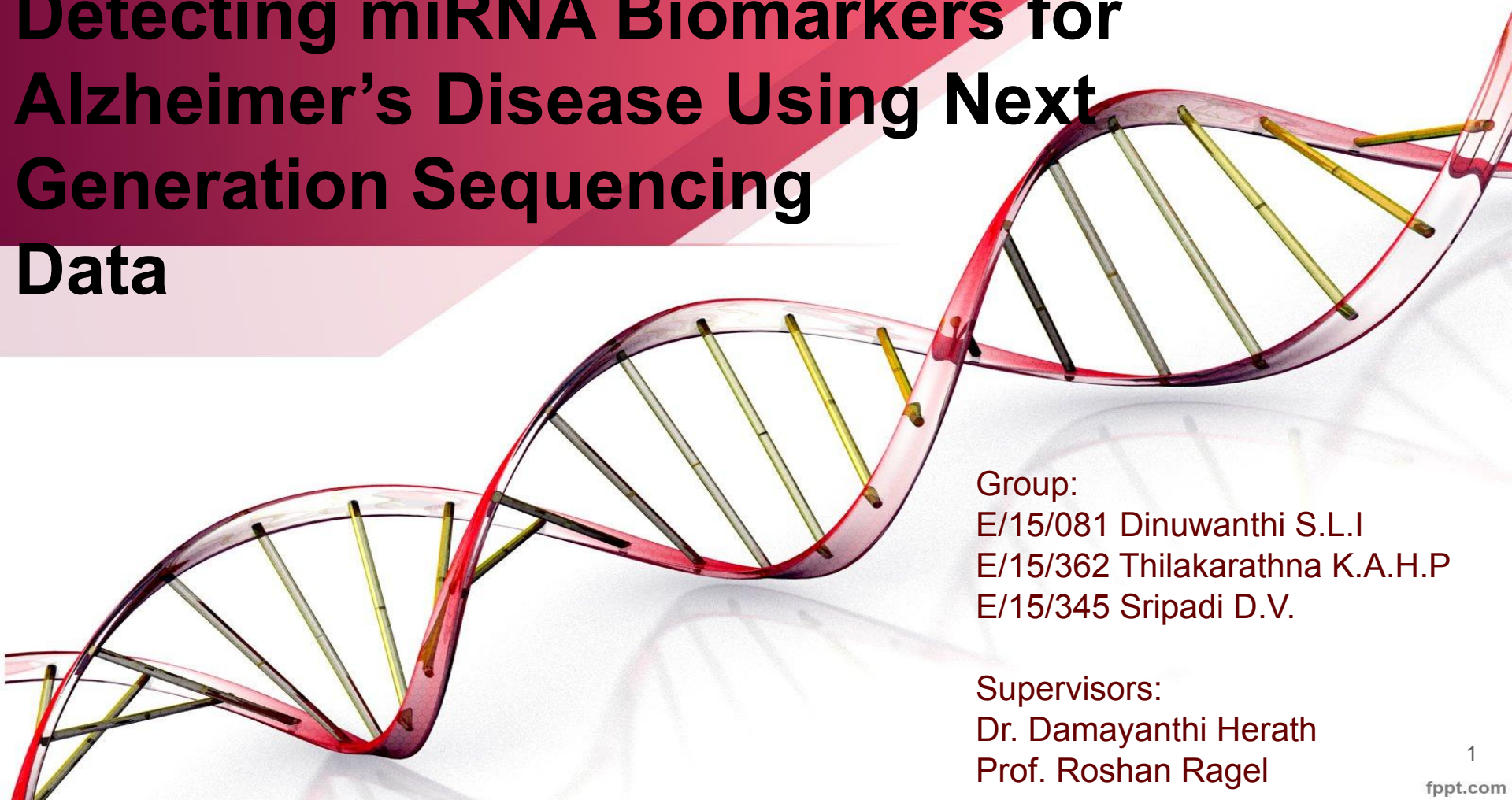


Detecting miRNA Biomarkers for Alzheimer's Disease Using Next Generation Sequencing Data



Group:

E/15/081 Dinuwanthi S.L.I

E/15/362 Thilakarathna K.A.H.P

E/15/345 Sripadi D.V.

