20], [27] so that further steps can be taken to delay the abnormality rate. The main disadvantages or rather

6

efforts that are needed to be taken into consideration while analyzing the brain samples is the uniqueness of a person's brain structure and that must be taken into account while analyzing MR Images.



*Figure 2: Characteristic patterns between Normal Control, sMCI, pMCI, and AD brain samples.*

Some of the analysis that are used on the radiological images include making use of technologies like Machine Learning[3]-[5], [7]-[9], [11], [14], [17], [19], [22], [23], [25]-[27] Deep Learning[6], [7], [16], [20], [24] and Neural Networks[25]-[27], which allows a system, to automatically classify images into whether the patient is likely to get affected with Alzheimer's disease or not.

Another major factor that has been thought to be a major factor in deciding upon early detection of Alzheimer's disease is the identification of variance in the Hippocampal volume[2], [4], [22], [24], [27]. The Hippocampus is a seahorse shaped formation in the brain, which plays a major role in the formation of new memories and the identification of

spatial memories. The change in hippocampal volume has been attributed to be vulnerable in the earliest stages of Alzheimer's disease. Studies have believed that change in hippocampal volume has been a more convenient way on identification of abnormality in the brain aging process, even more so than cognitive tests including Mini Mental State Examinations[2], [4], [14], [22], [23] and so on. In addition to simple volumetric analysis, inclusion of voxel based morphometry(VBM)[1], [7], [25], [27] or the analysis, that has helped in proving the concentration of tissue loss in the hippocampal region, in a major marker to the identification to the level of progress in AD. The main advantage of making use of the VBM technique is that major emphasis is not on a particular region of interest, rather that of the entire brain structure. Although VBM helps to identify bilateral grey matter loss, it cannot still be accurate to identify regional disparities of hippocampal atrophy.



*Figure 3: Image showing atrophy in Hippocampal region.*

Now that we have clearly established some of the more rudimentary causes of Alzheimer's disease. As the report will progress, more emphasis has been provided, on how hippocampal regional atrophy is calculated, or in other words, through what means, can a system learn to correctly identify the differences between a healthy controlled brain and a brain affected by AD. To this extent, there is a need to make use of Neuro Imaging. Some of the most famous types of Neuroimaging, include making use of MRI or Magnetic Resonance Imaging, PET or

Positron Emission Tomography or Cerebrospinal fluids(CSF) imaging. Out of these 3 different techniques, it is widely accepted for the sake of study that MRI is known for analysis of predictive AD diagnosis, as MRI imaging technique is known for its spatial resolution[1] which is high, and its accuracy in recognizing regional disparity as opposed by other techniques like Computed Tomography(CT), PET Scans or CSF Scans.

Another special feature that can be accessed mainly through MRI Scans is the ability to decide a particular Region of Interest (ROI)[2][24][7][26] like the earlier discussed Hippocampal region of the brain. This can give a more specialized approach towards studying the principle cause of a Neuro pathology in our case AD, as opposed to having to look at the changes at the entire brain structure as a whole, and aside from simply judging spatial or volumetric analysis of a brain sample, through ROI, and other methods like Voxel Based Morphometry, it is also possible to look at the concentration levels like Grey matter concentration and so on.

Aside from the medical imaging and brain scans, that clearly provide a particular view towards identifying abnormalities in the brain structure, there are also other factors that need to be considered some of which include cognitive patterns, a person's behavioral patterns, medical history and so on. Some neurological tests like the Mini Mental State Examinations(MMSE), have proven to be very beneficial in providing with some cognitive information that can help aid with the diagnosis of the Alzheimer's Disease. The main disadvantages of making use of the MMSE tests, are, they do not accurately provide information on the severity of the Alzheimer's disease. However information obtained from a neuroimaging e.g. MRI which can be integrated together to form an MRI Biomarker can be further integrated with cognitive Biomarkers like MMSE information to form an aggregate biomarker to remove the disadvantages of the MMSE tests.

In the field of Computer Science, a lot of study has been put into the application of Computer Systems in the Analysis of Medical Pathology and Medical Diagnostics. With the role of Artificial Intelligence, booming in the area of medicine, lot of study has been put into having computer systems learn to detect pathology through means of Neuroimaging and other techniques. Some of the techniques used in the detection of pathological evidences, in our specific scenario, the use of systems to detect if a patient is likely to get affected by Alzheimer's disease or not, include Machine Learning Techniques, and Deep Learning techniques, in tandem with proper Computer Vision and Image Processing techniques. Some of the most widely used image processing techniques include use of Region Optimization and Region of Interest. This in tandem with Machine Learning techniques like unsupervised Learning Techniques(PCA, k means clusters), supervised Learning techniques(SVM) along with Deep Learning Techniques(Convolutional Neural Networks) have helped with detecting patterns that can aid in the diagnosis of Alzheimer's disease in patients. A further study into the current techniques that have been used to detect Alzheimer's disease will be looked in detail upon the next section.

## 2.2.  Literature Survey:

As discussed in the previous section, many efforts have been taken into the analysis of Neuro Imaging techniques like as the use of MRI images to help aid with detection of Alzheimer's disease in patients. Generally, the aged brain takes 4-5 years to get completely affected, however it might take 11 years on a younger population, as discussed by[25]. According to [25] machine Learning techniques have been employed in the study on Alzheimer's Disease, some of which include Support Vector Machines and Extreme Learning Techniques. Another form of analysis[26] is through means of using cognitive ensemble trainings, which is integrating 2 or more Machine Learning techniques to generate a more efficient Machine Learning Technique.. [26] has used a Binary Coded Genetic Algorithm(BCGA) along with

Extreme Learning Machine to aid with Classification. The BCGA technique has been used to select prominent features from an VBM Mechanism, the selected features are then added as features into an ELM techniques, which would then result in a proper cognitive ensemble estimation, which is based on correctly classified outputs and misclassified outputs. This cognitive ensemble classifier helps in providing a much higher accuracy than a standard Extreme Learning Machine in terms of output as, there has been a substantial increase in the training and testing accuracy by about 20% difference(from 75% to 95%) by using a cognitive ensemble technique as opposed to using a BCGA or ELM training separately.

From experiments carried out by [2], the most widely used method of Support vector Machines was employed, due to the high dimensionality nature of MRI Image Analysis along with transition predictions. However some of the limitations of SVM itself have been countered due to the inefficiency of selecting s rightful kernel function, so using an SVM methodology, an accuracy of 51% to 73% was achieved for higher order classification groups. To cope up with such limitations, the Gaussian Discriminant Analysis(GDA) has been employed, which has a basis on a feature selection has claimed that the entorhinal cortex, one of the sub regions of the hippocampus, proves to play an essential role in detection of AD from MCI. Unlike the previous method of BCGA, [1] argues that as there is a greater shift with shape as AD progresses, a more ROI approach is considered unlike the VBM Method. Many performance features have been considered like Accuracy, F1 Score, using 80% of Subject noise free data, on a ten-fold cross validation and the below F1 Scores were obtained while classifying the data into HC vs MCI, AD, and MCI vs AD.

| Groups | CN vs. MCI and AD | | | MCI vs. AD | | |
|---|---|---|---|---|---|---|
| Decision Space | Left | Right | Comb. | Left | Right | Comb. |
| F1 % | 83.33 | 84.70 | 97.20 | 55.91 | 55.56 | 79.82 |
| ACC % | 74.00 | 75.00 | 96.00 | 68.00 | 68.29 | 87.43 |
| SEN % | 92.86 | 98.86 | 99.14 | 64.55 | 58.18 | 79.09 |
| SPE % | 30.00 | 19.33 | 88.67 | 69.58 | 72.92 | 91.25 |
| PPV % | 75.58 | 74.09 | 95.33 | 49.31 | 49.61 | 80.56 |
| NPV % | 64.29 | 87.88 | 97.79 | 81.07 | 79.19 | 90.50 |
| Number of the optimal features | 50 | 61 | L: 67 R: 1 | 6 | 1 | L: 27 R: 1 |

*Table 1: Performance Summary of a 10 fold cross validation after classifyingCN vs MCI and AD*

From [6], it can be seen that there are also a number of Deep Learning techniques that have been used to aid in the diagnosis of Alzheimer's Disease. [7] made use of SVM classifier, a technique similar to the paper published by[3], however, upon achieving a low accuracy rate of 84.4%, consorted to applying the convolutional Neural Networks technique CNN on various different types of image segmentations and chose the one with the highest accuracy rate of about 96%. Also a Bayes Classification was considered which provided a higher accuracy between CN and AD but a lower accuracy between MCI and non MCI. As far as CNN is considered, the use of SVM, focuses on high dimensionality and a more generalized effort whereas, the reverse is true in the case of CNN, where the image is repeatedly convoluted and pooled, so that the exact specific region of abnormality is captured. Hence application of CNN in this case proved far more effective than use of an SVM.

In their paper, [7] like [3] has made use of the concept of ensemble training. Here, they have adopted the use of a hybrid neural architecture like a Counter Propagation Network(CPN) to handle input space quantification. This along with ensemble trainings like s combination of machine learning techniques have provided a much higher accuracy of 77.5% as opposed to the use of standard CPN's. this method along with the proposed methods in [26] give a valid emphasis on the application of ensemble training techniques to provide better accuracy to detect the occurrence of Alzheimer's disease as opposed to the single ML technique strategy.

The method shown by [9] depends mostly on the risk factors associated with Alzheimer's Disease. In this paper, emphasis has been shown as to why it is not merely enough to study MRI, PET or CSF Scans, and why it is also important to have a general idea towards looking at neuro psychological factors, like family history, demographics and many other such important factors. In this paper, Machine Learning Techniques have been provided to rank the importance of some common risk factors, and how they can create a strong impact towards affecting the aging flow of the brain. Through this paper[8], it is shown that demographics as well as lifestyle play a very quintessential role in determining if a patient is likely to fall under the bracket of Alzheimer's Disease. Like as in earlier papers, there has been proven evidence to determine the important risk factors by training the network using ensemble training methods as opposed to standardized training methods.

The article by [23] focuses on a different approach of studying Computer Aided Diagnostics. In this approach, key slices have been considered, for testing purposes, upon which considerable eigen brain maps have been created for each brain subject. Furthermore, the most important map is considered, and a t test is carried out to consider feature selection. Once the relevant features have been gathered, standard machine Learning methods have been considered to perform classification. This paper necessitates the importance of having to use the concept of region of interest as opposed to other general methods including that of voxel based morphometry. Here again, a polynomial kernel was proved to provide better results to that of a monomial or a single training technique kernel.