Supervisors:

Dr. Damayanthi Herath

Prof. Roshan Ragel

Background



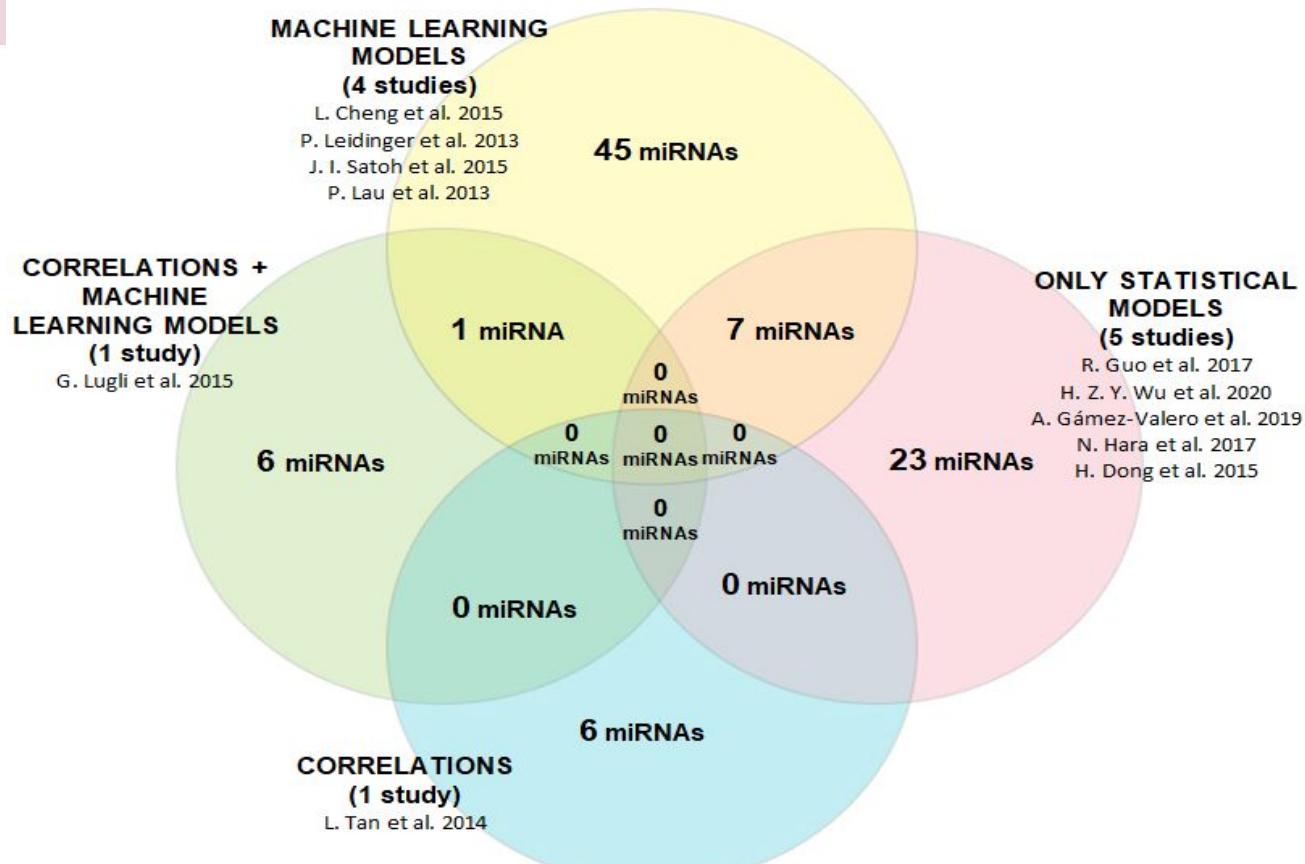
- What is Next Generation Sequencing ?
 - Sample preparation
 - Sequencing by machines
 - Data output
- What are miRNAs ?

Problem definition

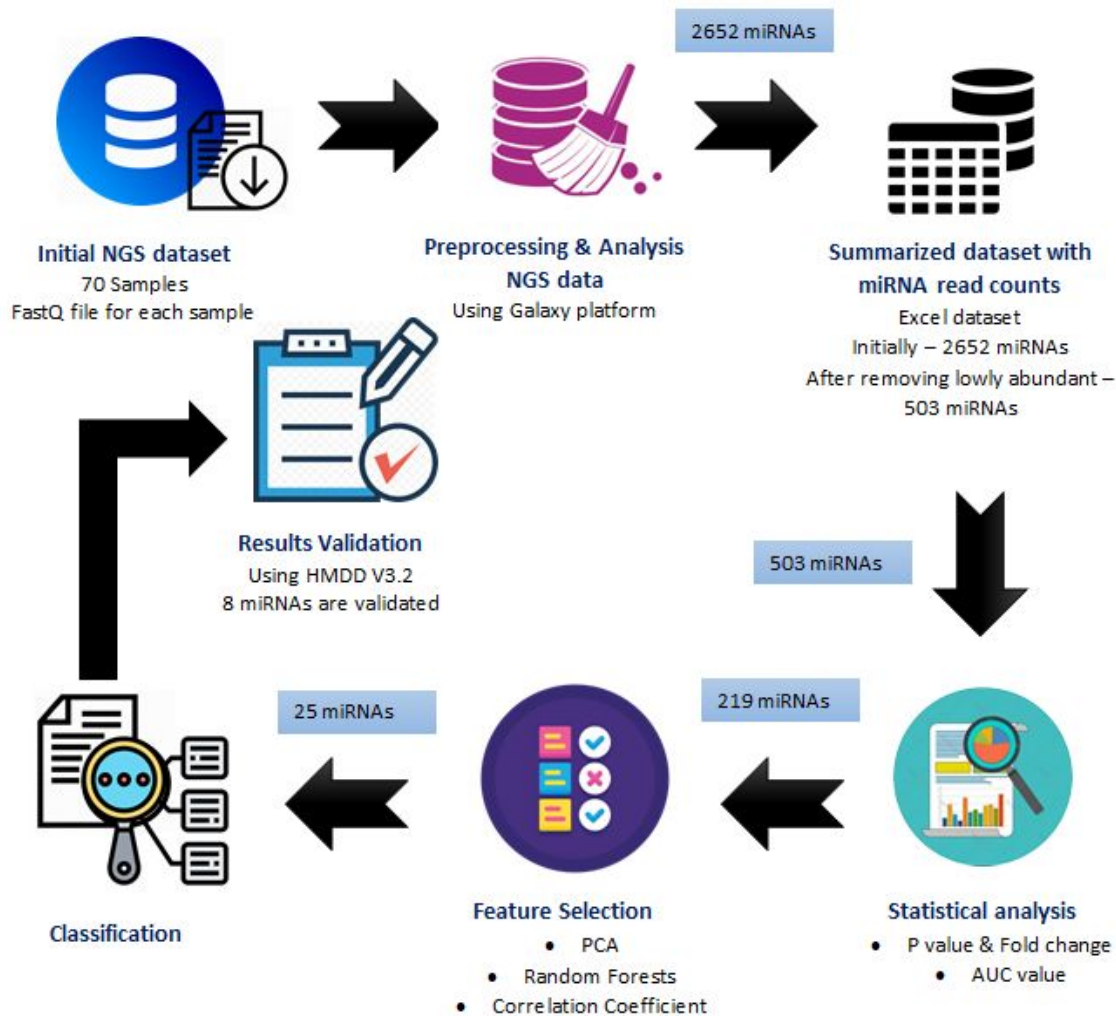


- Finding a method for identifying miRNA biomarkers for Alzheimer's Disease
- Need of new method because of the drawbacks of previously introduced methods

Related works



Methodology



Implementation choices: Tools



NGS data analysis

- Galaxy platform - web based, open source platform
- Galaxy tools
 - FastQC - Quality check
 - Trim Galore - remove adapters and low quality reads
 - Filter FASTQ - filter short read sequences
 - Bowtie2 - map sequence against reference genome
 - Htseq-count - identify read counts

Implementation choices: Tools



Statistical Analysis

- Python - clean syntax, straightforward semantics & Third-party toolkits
- Python libraries
 - Numpy
 - Pandas
 - Scipy
 - Scikit learn
 - Matplotlib

Implementation choices: Models



Classification

- Machine Learning models
 - Logistic regression
 - Linear SVM
 - Gaussian SVM
 - Naive Bayes
 - K Nearest Neighbour
 - Random Forests

Implementation choices: Methodologies



Preprocessing

- Trimming sequence data - remove adapters, indexes, low quality reads
- Filtering - remove short read sequences



Statistical analysis

- P value & Fold change - identify most significance miRNAs
- AUC - identify dysregulated miRNAs

Implementation choices: Methodologies contd.