As discussed earlier, much emphasis has been focused around the relation between the Hippocampal region of the brain to that of being the part or subregion that gets affected during the early onset of Alzheimer's Disease. In their article [7], Hippocampal Shape, volume and texture. Although this looks easy enough, the concept is quite complex in the exact study of the right point of region of interest. Although hippocampal region is able to

identify, the early onset of Alzheimer's through the presence of Neurofibrillary tangles, the cortical region plays, just as much as an important role in the formation of Neurofibrillary tangles. Hence a Hippocampal region has been used to aid in the detection of a normal aging(CN) brain to that of the Mild Cognitive Impairment(MCI) stage, whereas, the cortical thickness can be used at a later stage to gather if a brain in the MCI Stage is likely to be affected by AD. Some of the major sub regions including the hippocampal region, the amygdala and the cortical regions is shon in the figure below.



*Figure 4: MRI Biomarker showcasing specific regions in the Brain that could cause AD*

The hippocampal region[7] is a part of the Mesial Temporal Lobe, which has been attributed to show the first signs of formation of neurofibrillary tangles in the brain. The hippocampus is divided into 3 major regions, the CA regions(cornu ammonis regions), the dentate gyrus and the entorhinal cortex. Variation in the volumetric features of these subregions has been proved to be a beneficial factor in the identification of Alzheimer's. Through the experiments conducted[12], apart from CA delineation or considering the most important features targeting the disease through the CA region, other extracted features included surface mesh

modeling or identifying the volumetric variances in the Hippocampal region and Surface based measures.

Aside from the standard MRI Neuroimaging, studies have also been considered on other neuroimaging techniques such as that of FDG PET Scans and CSF Scans. FDG PET Scans have provided a good number of features regarding AD glucose uptake, and cerebral metabolic rates, which are affected even before the onset of AD. Factors like metabolism are very crucial in the identification of rate of mental activity in the brain. Hence parameters like this along with features like age, number of years of formal education and so on, have to be taken into account, which may prove beneficial to identify the root causes as to detecting the onset of the disease. Below figure shows how much of the brain and the subregions are affected due to the abnormality in the glucose uptake in the subregions of the brain. Here again as like [11], we can see the use of ROI's to provide better results in the detection of the disease as opposed to more crude methods like Voxel Based Morphometry(VBM). Further more, making use of methods like 5 fold cross validation gives us a more stable Learning curve.



*Figure 5: Figure showing abnormally low glucose uptake in subregions of the brain*

## 2.3. Critical Discussions:

The aim of the project is to have a computer system automatically detect the presence of Alzheimer's Disease from a patient's routine checkups, both through relational data means including factors like Mini Mental State Examinations, age related factors, gender, and so on, and through Medical Imaging methods like MRI Scanning and so on. In doing so, many experts have looked into different ideas to carry out proper analysis for early detection of Alzheimer's Disease. Some of which are elaborated below.

One major point that has been noted by many experts include the use of Ensemble Training methods[7], [3], or rather a mixture of Machine Learning techniques that have been proven to provide a better accuracy rate in the classification of a subject into an AD likely progressive or normal aging brain as and when compared with using a single ML Technique. One of the major techniques that has been used as proposed [20] include Making use of Support Vector Machines which is a more widely used technique as compared to other techniques. In their experiments [20], an accuracy of 95.03% was achieved, through using a cognitive ensemble classified as opposed to 78.4% that was achieved by using the best Extreme Learning technique. In the articles by [26], a proposed ensemble 3D Convolutional method was considered as compared to other standard techniques[ppr3, 26,27,28,29], an accuracy increase in the separation of MCI to AD transition of 97.1% was observed, as opposed to 91% observed by the standard methods. These proven examples along with many other examples gives us a strong standing statement on the role of ensemble training techniques being better effective in the role of early AD Detection.

The 2$^{nd}$ factor that can be considered is how to value the data. From articles [1], it can be considered a better viable option to consider making use of Neuroimaging information as opposed to simple neuro psychological or cognitive information. The major reason lies in the inaccurate data that has been gathered solely through means of neuropsychological or cognitive information. Furthermore, research and experiments[1] have been take into consideration to include both Neuroimaging gathered features along with age, lifestyle, neuropsychological and cognitive features, as both areas play an important role in extracting salient features While through a MRI image, it is possible to locate, volumetric changes or changes in tissue density, or tissue morphology, that can determine an early onset of the disease, it is still crucial to gather other information such as the reports from a Mini Mental State Examination(MMSE) or the person's age related factors. Other profound studies[22] that can further be adopted in future is inclusion of relevant information from scans like the FDG PET(Fluorodeoxyglucose Positron Emission Tomography) such as the concentration of metabolic activity around the most affected regions due to an abnormality in the glucose uptake in the brain.

A third major factor to look into is the type of Neuroimaging techniques[5] that has been considered for analysis. As discussed earlier, it is not merely the use of Magnetic Resource Imaging(MRI) or Functional Magnetic Resource Imaging(fMRI) that can provide the best and most salient features, but also neuroimaging techniques like the FDG PET and Cerebral Spinal Fluid(CSF) images, that can provide more important features. The major use of fMRI images is the lack of further human intervention involved, or a lack of efforts needed by a patient, and the reason for the high Spatial resolution and spatial information that is provided by fMRI images, that makes fMRI imaging techniques more widely used. Furthermore, the primary characteristic that can be identified by the Alzheimer's Disease is the shrinking in the brain volume and the abnormality in the hippocampal regions of the brain which makes it

more important to collect fMRI images into a standard space, and warp it to proper positions and perform volumetric analysis to detect abnormality in spatial characteristics. Alzheimer's Disease is formed by the creation of intracellular neurofibrillary tangles and extracellular β-amyloid plaques. A detailed description of the concentration of these tangles and plaques can be studied using CSF Scans, however MRI is still the more widely form of use in the Neuroimaging industry to aid with the detection of Alzheimer's Disease.

The 4[th] important thing to be considered is in tandem with the 1[st] point, as to which Machine Learning Technique would be most appropriate to be considered for analysis of both the neuroimaging stage as well as the neuro cognitive information. An MRI Biomarker or the output that will be extracted from an MRI Image can no doubt be in the form of 3D Voxel matrices, in which case it would be more viable to use a deep learning classification technique like the Convolutional Neural Networks or CNN. Information extracted from the entire brain structure, however, such a method would serve to be ill opposed when compared to other much better suited methods such as ROI analysis, wherein, features are gathered from specific regions of interest. In such a case, methods like Unsupervised Learning, or Semi supervised Learning[23] would prove to be more effective to classify a patient's brain into sMCI, and pMCI respectively.

# 3. Methodology

## 3.1. Experimental Overview

The model which we wish to accomplish is to have a system learn to detect whether a patient is likely to be affected by Alzheimer's Disease or not. To present this model, we have looked into several plus points and advantages, most of which were discussed in the previous section. However, taken all factors discussed in the previous discussions. The model we intend to build will work in a 2 Stage approach. As we looked before, it is not merely enough for a neuroimaging technique to help identify the presence of the Alzheimer's Disease, taking this into account, our model wishes to take MRI parameters, and perform an unsupervised Learning algorithm develop an MRI Biomarker, that will play a pivotal role in the $2^{nd}$ Stage.

In the $2^{nd}$ Stage, we will use this MRI Biomarker to associate with some of the more Neuro psychological or Neuro Cognitive Methods, to develop an Aggregate Biomarker.

As far as previous discussions were considered, we will be employing the use of an MRI Bio Marker, as opposed to other biomarkers like a CSF Marker or PET Marker. Also as we saw a plus point that by developing a more efficient algorithm by making use of Ensemble training techniques. Our system will employ an ensemble training technique in the $2^{nd}$ stage. In this stage, the output generated from the $1^{st}$ stage that is the MRI Biomarker along with cognitive information, involving a person's age, sex, lifestyle and ethnicity, and their results in cognitive examinations like Mini Mental State Examination will be subjected to an ensemble training technique, which will altogether give the desired output as to whether the patient is likely to fall under a normal aging process(Controlled Normal) or suffer from Alzheimer's Disease.

However, before, we plan on carrying out any machine Learning techniques, we have also placed a lot of emphasis on how we intend on getting all the information from MRI Scans

into their relevant form. To this, we have to perform the necessary Image preprocessing techniques, along with a proper region of interest. From previous study. It has been noted, that the Hippocampal region plays a pivotal role in the early detection of Alzheimer's disease, hence to further the efficiency of our experiments, we will consider the Hippocampal regions as our regions of interest.

Hence when we work towards achieving our MRI Biomarker, features attributed to the Hippocampus as opposed to the entire brain structure. Also the number of features that can be acquired from the Hippocampal region are several in number so we will carry out experiments to decide the relevant features, most commonly tissue quantity and Hippocampal volume to help assess the detection of the disease.

An extended overview as to the techniques we have assessed to aid us in selecting the right model, along with the ways in which we have tended to approximate our selected features will be explained in further in the following sections Also an overview of the tools, which we intend to use to preprocess our information as well as select our ROI will be listed in the below sections.