Feature selection

- PCA
- Random Forest
- Correlation coefficient



Classification

- Classification - Support vector machine, Logistic regression and Random Forest

● NGS data analysis



```
TCCTGTACTGAGCTGCCCCGAGATGGAATTCTCGGGTGCCAAGGAACTCC
TCCTGTACTGAGCTGCCCCGAGGGGAATTCTGGGGTGCCAAGGAACCCCA
TCCTGTACTGAGCTGCCCCGAGTGAATTCTCGGGTGCCAAGGAACTCCA
```

First 3 sequences from fastQ file



Trimming using
Trim Galore

```
TCCTGTACTGAGCTGCCCCGAGA
TCCTGTACTGAGCTGCCCCGAGGGGAATTCTGGGG
TCCTGTACTGAGCTGCCCCGAG
```

After trimming



Filtering using
Filter FASTQ

```
TCCTGTACTGAGCTGCCCCGAGA
TCCTGTACTGAGCTGCCCCGAG
```

After filtering
Second sequence was removed



Mapping against
reference genome

Reference genome SN:chr1

```
..TGCCCCGAGTCCTGTACTGAGCTGCCCCGAGTCCTGTACTGA...
```

```
TCCTGTACTGAGCTGCCCCGAGA
```

Seq 1
SN:chr1

Mapping against
miRBase



miRBase sequence (From has.gtf file) Name=hsa-mir-200a

```
..TGCCCCGAGTCGTTCCGTCTCTGTACTGAGCTGCCCCGAGATCCT...
```

```
TCCTGTACTGAGCTGCCCCGAGA
```

Seq 1

Seq 1 identified as hsa-mir-200a

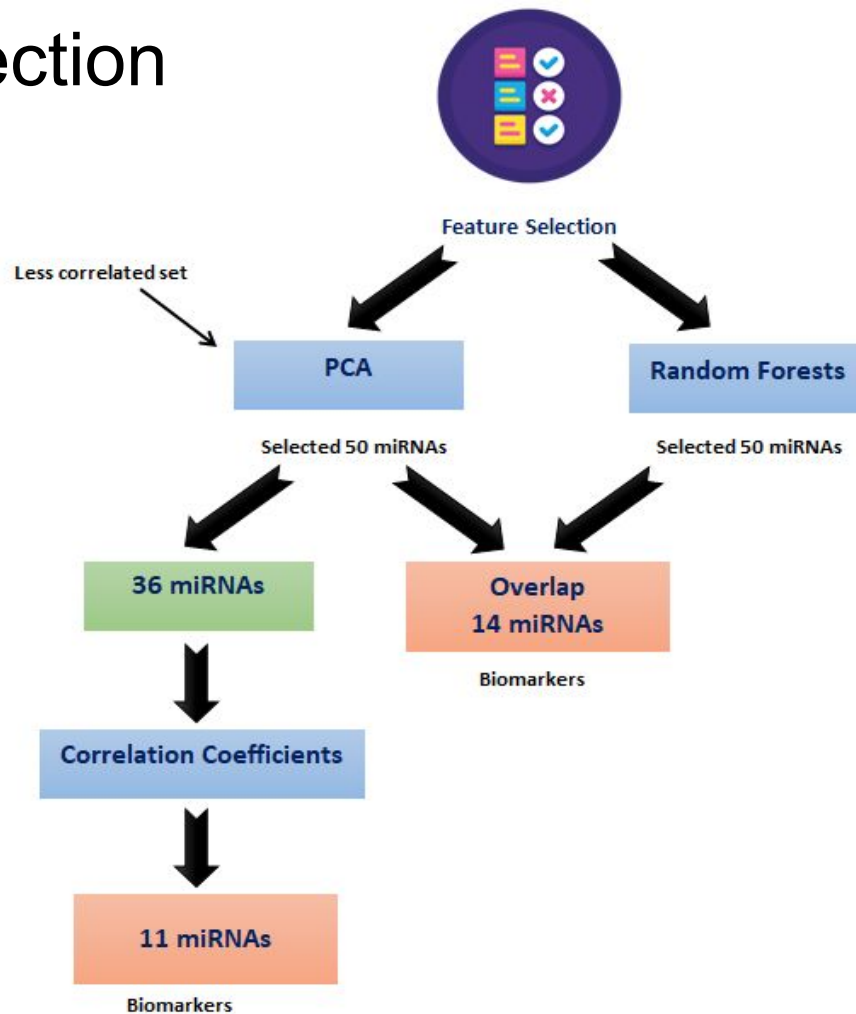


Count reads of
each miRNA

miRNA name	Read counts
hsa-miR-106b-3p	2622
hsa-miR-107	594
hsa-miR-10a-3p	1
hsa-mir-200a	45
hsa-miR-1269a	0