## 3.2. Project Planning and Methodology

### 3.2.1. Work Breakdown Structure:



*Figure 6: Work Breakdown Structure of Our Project*

### 3.2.2. Gannt Chart:

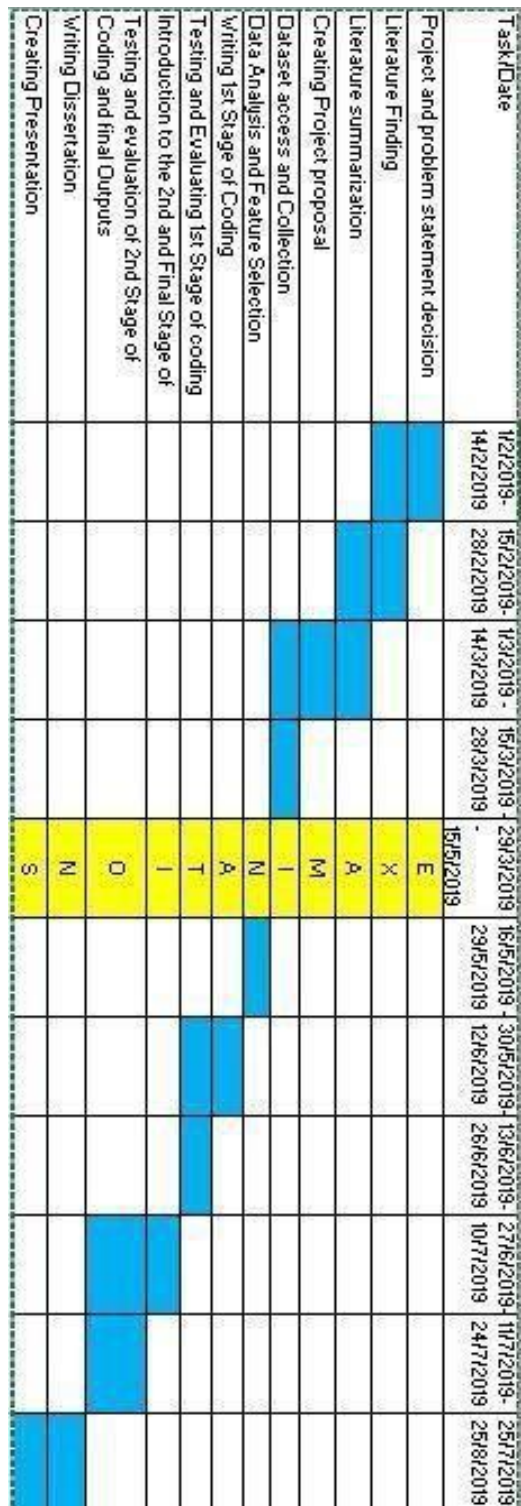| Task/Date | 1/2/2019 - 14/2/2019 | 15/2/2019 - 28/2/2019 | 1/3/2019 - 14/3/2019 | 15/3/2019 - 28/3/2019 | 29/3/2019 - 15/5/2019 | 16/5/2019 - 29/5/2019 | 30/5/2019 - 12/6/2019 | 13/6/2019 - 26/6/2019 | 27/6/2019 - 10/7/2019 | 11/7/2019 - 24/7/2019 | 25/7/2019 - 25/8/2019 |
|---|---|---|---|---|---|---|---|---|---|---|---|
| Project and problem statement decision | ■ | | | | E | | | | | | |
| Literature Finding | ■ | ■ | | | X | | | | | | |
| Literature summarization | | ■ | ■ | | A | | | | | | |
| Creating Project proposal | | | ■ | | M | | | | | | |
| Dataset access and Collection | | | ■ | ■ | I | | | | | | |
| Data Analysis and Feature Selection | | | | | N | ■ | ■ | | | | |
| Writing 1st Stage of Coding | | | | | A | | ■ | | | | |
| Testing and Evaluating 1st Stage of Coding | | | | | T | | | ■ | | | |
| Introduction to the 2nd and Final Stage of coding | | | | | I | | | | ■ | | |
| Testing and evaluation of 2nd Stage of Coding and final Outputs | | | | | O | | | | ■ | ■ | |
| Writing Dissertation | | | | | N | | | | | | ■ |
| Creating Presentation | | | | | S | | | | | | ■ |

*Figure 7: Gannt Chart*

### 3.2.3. Proposed Methodology:

In this section, we will look into the proposed software engineering methodologies thst are generally carried out on research projects in the IT Industry. We will further look into each of the models and state the advantages and disadvantages as to why we used our proposed model. In the end, the waterfall model has been selected as the most apt model to handle the proposed model.

**Software Engineering Methodologies:**

When any research project is to be implemented, it always bears importance upon how the initial planning, methodology, and metrics are going to be visualized. "Trying to manage a project without project management is like trying to play a football game without a game plan". Hence when we work out a model to help detect the early onset of Alzheimer's Disease in patients, a lot many methodologies must be considered to select the perfect methodology. Some of the models that we have considered include the Waterfall Model, the evolutionary prototyping model, the V-Shaped Model, the incremental Model, the Spiral model and the Agile Model.

As we delved deeper into analyzing each model, we ruled out certain features, which in the end helped decide the proposed model.

**The V-Model**: Initially, we looked at one of the most well sought after models, the V MODEL. Like the Waterfall model, also called the verification and validation model. The main feature of this model is that it runs through means of phases, wherein each phase needs to be fulfilled before proceeding to a different phase.

When we considered our model. We looked into a more straightforward, training and testing model, wherein we learn from training data and we implement the model on testing

data.Validation in our case was considered after checking the accuracy of testing data from using a proposed model and hence such a method, requires a sequential model rather than a parallel model like the V Shaped model, and hence this model was rejected.

Also the fact that this model looks into coding at the very end, and testing and validation at the early stages, this did not suit our propose model, wherein, evaluation was only done after coding. Hence we ruled out the use of the V Model.

**Incremental Model:** The incremental model is quite similar to the Waterfall model. Also like the V model, this model works in the form of phases. In this model, each phase cycle carries out a requirement cycle, and at each new phase a later requirement is considered. As in the case of our model, we just have a single requirement, that is having the system work on different modes of Machine Learning models and learn to identify the progression of AD in an MRI sample along with more characteristic cognitive features. A sequential model was the major requirement and hence even this model was rejected.

**Spiral Model:** In the planning stage, major emphasis is provided on the risk analysis and risk planning. Despite this model being very effective in the calculation of risks and this model helps dealing with large scale projects. Risk was not a major emphasis into our project, hence even this model was not considered.

**Agile Model:** In today's era of the information age, a lot of emphasis has been catered on the rise of technological advancements, it is conidered normal, if you have your client approach, you, saying he needs to have a variation, in his original set of requirements. Furthermore, such is the complexity of today's projects, that a normal project manager, would insist on considering making user of a hierarchical model, rather than use simply a single model like a V-Shaped Model. Hence it is of utmost importance that at the very earliest Project planning stage, we should have a model that has a high confidence level to base with the constant

ever-demanding nature of requirements, code changes, code optimization, design and model evolutions and model evaluations, test upgrades and changes in deployment specifications.

Most of the traditional software engineering projects that are undertaken in organizations make use of a sequential step wise manner of procedure, commonly called the waterfall model but what differentiates the waterfall model from some of the more advanced methodologies is its lack in adapting to the ever-demanding changes. Hence an approach, we initially focused on in fulfilling for our project was the Agile Methodology.

An agile model can be thought of as a group of waterfall models. As discussed before, agile models are a fix to the ever demanding changes and additions to requirements. Hence when we look at the agile model, it is like a group of waterfall models which work in sync with each other, each handling the old and new set of requirements. The agile method in general, is not structurally dependent like the waterfall model. There is instead more emphasis provided in trying to communicate the exact requirements that have been received from stake holders, priority and impact of the new requirements on the original design and so on. Hence, the Agile method works more on a more communicative and real time basis to help cope with the change in demands of requirements.

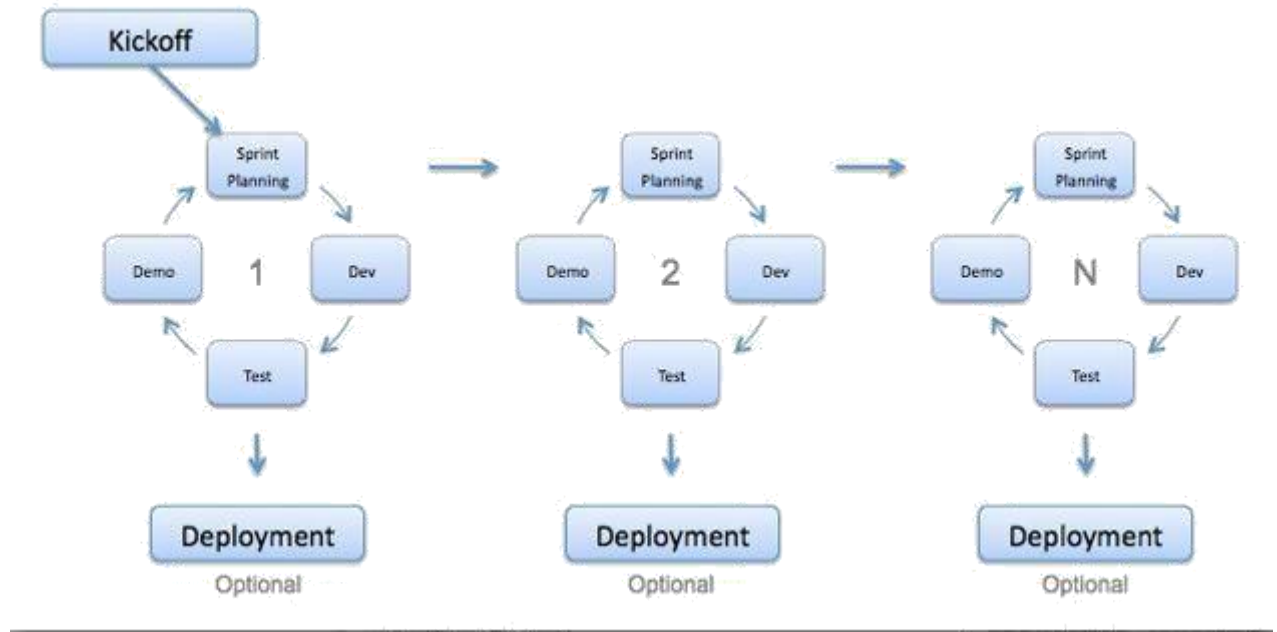An example of the agile method process can be seen in the below figure:

*Figure 8: An example showing the general Agile Model Process*

In the agile model, each cycle is called a sprint. A Sprint master manages the whole Sprint structure. A project owner allots tasks to project engineers, The project engineers carry out their individual tasks and peer review all completed tasks to get the required feedback.

In our project, emphasis is provided into trying to achieve a proper ML model that can help aid in the detection of Alzheimer 's disease through a person's MRI Scans and neuropsychological information. In the world of Machine Learning, the object in consideration are not associated with single tasks. Our original model and requirements would have been to try and detect the Alzheimer's disease through means of MRI Scans. And later on the concept of aiding the process of detection through means of Neuropsychological factors meant, the requirements were changed. Furthermore, the model could also come up with deciding further features, or implementing scans to our model to ensure detection of Alzheimer's disease, and hence this project though initially planned, was considered for an agile approach.