After counting reads for all miRNA

● Feature Selection



Results Obtained From Each Stage of Analysis



- Filtered miRNAs Using Significance Values and Fold Change
 - 228 miRNAs

Filtered miRNAs using significance value and fold change are: ['hsa-mir-30a:hsa-miR-30a-3p', 'hsa-mir-550a-1:hsa-miR-550a-3p', 'hsa-mir-29a:hsa-miR-29a-3p', 'hsa-mir-628:hsa-miR-628-3p', 'hsa-mir-26a-2:hsa-miR-26a-5p', 'hsa-mir-106b:hsa-miR-106b-5p', 'hsa-mir-4781:hsa-miR-4781-3p', 'hsa-mir-10b:hsa-miR-10b-5p', 'hsa-mir-215:hsa-miR-215', 'hsa-mir-548aj-2:hsa-miR-548g-5p', 'hsa-mir-181a-1:hsa-miR-181a-3p', 'hsa-mir-548x:hsa-miR-548ar-5p', 'hsa-mir-548k:hsa-miR-548av-5p', 'hsa-mir-199a-1:hsa-miR-199a-3p', 'hsa-mir-30e:hsa-miR-30e-3p', 'hsa-mir-4508:hsa-miR-4508', 'hsa-mir-548aj-2:hsa-miR-548x-5p', 'hsa-mir-371b:hsa-miR-371b-5p', 'hsa-mir-5001:hsa-miR-5001-3p', 'hsa-mir-16-2:hsa-miR-16-2-3p', 'hsa-mir-128-2:hsa-miR-128', 'hsa-mir-486:hsa-miR-486-3p', 'hsa-mir-4482-1:hsa-miR-4482-3p', 'hsa-mir-941-4:hsa-miR-941', 'hsa-mir-550a-1:hsa-miR-550a-5p', 'hsa-mir-199a-2:hsa-miR-199b-3p', 'hsa-mir-144:hsa-miR-144-5p', 'hsa-let-7f-2:hsa-let-7f-5p', 'hsa-mir-126:hsa-miR-126-5p', 'hsa-mir-191:hsa-miR-191-3p', 'hsa-mir-10a:hsa-miR-10a-5p', 'hsa-mir-98:hsa-miR-98', 'hsa-mir-548x:hsa-miR-548x-5p', 'hsa-mir-363:hsa-miR-363-3p', 'hsa-mir-548h-1:hsa-miR-548h-5p', 'hsa-mir-223:hsa-miR-223-3p', 'hsa-mir-5690:hsa-miR-5690', 'hsa-mir-199b:hsa-miR-199b-3p', 'hsa-mir-3200:hsa-miR-3200-3p', 'hsa-mir-424:hsa-miR-424-3p', 'hsa-mir-644b:hsa-miR-644b-3p', 'hsa-mir-548h-5:hsa-miR-548h-5p', 'hsa-mir-18a:hsa-miR-18a-5p', 'hsa-mir-548g:hsa-miR-548x-5p', 'hsa-mir-548g:hsa-miR-548g-5p', 'hsa-mir-21:hsa-miR-21-5p', 'hsa-mir-99b:hsa-miR-99b-5p', 'hsa-mir-25:hsa-miR-25-3p', 'hsa-mir-937:hsa-miR-937', 'hsa-mir-1180:hsa-miR-1180', 'hsa-mir-30c-1:hsa-miR-30c-5p', 'hsa-let-7a-1:hsa-let-7a-5p', 'hsa-mir-941-1:hsa-miR-941', 'hsa-mir-660:hsa-miR-660-5p', 'hsa-mir-421:hsa-miR-421', 'hsa-mir-374a:hsa-miR-374a-5p', 'hsa-mir-328:hsa-miR-328', 'hsa-mir-151a:hsa-miR-151a-5p', 'hsa-mir-548x:hsa-miR-548aj-5p', 'hsa-mir-101-2:hsa-miR-101-3p', 'hsa-mir-28:hsa-miR-28-3p', 'hsa-mir-139:hsa-miR-139-5p', 'hsa-mir-2110:hsa-miR-2110', 'hsa-let-7g:hsa-let-7g-5p', 'hsa-mir-550a-3:hsa-miR-550a-3-5p', 'hsa-mir-548aj-2:hsa-miR-548ar-5p', 'hsa-mir-144:hsa-miR-144-3p', 'hsa-mir-548e:hsa-miR-548e-5p']