However as discussed before, a lot of approached even in the field of Machine Learning can be adopted, or enhanced to a model to help aid with the detection of the disease. For instance, we could, instead of adopting the VBM approach of characterizing the model, through the volumetric and intensity information of voxels, employ to try using other brain functions associated with fMRI. In practicality, a lot of research has yet been need to consider, for helping with neuroimaging pathological diagnostics, and hence, as we decide to create a simplistic, and single requirement model, we decided to abandon the approach towards a agile model.

**Waterfall Model:** One of the most common Software Engineering methodology to handle simple projects with lower number of requirements is the waterfall model. As stated before in a waterfall model, the requirements are fixed and known at the $1^{st}$ stage of the project. Hence as the requirements are specified just in the $1^{st}$ stage of the model, it follows a more contractual approach to project planning. In the waterfall model, each stage is important. Right from requirements gathering, planning, designing, implementing, testing and deploying stages, every stage is carried out in a sequential fashion. In other words, the waterfall model is called as the linear sequential lifecycle model.

An example of the waterfall model is shown below.
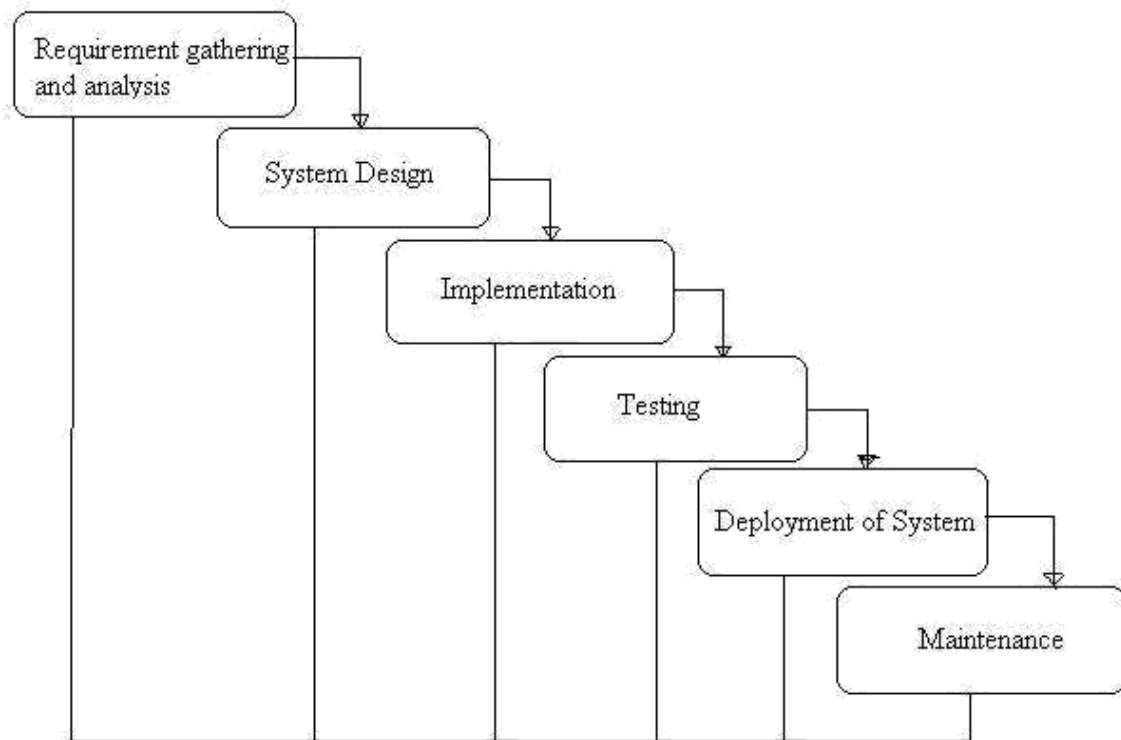
**General Overview of "Waterfall Model"**



*Figure 9: An example of a Waterfall approach*

From figure 8, it can be noticed that the Waterfall model is divided into several phases or stages:

In the requirements and information gathering stage, much of the requirements and elicitation properties are gathered. In this stage, the primary stakeholders, have meetings and sessions with the other members of the project like the project manager, team lead and so on, on the requirements that are need to be performed as an outcome to the project. As this stage directly impacts the project outcome, it can be noted, that this is one of the most important stages in the Software Development Lifecycle. In the design analysis stage, more emphasis will be put on the type of model to be created, like as in the Use cases, the dataflow diagrams and in ourcase, the work flow diagrams for our experimental models as shown in section 3.1.

In the implementation stage, the codes are set into use, as our modelscreated in the design stage are realized, and the required output is generated Further mode in the implementation stage, some basic testing of the code, like the unit testing parameters are also considered. In the testing phase, detailed testing techniques are taken into account, test cases or evaluation reports are built and a test plan is created.

In the deployment stage, the project is deployed to the stakeholders and future maintenance, handling of project bugs, and fixes are put into account.

We designed our system with the waterfall model, because of its simplistic approach in use, and having our requirements exactly specified at the earliest stage made it plausible to choose this model.

## 3.3. Experimental Setup

In this section, major focus on the Software tools that were used for the implementation of our project will be discussed along with how these tools were used in the implementation of our model.

### 3.3.1. Datasets Used

One of the leading Initiatives, that has solely been created to conduct study exclusively in the area of Alzheimer's disease is the ADNI or Alzheimer's Disease Neuro Imaging Initiative (Alzheimer's Disease Neuroimaging Initiative, 2019). The Alzheimer's Disease Neuroimaging Initiative (ADNI) brings together, Neuroscientists and researchers with data, as a collaborative effort to help aid in progressive research in the detecti on and diagnostics of the Alzheimer's Disease. ADNI researchers collect and use data from Magnetic Resonance Imaging and Positron Emission Tomographical images, cognitive tests, CSF and genetics and blood biomarkers to predict the onset of the disease. Study resources from North American ADNI study along Alzheimer's disease patients, MCI subjects, Neuropsychological and Neurocognitive study are provided through the ADNI Website, which we have used for our study.

For our study, we have made use of 2 levels of information, one is the        Study data, which will provide more emphasis on the neuro cognitive information which will be required in

the $2^{nd}$ phase of our 2 phase project, and we have also taken into account the fMRI images , which we will be using for our $1^{st}$ phase of analysis. For analysis, we have made use of data recorded from 150 subjects, each having had multiple MRI Screening sessions over a period of 2 years from Initial screening. The kind of MRI Images that have been considered inclue the 1.5T type images, altogether a little close to around 700 images will be analyzed in general. Each of these subjects have been classified under 3 target stages of the disease, the Controlled Normal(CN), Alzheimer's Disease(AD) and Late Mild Cognitive Impairment(lMCI) or rather Progressive Mild Cognitive Impairment(pMCI). Anumber of MRI Scans have been provided per subject over the 2 year period, showing various changes that occur in the brain structure which might lead to the potential onset of the disease.

Apart from the MRI Image dataset, which we have used, we have also considered a dataset involving the same subjects, which have been carried out be a previous challenge. This dataset provides information regarding, the subject's gender, initial age, Mini Mental S tate Ecamination Scores and so on. Using the MRI Scansfrom the $1^{st}$ phase, we will generate an MRI Biomarker, which will use the dataset from the challenge dataset for the $2^{nd}$ phase, to create the aggregate biomarker.

### 3.3.2. Hardware requirements

When it comes to hardware requirements, we will have to take into account that much of what the system wishes to account for is handled through Image processing in the form of MRI Processing and in regards to Relational Database information. Memory requirements by our system should be able to take into account the vast amount of data collected through the form of 3D Nifti images, along with further information from preprocessing the image into our system memory. In order to avoid the memory overload. A stringent method of gathering information from an external source like an

external hard disk holding much of the relevant data information was adopted to fulfil memory requirements.

Apart from memory, the next important concern to any Machine Learning Project is the efficient use of processing power. In terms of Machine Learning, especially with copious amount of data, a lot of processing power is required, and hence, making use of a higher end GPU, like as in the GPU provided by Google Collab was used to monitor Machine Learning.

Apart from these 2 connections, a good Internet connection to acquire data, and a good file management system to store and retrieve information was needed for the efficiency of building and executing our model.

### 3.3.3. Software Requirements:

To handle the processing of Images, as well as data assembly and Machine Learning, primarily, 2 forms of programming languages were used. **Matlab** and **Python 3.6.** There are many special softwares that are employed by Matlab to help aid with the Image preprocessing of fMRI data, some of which include **Statistical Parametric Mapping(SPM)**, Analysis of **Functional Neuro-Images(AFNI)**, **Free Surfer**, and the **FMRIB Software Library(FSL)**. Of all these tools, SPM works directly in tandem with that of the Matlab Software, and also works independently in different platforms, as opposed to FSL, Free Surfer and AFNI, which work solely on Linux, CentOS and MacOS Platforms. The Operating System, which we corresponded to use for our research was **Windows 7(64 Bit Version)**. So seeing as Statistical Parametric Mapping was the best tool, that we could employ for our research. We made use of SPM.

"Statistical Parametric Mapping(SPM) refers to the construction and assessment of spatially extended statistical processes used to test Hypothesis about functional imaging data". SPM, currently on its $12^{th}$ version, can be downloaded free of cost to any local machine and upon

setting the unzipped SPM folder as path is Matlab, allows, SPM software to be loadedinto matlab, simply by typing spm in the command window. Figure 10 shown below gives an overview of SPM, home screen.
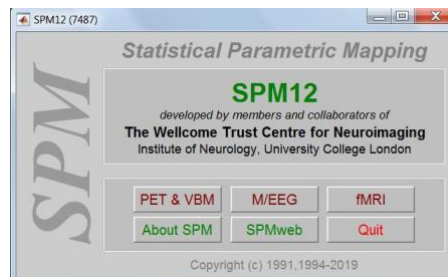


*Figure 10: An overview of the SPM 12 Home Screen*

As can be see, the home screen contains options to select different types of Neuroimaging techniques like Positron Emission Tomography(PET) and Voxel Based Morphometry(VBM), M/EEG(Electroencephalogram) and so on. As we require to work on fMRI images, we select the fMRI button. As can be seen below, upon selecting the fMRI button a new page with 3 windows shows up. The top left screen allows us to choose from different Image preprocessing techniques that will be carried out on the fMRI images like Time Slicing, Realignment, Coregistration, Normalization and Smoothing, along with co comparing with images, using Check Reg, First and Second Level Analysis and the Display buttons. The bottom Left window, is used to show transition periods, or rather upon running a certain st of Analysis, a thermometer like Loading window popups in this window, giving account to how far, a certain number of images have been analyzed. The window to the right holds the SPM Logo, and in times, of Displaying a Realignment graph, or Display an image. This screen, shows the graph and the images.