● Filtered miRNAs Using AUC Analysis

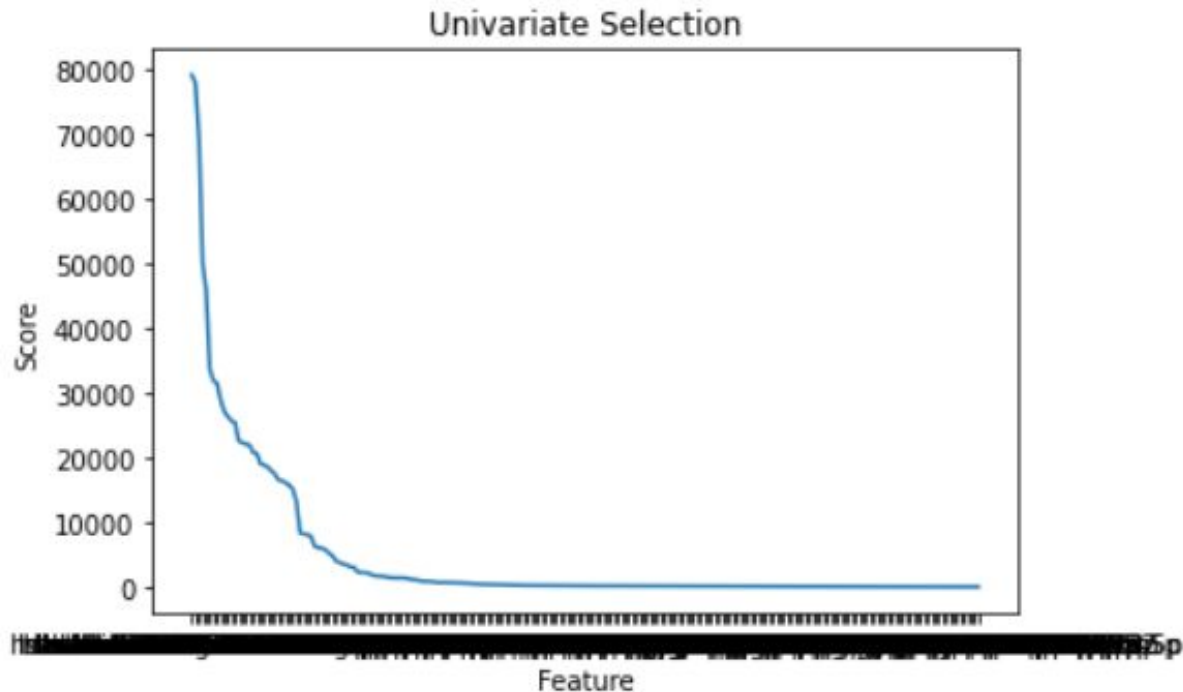
- 219 miRNAs
 - 179 down-regulated miRNAs
 - 40 up-regulated miRNAs