*Figure 11: The SPM fMRI Home Screen with its 3 windows*

fMRI images can come in one of 2 formats, the DICOM(Digital Imaging and Communications in Medicine) and NII(Nifti) Images. For themost part of the analysis, the NII Images play a big role, and hence from the top left screen in Figure 11, there is also an option to convert DICOM(.dcm) images to Nifti((.nii) images.

For our research, we handled image preprocessing on Nifti files, using the SPM Tool. For our images, we first realigned the images together to align all images to the mean image of all images. To do so, we selected the drop down button next to Realign and selected the Realign(est and Res.) combo list which prompted to opening a new window as shown in Figure 12.

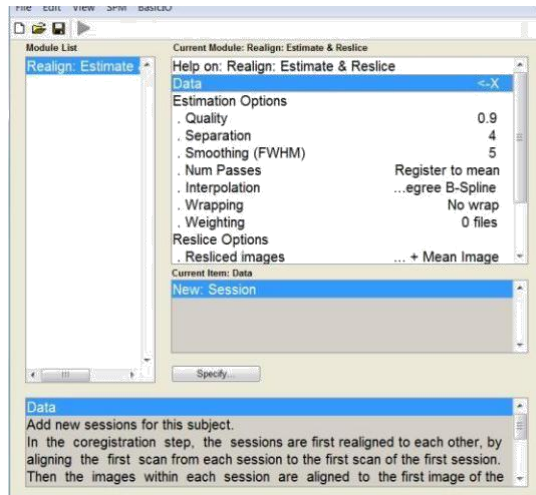*Figure 12: Image Realignment Window*

Here, we created a new session on the data field and selected the images which we needed to realign, and selecting other parameters like Interpolation using a $5^{th}$ Degree B Spline, Wrapping along the Y axis, and so on, we realigned all the images to the mean image.

Upon Reimaging, we set the new Origin for all the images. This is important, as for us, to later on co-register and normalize the images into a 3D Space. The origin was set at the center of the brain along the Anterior Commissure region. Setting the origin upon the images, the next step was to co register all images along the mean image. The co registered images, not only points out realignment but also maps, the interior sub regions of the brain to the closes point possible to build the best sort of analysis of a brain structure. Figure 13 shows us the window for co registered estimation and resliced using SPM.
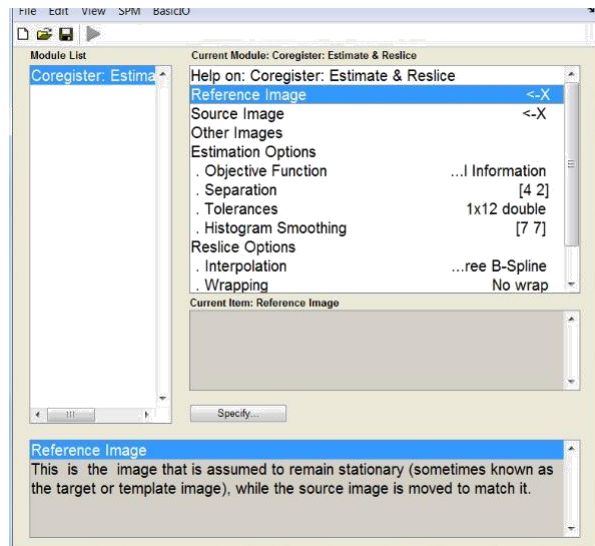
*Figure 13: Co registration Window in SPM*

Once again like in realignment, we enter parameters pertaining proper co registration to occur in the image, like selecting the interpolation, reference and source images.

Now that we had the images properly aligned in a co registered space. We had set the images in a standard template 3D Space. The most common space that we considered was the MNI(Montreal Neurological Institute) Space. Hence the next step, we considered was normalizing the model into a 3D MNI Space. Once again, upon clicking the Normalization tab in the top left window of the SPM fMRI Home page, we entered the required parameters and we were able to generate a normalized image of the brain structures in MNI 3D Space. The last step which was considered for our preprocessing step using SPM was smoothing. Smoothing is just the unwarping of the brain mage along a Gaussian Kernel of a Specified width. Like as in the previous preprocessing techniques, clicking on the Smooth Tab, a new Smooth window was opened, wherein we set the required parameters, and had the image smoothened into a required space.

Upon preprocessing, we transformed the image from an individual MRI, into a realigned, co registered normalized and smoothened image. The image at each of the preprocessing stages can be shown in the Figure below.



*Figure 14: Original Image, Realigned, Co registered, Normalized and Smoothened Image.*

*Once* the Image was properly smoothened, we needed to acquire a region of interest, instead of having to study the entire brain structure, hence, we needed to select the ROI in our case, the Hippocampal Regions. As the MRI Scans, does not have the required Labelling concept. We make use of a tool **MRIcron**, which is downloaded with a whole range of templates, than can be adapted to the MNI coordinates to find specific regions of interest. Under the MRICron Software, the AAL Template provides a substantial segementation of the brain into its corresponding Subregions.In the tool, upon opening view, we can set the coordinates into MNI Coordinates to reflect upon the preprocessing space, which we looked into earlier. Once

that is done, the corresponding label to a particular point the MNI axis can be found in the top left corner of the application, as shown in figure 15.



*Figure 15: MRICron Software showing region of Hippocampal volume.*

Using the coordinates that we have gathered. The ROI can be calculated from all images using a tool like **MarsBar**, which again works as Add In to Matlab's SPM. Upon installing the marsbar package into the SPM Toolbox, the Marsbar software can be accessed through SPM directly. Upon Selecting, the Marsbar too from SPM, a new window opens, allowing us to look into ROI definitions as shown in Figure 16.



*Figure 16: Marsbar Home Window*

Under the ROI Definitions, we can build ROI's based on a cluster of voxels, which we can generate by selecting the appropriate contrast regions, or we, could create a spherical region

as our ROI, or a cuboidal region as our ROI. As in our scenario, we are looking into presence of Alzheimer's disease in a brain sample which has been characterized with a change in Hippocampal volumes, it was not appropriate to make use of clusters, nor spherical region, rather we decided to go with the cuboidal region of Interest. Applying the coordinates generated through the MRIcron Software, we were able to create the necessary cuboidal regions encompassing the left and right Hippocampal regions of the brain, as shown in Figure 17.



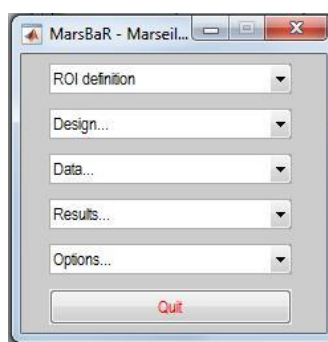*Figure 17: ROI Region Selection Left and Right Hippocampal Regions*

Considering, we have now obtained the information of ROI, using SPM, our next stp involves, making use of Python, and Python Libraries, to calculate the intensity values and Hippocampal volumes for all of our brain samples along the selected Region Of Interest.

The python version, which we have partaken in our application is **Python 3.6**, we have made use of **Jupyter Notebook** to run our code. Using the packages in Python like Nibabel, we are able to read fMRI Nifti images, in a 3D Matrix along with its intensity values. Making use of packages in python like Scikit Learn and Numpy, we have applied Unsupervised and

ensemble Machine Learning algorithm into our system to aid in the detection of Alzheimer's Disease.

## 3.4. Experiments Carried Out.

As discussed in section 3.1, our model consists of a 2 phase system. In the first phase image information would be gathered and a Dimensionality Reduction Algorithm would be carried out to provide some relational image information. Whereas in the $2^{nd}$ phase, the output or the relational image information gathered from the first phase would be attributed with cognitive information from patients through the life style and neuropsychological tests. Then machine learning will be carried out by the system on the entire dataset. And through pattern recognition, it would be considered if a patient is likely to ail from Alzheimer's Disease or not in the long run.A working model of our 2 phase system is shown below(Figure 19).
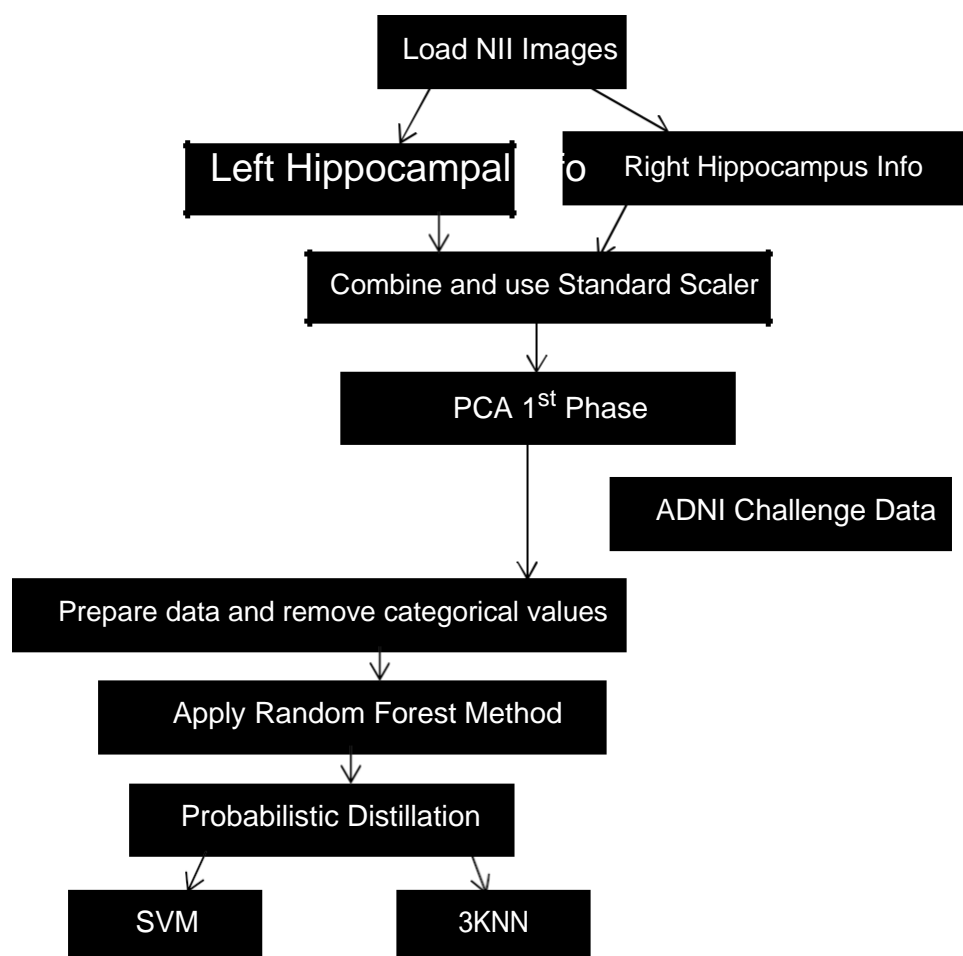


*Figure 18: Working Model of 2 phase model*

For our experiments, 150 subjects were selected through the ADNI Initiative. For each of the subjects, an MRI initial screening session was conducted, and over the span of 2 years, numerous MRI Sessions were conducted, to gauge the changes in the aging pattern of the Human brain. Apart from this, other personal information and lifestyle preferences of the patients along with the MMSE Cognitive examination scores were conducted.

Much of the emphasis of the Image processing experiments are software GUI based, and hence a major account on the same is in Section 3.3.3. However much of this section will focus on the Machine Learning techniques used to generate our result.

### 3.4.1. Implementation Of First Phase Of Model:

**Principal Component Analysis:**

For the phase 1, of our project, the average intensities along the ROI Regions of Hippocampus are gathered along with the voxel count for the Hippocampal volumes. As a major drawback of machine Learning includes the processing time, which is generally huge in the gathering of Image data, one of the major techniques which is used to aid in machine learning is the use of principal component analysis. Principal component Analysis is a Statistical procedure that aids in transforming linear correlated components in a database into uncorrelated principal components. In this model, each principal component is arranged in a fashion such that the $1^{st}$ component has the highest level of variance, and the components from 2nd have the highest variance pertaining to a constraint raised y previous components.

As our model accounts for intensities observed from a ROI region form the hippocampus, it is necessary to reduce the information provided by the Hippocampus, into respective Eigen vectors or principal components that can aid in future processing and Machine Learning of the image information.

To carry out PCA Analysis, the image information in consideration must be normalized. This has been carried out using Python sklearn's preprocessing technique Standard Scaler, which normalizes the image information into ranges between 0 to 1. This can be shown in the below figure 20.

```
from sklearn.preprocessing import StandardScaler

features=["Left_ROI_Avg Intensities","Left_Hippocampal_Volume","Right_ROI_Avg_Intensities", "Right_Hip
x=df2.loc[:,features].values
x=StandardScaler().fit_transform(x)
print(x)
```

```
[[ 1.38917487  1.5894731   1.57080637 -0.32827867]
 [ 1.26052281  2.27743259  1.52029483 -0.29000782]
 [ 1.38676093  1.59851605  1.57224538 -0.32948216]
 ...
 [-0.59512657  0.54918654 -0.5858269  -0.43972146]
 [-0.35875949 -0.75334545 -0.30840502 -1.39625204]
 [-0.40458473 -0.4834483  -0.39932544 -1.10284218]]
```

*Figure 19: Code Snippet Showing use of Standard Scaler for Normalization*

Once the images have been normalized, sklearn provides the PCA method under the decomposition class, that when provided with the number of principle components, reduces the number of features into the principle components. A code showing the utilization of PCA Algorithm is shown below.

```
pca=PCA(n_components=2)
PC=pca.fit_transform(x)

    Subject_ID       PC1        PC2
0    002_S_0295   2.370924  -0.693114
1    002_S_0295   2.372282  -1.203138
2    002_S_0295   2.372031  -0.698526
3    002_S_0295   2.407100  -0.976808
4    002_S_0413   1.672813  -0.126494
5    002_S_0413   1.236724   0.096254
6    002_S_0413   1.671244  -0.111068
7    002_S_0413   1.566481   0.003706
8    002_S_0619  12.930872   3.287687
9    002_S_0619  11.014633   2.667719
10   002_S_0619  13.146209   3.300931
11   002_S_0619   3.274378   0.367081
12   002_S_0619  12.479669   3.292650
13   002_S_0685   1.396280  -0.275740
14   002_S_0685   1.947561  -0.682646
```

*Figure 20: Application Of PCA Analysis dimensionality reduction to the Image Information*

Upon application, our MRI Image information like the Left and Right Hippocampal volumes along with the average ROI Intensities are reduced into its principal components to form an MRI Biomarker which can be used for analysis in the Stage 2.

### 3.4.2. Implementation Of $2^{nd}$ and Final Phase Of Model

In the 2nd phase, we combine personal information, lifestyle preferences and neuropsychological tests, which are then applied to Machine Learning Techniques, to determine if a patient is neurologically aging normally, or if he/she is in the late stage of MCI, or if he is likely to be affected by AD.

**Feature Selection**

The data which we derive for the $2^{nd}$ phase is from the ADNI Challenge dataset. The dataset provides a lot of features, however, all we will consider are the cognitive information including MMSE Score results, a person's age, education level, gender, race and ethnicity. Making use of Python's Pandas data frame structure, we will be able to acquire the relevant information. Furthermore, before performing machine Learning techniques, it should be noted that Python's scikit learn package, which supports Machine Learning, does not categorical data, and hence there is need of performing encoding techniques which can be beneficial to convert categorical data into numerical data. To do this, we make use of One Hot encoding, which converts a feature of values into multiple feature value columns as shown below, this feature value column works as a binary list plotting 1 if a particular tuple records that feature value as true else 0.

A label encoder has been used to convert the target values into numerical label numbers.

**Random Forest Technique:**

The first machine learning technique that we will apply on our system is the Random Forest technique. The random forest technique is an ensemble training technique that works like the application of a group of decision tree algorithms. In the case of a random forest algorithm, the outputs generated by the individual decision trees are taken together, and the modal class, or the class with the highest probability of being the right class is considered the predicted value of the output of the Random Forest Method. It helps in reducing the habit of Overfitting of training data with test data, and hence it is decided to implement this level on ensemble Machine Learning technique for our aggregate biomarker. The fundamental use of Random Forest method is as shown below.

```
In [4]: x_train,x_test,y_train,y_test=train_test_split(X,Y,test_size=0.2)
        clf = RandomForestClassifier(n_estimators=100)
        clf=clf.fit(x_train,y_train)

        scores = cross_val_score(clf, x_train, y_train, cv=5)
        print("Accuracy of a Random Forest is:",round(np.mean((scores*100)),2))

        yp=clf.predict(x_test)
        print("Accuracy of a Random Forest predicted over actual is:",round(((accuracy_score(y_test,yp))*100),
        listt=list(le.classes_)
        y_test=list(le.inverse_transform(y_test))
        yp=list(le.inverse_transform(yp))
        #Evaluation
        print('The confusion Matrix is as Below:\n')
        print(confusion_matrix(y_test,yp,labels=listt))
        print('\nThe Evaluation Report is as Below:\n')
        print(classification_report(y_test,yp,target_names=listt))
```

*Figure 21: Code Snippet Showing use of Random Forest*

*Classification* **Probabilistic Binning:**

Employing the fundamentals of what the Random Forest Method was characterized for. Another method which focuses on using the outputs from the Random forest method in addition to our dataset, which would provide a better result in helping aid in providing better accuracy and evaluation. In the method of probabilistic binning, for the first step, we create a probabilistic map with the outputs generated from the random forest method before. A better example of a probabilistic map can be shown below. The probabilistic map, takes in a

43

probabilistic value for each class of the output for a particular subject. This map for all the classes, becomes the additional features in our original dataset.

Once the probabilistic maps have been constructed, we will use the method of data binning. Data binning is a form of data preprocessing technique. This technique helps aid in rounding off values from features to remove minor observational errors. As can be gathered from the output of a probabilistic map, a feature would contain probabilistic values governing chance of detecting a particular class value and hence, data binning has been used to round these probabilistic values to give us single decimal feature values.

Making use of the predict_proba() function, we will be able to create probabilistic mapping. A detailed code on probabilistic mapping and data binning is shown below.

```
In [5]: yproba=clf.predict_proba(X)

        binlab=[0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]

        bins=[0.0,0.1,0.2,0.3,0.4,0.5,0.6,0.7,0.8,0.9,1.0]
        X['AD']=pd.cut(yproba[:,-3],bins,labels=binlab, include_lowest=True)
        X['CN']=pd.cut(yproba[:,-2],bins,labels=binlab, include_lowest=True)
        X['LMCI']=pd.cut(yproba[:,-1],bins,labels=binlab, include_lowest=True)

        print(X)
        print(Y)
```

*Figure 22: Application Of Probabilistic Binning*

**Support Vector Machines in addition to Probabilistic Binning:**

One of the most widely used Machine Learning techniques that has been used earlier to detect the presence of AD is Support Vector Machines. Support vector machine is a Supervised Machine Learning Technique that is used in classification and regression examples. In a support vector machine a clear margin called as a support vector divides the ambiguity in probabilistically determining the class, a particular subject should fall under. Hence making use of this method in addition to probabilistic binning will help boost in the general increase in prediction.