Selected miRNAs using AUC analysis: ['hsa-mir-30a:hsa-mir-30a-3p', 'hsa-mir-550a-1:hsa-mir-550a-3p', 'hsa-mir-29a:hsa-mir-29a-3p', 'hsa-mir-628:hsa-mir-628-3p', 'hsa-mir-26a-2:hsa-mir-26a-5p', 'hsa-mir-106b:hsa-mir-106b-5p', 'hsa-mir-4781:hsa-mir-4781-3p', 'hsa-mir-10b:hsa-mir-10b-5p', 'hsa-mir-215:hsa-mir-215', 'hsa-mir-548aj-2:hsa-mir-548g-5p', 'hsa-mir-181a-1:hsa-mir-181a-3p', 'hsa-mir-548x:hsa-mir-548ar-5p', 'hsa-mir-548k:hsa-mir-548av-5p', 'hsa-mir-199a-1:hsa-mir-199a-3p', 'hsa-mir-30e:hsa-mir-30e-3p', 'hsa-mir-4508:hsa-mir-4508', 'hsa-mir-548aj-2:hsa-mir-548x-5p', 'hsa-mir-371b:hsa-mir-371b-5p', 'hsa-mir-5001:hsa-mir-5001-3p', 'hsa-mir-16-2:hsa-mir-16-2-3p', 'hsa-mir-128-2:hsa-mir-128', 'hsa-mir-486:hsa-mir-486-3p', 'hsa-mir-4482-1:hsa-mir-4482-3p', 'hsa-mir-550a-1:hsa-mir-550a-5p', 'hsa-mir-199a-2:hsa-mir-199b-3p', 'hsa-mir-144:hsa-mir-144-5p', 'hsa-let-7f-2:hsa-let-7f-5p', 'hsa-mir-126:hsa-mir-126-5p', 'hsa-mir-191:hsa-mir-191-3p', 'hsa-mir-10a:hsa-mir-10a-5p', 'hsa-mir-98:hsa-mir-98', 'hsa-mir-548x:hsa-mir-548x-5p', 'hsa-mir-363:hsa-mir-363-3p', 'hsa-mir-548h-1:hsa-mir-548h-5p', 'hsa-mir-223:hsa-mir-223-3p', 'hsa-mir-5690:hsa-mir-5690', 'hsa-mir-199b:hsa-mir-199b-3p', 'hsa-mir-3200:hsa-mir-3200-3p', 'hsa-mir-424:hsa-mir-424-3p', 'hsa-mir-644b:hsa-mir-644b-3p', 'hsa-mir-548h-5:hsa-mir-548h-5p', 'hsa-mir-18a:hsa-mir-18a-5p', 'hsa-mir-548g:hsa-mir-548x-5p', 'hsa-mir-548g:hsa-mir-548g-5p', 'hsa-mir-21:hsa-mir-21-5p', 'hsa-mir-99b:hsa-mir-99b-5p', 'hsa-mir-25:hsa-mir-25-3p', 'hsa-mir-937:hsa-mir-937', 'hsa-mir-1180:hsa-mir-1180', 'hsa-mir-30c-1:hsa-mir-30c-5p', 'hsa-let-7a-1:hsa-let-7a-5p', 'hsa-mir-660:hsa-mir-660-5p', 'hsa-mir-421:hsa-mir-421', 'hsa-mir-374a:hsa-mir-374a-5p', 'hsa-mir-151a:hsa-mir-151a-5p', 'hsa-mir-548x:hsa-mir-548aj-5p', 'hsa-mir-101-2:hsa-mir-101-3p', 'hsa-mir-28:hsa-mir-28-3p', 'hsa-mir-139:hsa-mir-139-5p', 'hsa-mir-2110:hsa-mir-2110', 'hsa-let-7g:hsa-let-7g-5p', 'hsa-mir-550a-3:hsa-mir-550a-3-5p', 'hsa-mir-548aj-2:hsa-mir-548ar-5p', 'hsa-mir-144:hsa-mir-144-3p', 'hsa-mir-548e:hsa-mir-548e', 'hsa-mir-3074:hsa-mir-3074-5p', 'hsa-mir-1294:hsa-mir-1294', 'hsa-mir-19a:hsa-mir-19a-3p', 'hsa-mir-199a-1:hsa-mir-199b-3p', 'hsa-mir-17:hsa-mir-17-3p', 'hsa-mir-340:hsa-mir-340-3p', 'hsa-mir-3158-2:hsa-mir-3158-3p', 'hsa-mir-548x:hsa-mir-548g-5p', 'hsa-mir-331:hsa-mir-331-3p', 'hsa-mir-4742:hsa-mir-4742-3p', 'hsa-mir-4482-2:hsa-mir-4482-3p', 'hsa-mir-548h-2:hsa-mir-548h-5p', 'hsa-mir-324:hsa-mir-324-5p', 'hsa-mir-18b:hsa-mir-18b-5p', 'hsa-mir-659:hsa-mir-659-5p', 'hsa-mir-532:hsa-mir-532-5p', 'hsa-mir-3158-1:hsa-mir-3158-3p', 'hsa-mir-671:hsa-mir-671-3p', 'hsa-let-7d:hsa-let-7d-3p', 'hsa-mir-29b-2:hsa-mir-29b-3p', 'hsa-mir-548h-4:hsa-mir-548h-5p', 'hsa-mir-3615:hsa-mir-3615', 'hsa-mir-4746:hsa-mir-4746-5p', 'hsa-mir-210:hsa-mir-210', 'hsa-mir-3127:hsa-mir-3127-3p', 'hsa-mir-130b:hsa-mir-130b-3p', 'hsa-mir-550a-2:hsa-mir-550a-5p', 'hsa-let-7a-2:hsa-let-7a-5p', 'hsa-mir-148a:hsa-mir-148a-3p', 'hsa-mir-190a:hsa-mir-190a', 'hsa-mir-1304:hsa-mir-1304-3p', 'hsa-mir-4792:hsa-mir-4792', 'hsa-mir-106b:hsa-mir-106b-3p', 'hsa-mir-29b-1:hsa-mir-29b-3p', 'hsa-mir-30b:hsa-mir-30b-5p', 'hsa-mir-378d-2:hsa-mir-378d', 'hsa-mir-20a:hsa-mir-20a-5p', 'hsa-mir-323b:hsa-mir-323b-3p', 'hsa-mir-548ar:hsa-mir-548ar-5p', 'hsa-mir-378i:hsa-mir-378