44

A below code snippet includes how support vector mechanism provided by scikit learn helps in the machine learning process.

```
In [6]: x_train2,x_test2,y_train2,y_test2=train_test_split(X,Y,test_size=0.2)

        clf2 = svm.SVC(gamma=0.01, C=100.)
        clf2=clf2.fit(x_train2,y_train2)

        scores = cross_val_score(clf2, x_train2, y_train2, cv=5)
        print("Accuracy of Random Forest + Probabilistic Distillation + SVM  is:",round(np.mean((scores*100)),
        
        yp2=clf2.predict(x_test2)
        print("Accuracy of a Random Forest + Probabilistic Distillation + SVM predicted over actual is:",round
        listt=list(le.classes_)
        y_test2=list(le.inverse_transform(y_test2))
        yp2=list(le.inverse_transform(yp2))
        #Evaluation
        print('The confusion Matrix is as Below:\n')
        print(confusion_matrix(y_test2,yp2,labels=listt))
        print('\nThe Evaluation Report is as Below:\n')
        print(classification_report(y_test2,yp2,target_names=listt))
```

*Figure 23: Random Forest Probabilistic Distillation Method Implementation +*
*SVM* **K Nearest Method in addition to probabilistic binning:**

The k nearest method or the KNN method is a non-parametric method, widely used in classification and regression problems. In this method, a relationship is developed between each of the subject values, and each of the feature classes are classified as based on the closest member relationship amongst the individual features. Although not widely used in most neuroimaging methods. In this case, major emphasis has been placed in comparison of a subject's MRI Scan over a time period of 2 years, and hence it is necessary to judge the relationship between each of the individual MRI screening sessions, to determine the aging process of the brain. Also this in addition to probabilistic binning, can prove to provide better results in the prediction of AD. A simila**r** implementation of figure 23 is used for KNN technique with probabilistic binning.

## 3.5. Cross Validation and Classification Report

A code that generates results is a fine code, but the real question that comes in our mind is how accurate, a particular result is in pertaining to a certain problem. To answer this question is simple. It

is not merely enough to generate the right results, but is also important to evaluate a result under a particular standard. When it comes to machine learning, first and foremost, in the case of a Supervised Llearning Technique, it is important to see how will a Machine Learning technique can be used to learn a particular training dataset. Furthermore, it is also important to know, if the training and test data are biased in terms of the split, and if so, how does this bias affect, the proper prediction of the results. And last and foremost, how effective are the predicted test class output values in accordance to the original test class output values.

In our case, we have an entire dataset, which is then divided into training and test dataset. To help with reducing the bias, we have adopted the use of stratification, wherein, we have divided the dataset into 80% training data and 20% test data. Upon training the data, we were able to acquire an accuracy of training the data on the dataset. Moreover, we used the trained classifier to predict the output classes from our test dataset. These predicted values were then compared with the original output of our test dataset, to provide accuracy in prediction.

In addition to providing accuracy in training and in prediction, we looked into other metrics, to decide the truth value in the classification of output, we acquired this through creating a confusion matrix. A confusion matrix is a table used to judge the performance parameters of a classifier. It compares the mislabeled values as compared to correctly labeled outputs to determine the efficiency of our Machine Learning Techniques. In addition to confusion matrix, we have also generated a classification report, that evaluates other performance metrics like that of preciosn, recall and f1 scores. The below descriptions provide the evaluation results pertaining to the different Machine Learning techniques put into consideration.

**Random Forest Evaluation:**

For our first Machine Learning technique as provided in the previous section, a random forest classification technique was carried out. Below is a table showing the evaluated results observed by the classification report and confusion matrix.

```
Accuracy of a Random Forest is: 97.3
Accuracy of a Random Forest predicted over actual is: 97.86
The confusion Matrix is as Below:

[[18  0  1]
 [ 0 41  2]
 [ 0  0 78]]

The Evaluation Report is as Below:

              precision    recall  f1-score   support

          AD       1.00      0.95      0.97        19
          CN       1.00      0.95      0.98        43
        LMCI       0.96      1.00      0.98        78

   micro avg       0.98      0.98      0.98       140
   macro avg       0.99      0.97      0.98       140
weighted avg       0.98      0.98      0.98       140
```

*Figure 24: Classification Report and Confusion Matrix for Random Forest Method*

Upon having the classifier training our dataset, a training accuracy of 97.3% was achieved. This shows the efficiency of application of this ML Technique on our dataset. Furthermore upon predicting the values for our test dataset against the original dataset, an accuracy of 97.86% was achieved. This proves that the Random Forest method is an effective method in the prediction of Alzheimer 's disease, given our dataset. Also upon employing our training method a weighted average precision, recall and F1-Score of 99% was achieved.

**Support Vector Machine With Probabilistic Distillation:**

As discussed before, the techniques of probabilistic mapping was employed upon the outputs which were received from the Random Forest techniques. This along with data binning techniques was fed into a Support Vector classifier, which generated the below classification reports and confusion matrix.

```
Accuracy of Random Forest + Probabilistic Distillation + SVM  is: 97.12
Accuracy of a Random Forest + Probabilistic Distillation + SVM predicted over actual is: 97.14
The confusion Matrix is as Below:

[[18  0  0]
 [ 0 43  4]
 [ 0  0 75]]

The Evaluation Report is as Below:

              precision    recall  f1-score   support

         AD       1.00      1.00      1.00        18
         CN       1.00      0.91      0.96        47
       LMCI       0.95      1.00      0.97        75

  micro avg       0.97      0.97      0.97       140
  macro avg       0.98      0.97      0.98       140
weighted avg       0.97      0.97      0.97       140
```

*Figure 25: Classification Report and confusion matrix for SVM with probabilistic binning*

As this ensemble of Machine Learning techniques were carried out on the training dataset, an accuracy of using the method on our current dataset provided an accuracy of 97.12% on our training dataset, and an accuracy of prediction over original test dataset of 97.14% was achieve. Also a weighted mean average of precision, recall and F1 Score of 97% was achieved. Ths goes to showing that as compared to a technique like as in Support Vector Machine, a distillated approach of distilling a random forest method with a widely used technique like support vector machine lowers the accuracy only slightly. But this again gives us an insight on the user of an ensemble training technique against a more distillated ensemble with standard machine learning technique.

**K Nearest Neighbor With Probabilistic Distillation:**

As already provided earlier, the K Nearest neig hbor method, does not help to easily determine the patterns in image processing. However, this method is used to classify a relationship between ckuster of subject MRI's over the 2 year period. Below is the classification report and confusion matrix generated upon utilization of this method on our dataset.

```
Accuracy of Random Forest + Probabilistic Distillation + 3KNN  is: 91.92
Accuracy of a Random Forest + Probabilistic Distillation + 3KNN predicted over actual is: 42.14
The confusion Matrix is as Below:

[[ 6  5 15]
 [ 4 12 24]
 [10 23 41]]

The Evaluation Report is as Below:

             precision    recall  f1-score   support

         AD       0.30      0.23      0.26        26
         CN       0.30      0.30      0.30        40
       LMCI       0.51      0.55      0.53        74

  micro avg       0.42      0.42      0.42       140
  macro avg       0.37      0.36      0.36       140
weighted avg       0.41      0.42      0.42       140
```

*Figure 26: Classification Report and Confusion matrix for 3KNN with Probabilistic Binning*

Upon having our dataset tested along with probabilistic distillation on the 3 k Nearest Neighbour methodology. A standard training accuracy of 91.92% was achieved, which although low as compared to the Support Vector Machine and Random Forest technique still gives a substantially good use of this technique to our dataset. From our earlier discussion it can be seen that the accuracy in predicting the classes for the test dataset is considerably low achieving a 42% in detection of the disease. This creates a variance in the structure of weighted average precision, recall and F1 Scores of 42% in the utilization of this technique on our dataset.

# 4. Results and Discussion

The initial aim of our research was to generate a model that would prove to successfully have a system be able to classify if a particular subject would likely fall prey to Alzheimer's disease or if his/her brain, would age in the normal process, ignoring the progress towards the disease altogether. From our various discussions gathered from previous experts in the area, we came into conclusion of many such factors that provide a key role in the early detection of the disease in a subject's brain scan along with an intermediate prodromal stage known as MCI.

Some of such factors included the need to partake in the data received and recorded both from image data gathered from Neuro imaging techniques like that of MRI, PET and CSF Scans as well as the data recorded from cognitive information like as in lifestyle and psychological examinations. Employing this approach, we had our system gather information from fMRI dataset, along with ADNI's own challenge dataset. Due to reasons of MRI providing more spatial information as compared to the more functional attributes that were gathered from PET and CSF Scans, we believed in making use of fMRI for our analysis.

Providing fMRI analysis allowed us to gather the role of the Hippocampal region to help aid in the detection of the Disease. Information gathered from the voxel count of the Hippocampal region along with grey matter and white matter concentration in the particular regions both for the Left and Right Hippocampal regions allowed us to aid in the identification of the disease along with other factors like a patient's level of education MMSE test scores and so on.

Another important factor that provided with major emphasis of using ensemble training techniques as opposed to making use of the more standardized or adapted Machine Learning techniques. From our evaluation, we were able to see that making use of the ensemble

technique of Random Forest Evaluation we were able to get an estimated weighted average precision of 97%-99%. Furthermore, an AUC Level of 95%-99% has been achieved as compared to using adapted ensemble training techniques like Random Forest technique in distillation with SVM and the nearest neighbor method. We also looked into the effectiveness of certain machine learning techniques preferentially that of Support Vector Mechine which provided better results as opposed to other ML techniques like the 3KNN Methods.

Also to bide time in the process of Machine Learning, we focused into reducing the dimensionality constraint that proves to play an effective role in simplifying the dataset to easily perform the required Machine Learning techniques. To this end, we performed an Eigen transformation of the image information provided using principal component analysis to standardize and correlate the regions of interest.

In overall our model has effectively been able to correlate the role on Hippocampal volume and the application of ensemble training techniques to prove in the effectiveness of having a system learn to identify the presence of AD in a subject's brain. But we also have to take into account the solidarity of our statement. As far as our model has been concerned, we have only been able to make accounted predictions through previous research in the role of Hippocampal region of interest to aid in the early detection of AD. Moreover, we did not actively look into the role of Voxel based morphometry in the detection of the disease. Moreover, as discussed earlier many other methods like as in looking into neurofibrillary tangles and β-Amyloid concentration which can be gathered from other imaging techniques like PET Scans and CSF concentration can look into other aspects of detection of the disease like protein concentration, metabolic rate, and glucose concentration. In general there are a lot of other factors aside from hippocampal region that can help aid in the detection of the disease.

However it can also be gathered from our results that the Hippocampal region and the hippocampal volume both are very efficient means of detecting the onset of the diseas.

# 5. Conclusion

When our research was initiated, we were faced with proposing the question, can a system learn with as little human intervention and through means of functional brain scans, to identify if a patient is likely over a periodof time to be diagnosed with Alzheimer's Disease or not? To give an answer to this question, it would be yes. Many contributed research have proven the emphasis of applying machine learning and neural network techniques to aid in the detection of the disease.

In our experiments, we had patients undergo a limited amount of MRI Screening tests, along with providing their personal information such as age of taking the test as well as gender, race and MMSE Score results as such. Having as little as 4-5 sessions of MRI Screening and a limited number of tests was able to provide us with an effective method of able to detect the onset of the disease. Although effective as it may be, our model is still capable of being quite imprecise in the detection of the disease. Research in the area has only touched the tip of the iceberg in claiming certain factors that may aid in the detection of disease. For instance, while hippocampal volume detection as claimed by our system can aid in providing 97% accuracy in detecting the onset of the disease, there could be other methods like identifying glucoseconcentration and metabolic activity in the brain thst could completely refute our findings.

Future scope in the area can be built upon studying the inner workings of the cortical subregions of the brain, and deciding machine learning techniques that can find a pattern in detecting which subregion ought to provide the maximum weight towards chanf=ges in neuronal activity.

Other inclusive studies can be making use of other features and neuro imaging techniques like PET Scans and CSF Scans in detecting the early onset of the disease. In our model we looked into some of the most widely used techniques including SVM, Random Forest and K Nearest Neighbour techniques. The future brings with it, a lot more efficient models, that can one day potentially be used to help with aiding the detection of the disease. Also upon training, we checked the data as of 150 subjects, providing a higer number, can prove to be more beneficial in aiding in the Machine Learning process.

In retrospective claims, we can gather that substantial research has proved into aiming to detect the onset of the disease, both through neurology and through artificial intelligence and machine learning. However, the claims cease to end, and there are lot more untapped areas that can promote to detecting the disease and helping to find a potential cure to this deadly disease.

# 6. References

[1]    K. S. Biju, S. S. Alfa, K. Lal, A. Antony, and M. K. Akhil, "Alzheimer's Detection Based on Segmentation of MRI Image," Procedia Comput. Sci., vol. 115, pp. 474–481, 2017.

[2]    C. Fang et al., "A novel Gaussian discriminant analysis-based computer aided diagnosis system for screening different stages of Alzheimer's disease," Proc. - 2017 IEEE 17th Int. Conf. Bioinforma. Bioeng. BIBE 2017, vol. 2018-Janua, pp. 279–284, 2018.

[3]    H. Yu, X. Lei, Z. Song, C. Liu, and J. Wang, "Supervised Network-based Fuzzy Learning of EEG Signals for Alzheimer's Disease Identification," IEEE Trans. Fuzzy Syst., vol. PP, no. c, pp. 1–1, 2019.

[4]    S. Li et al., "Hippocampal shape analysis of Alzheimer disease based on machine learning methods," Am. J. Neuroradiol., vol. 28, no. 7, pp. 1339–1345, 2007.

[5]    S. Joshi, D. Shenoy, G. G. Vibhudendra Simha, P. L. Rrashmi, K. R. Venugopal, and L. M. Patnaik, "Classification of Alzheimer's disease and Parkinson's disease by using machine learning and neural network methods," ICMLC 2010 - 2nd Int. Conf. Mach. Learn. Comput., pp. 218–222, 2010.

[6]    S. Sarraf and G. Tofighi, "Deep learning-based pipeline to recognize Alzheimer's disease using fMRI data," FTC 2016 - Proc. Futur. Technol. Conf., no. December, pp. 816–820, 2017.

[7]    S. Wang, H. Wang, Y. Shen, and X. Wang, "Automatic Recognition of Mild Cognitive Impairment and Alzheimers Disease Using Ensemble based 3D Densely Connected Convolutional Networks," Proc. - 17th IEEE Int. Conf. Mach. Learn. Appl. ICMLA 2018, pp. 517–523, 2019.

[8]    J. Ye, T. Wu, and J. Li, "Machine Learning Approaches for the Neuroimaging Study of Alzheimer ' s Disease," no. April, pp. 77–79, 2011.

[9]    M. Mahyoub, M. Randles, T. Baker, and P. Yang, "Comparison analysis of machine learning

algorithms to rank Alzheimer's disease risk factors by importance," Proc. - Int. Conf. Dev. eSystems Eng. DeSE, vol. 2018-Septe, pp. 1–11, 2019.

[10]    C. R. Jack et al., "Magnetic resonance imaging in Alzheimer's Disease Neuroimaging Initiative 2," Alzheimer's Dement., vol. 11, no. 7, pp. 740–756, 2015.

[11]    D. Junwei and H. Qiu, "Prediction of MCI to AD conversion using Laplace Eigenmaps learned from FDG and MRI images of AD patients and healthy controls," 2017 2nd Int. Conf. Image, Vis. Comput. ICIVC 2017, pp. 660–664, 2017.

[12]    T. Tong, K. Gray, Q. Gao, L. Chen, and D. Rueckert, "Multi-modal classification of Alzheimer's disease using nonlinear graph fusion," Pattern Recognit., vol. 63, no. October 2016, pp. 171–181, 2017.

[13]    D. Zhang and D. Shen, "Multi-modal multi-task learning for joint prediction of multiple regression and classification variables in Alzheimer's disease," Neuroimage, vol. 59, no. 2, pp. 895–907, 2012.

[14]    E. Moradi, A. Pepe, C. Gaser, H. Huttunen, and J. Tohka, "Machine learning framework for early MRI-based Alzheimer's conversion prediction in MCI subjects," Neuroimage, vol. 104, pp. 398–412, 2015.

[15]    J. Peng, X. Zhu, Y. Wang, L. An, and D. Shen, "Structured sparsity regularized multiple kernel learning for Alzheimer's disease diagnosis," Pattern Recognit., vol. 88, pp. 370–382, 2019.

[16]    S. Korolev, A. Safiullin, M. Belyaev, and Y. Dodonova, "Residual and plain convolutional neural networks for 3D brain MRI classification," Proc. - Int. Symp. Biomed. Imaging, pp. 835–838, 2017.

[17]    E. Moradi, J. Tohka, and C. Gaser, "Semi-supervised learning in MCI-to-ad conversion prediction - When is unlabeled data useful?," Proc. - 2014 Int. Work. Pattern Recognit.

Neuroimaging, PRNI 2014, pp. 1–4, 2014.

[18]   O. Ben Ahmed, J. Benois-Pineau, M. Allard, G. Catheline, and C. Ben Amar, "Recognition of Alzheimer's disease and Mild Cognitive Impairment with multimodal image-derived biomarkers and Multiple Kernel Learning," Neurocomputing, vol. 220, pp. 98–110, 2017.

[19]   V. Kebets et al., "Predicting Pure Amnestic Mild Cognitive Impairment Conversion to Alzheimer's Disease Using Joint Modeling of Imaging and Clinical Data," Proc. - 2015 Int. Work. Pattern Recognit. NeuroImaging, PRNI 2015, pp. 85–88, 2015.

[20]   A. Majumdar and V. Singhal, "Noisy deep dictionary learning: Application to Alzheimer's Disease classification," Proc. Int. Jt. Conf. Neural Networks, vol. 2017-May, pp. 2679–2683, 2017.

[21]   L. He et al., "Multi-way multi-level Kernel modeling for neuroimaging classification," Proc. - 30th IEEE Conf. Comput. Vis. Pattern Recognition, CVPR 2017, vol. 2017-Janua, pp. 6846–6854, 2017.

[22]   L. Sørensen et al., "Differential diagnosis of mild cognitive impairment and Alzheimer's disease using structural MRI cortical thickness, hippocampal shape, hippocampal texture, and volumetry," NeuroImage Clin., vol. 13, pp. 470–482, 2017.

[23]   Y. Cabrera-León, P. G. Báez, J. Ruiz-Alzola, and C. P. Suárez-Araujo, "Classification of Mild Cognitive Impairment Stages Using Machine Learning Methods," INES 2018 - IEEE 22nd Int. Conf. Intell. Eng. Syst. Proc., pp. 000067–000072, 2018.

[24]   K. A. N. N. P. Gunawardena, R. N. Rajapakse, and N. D. Kodikara, "Applying convolutional neural networks for pre-detection of Alzheimer's disease from structural MRI data," 2017 24th Int. Conf. Mechatronics Mach. Vis. Pract. M2VIP 2017, vol. 2017-Decem, pp. 1–7, 2017.

[25]   G. Fiscon et al., "Alzheimer's disease patients classification through EEG signals processing," IEEE SSCI 2014 - 2014 IEEE Symp. Ser. Comput. Intell. - CIDM 2014 2014 IEEE Symp.

Comput. Intell. Data Mining, Proc., pp. 105–112, 2015.

[26]     V. Sachnev and B. S. Mahanand, "A Cognitive Ensemble Classifier based on Risk Sensitive Hinge Loss Function for Alzheimer's Disease diagnosis in early stages," Proc. 2018 IEEE Symp. Ser. Comput. Intell. SSCI 2018, pp. 807–812, 2019.

[27]     Y. Zhang et al., "Detection of subjects and brain regions related to Alzheimer's disease using 3D MRI scans based on eigenbrain and machine learning," Front. Comput. Neurosci., vol. 9, no. June, pp. 1–15, 2015.

# 7. Bibliography

[1]     Alzheimer's Disease Neuroimaging Initiative. (2019). *Home*. Retrieved from
        ADNI: https://adni.loni.usc.edu

[2]     Cuingnet, R., Gerardin, E., Tessieras, J., & Auzias, G. (2009). Automatic
        classification of patients with Alzheimer's disease from structural MRI: A. *Elsevier*.

[3]     Gupta, A., Ayhan, M. S., & Maida, A. S. (n.d.). Natural Image Bases to
        Represent Neuroimaging Data.

[4]     Hosseini-Asl, E., Gimel'farb2, G., El-Baz3, A., & Initiative, f. t. (n.d.). ALZHEIMER'S
        DISEASE DIAGNOSTICS BY A DEEPLY SUPERVISED ADAPTABLE 3D.

[5]     Huang*, A., Abugharbieh, R., & Tam, R. (2009). A Hybrid Geometric–Statistical Deformable
        Model. *IEEE TRANSACTIONS ON BIOMEDICAL ENGINEERING, VOL. 56, NO. 7, JULY*.

# 8. Appendices

**Acquiring Image Data.ipynb**  **Creating 1st Phase Analysis.ipynb**  **Acquiring Data Phase 2.ipynb**  **2nd Phase Machine Learning.ipynb**

.

Attached is the code file in ipynb format.