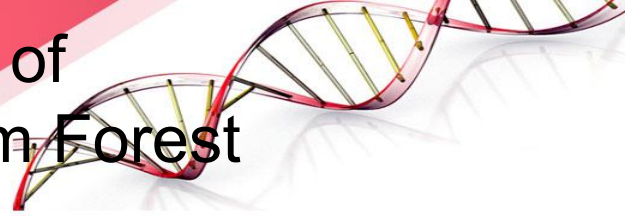
- Univariate selection to identify most related miRNAs



- 50 miRNAs



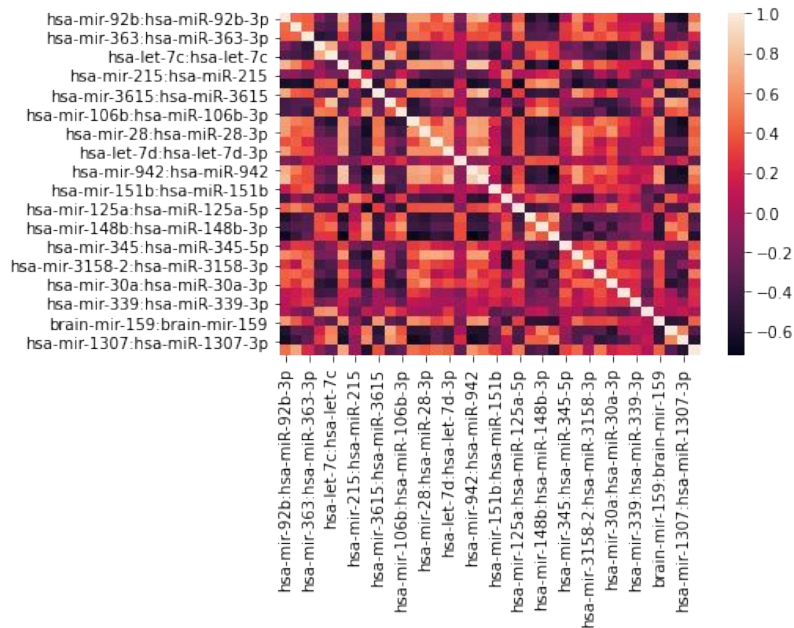
- Visualization of the overlap between set of miRNAs selected from PCA and Random Forest



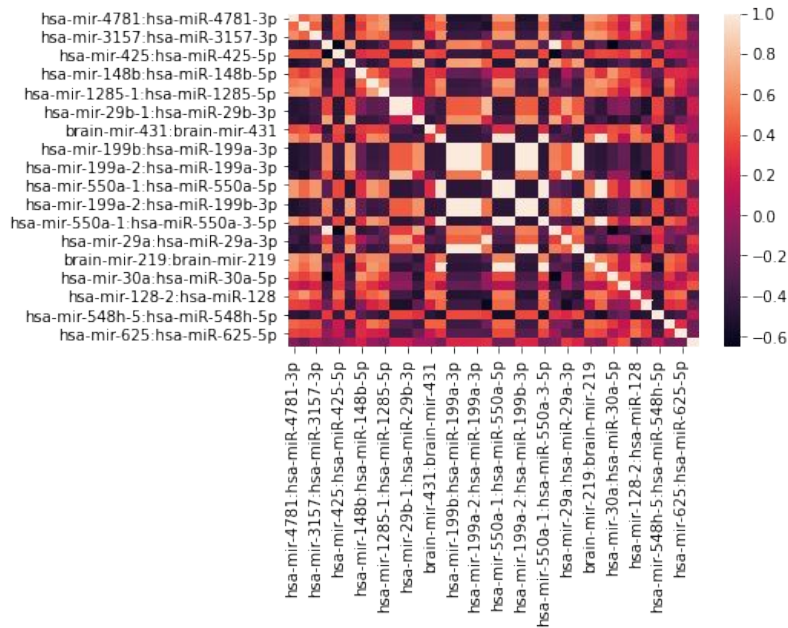
● Correlation heatmaps



PCA



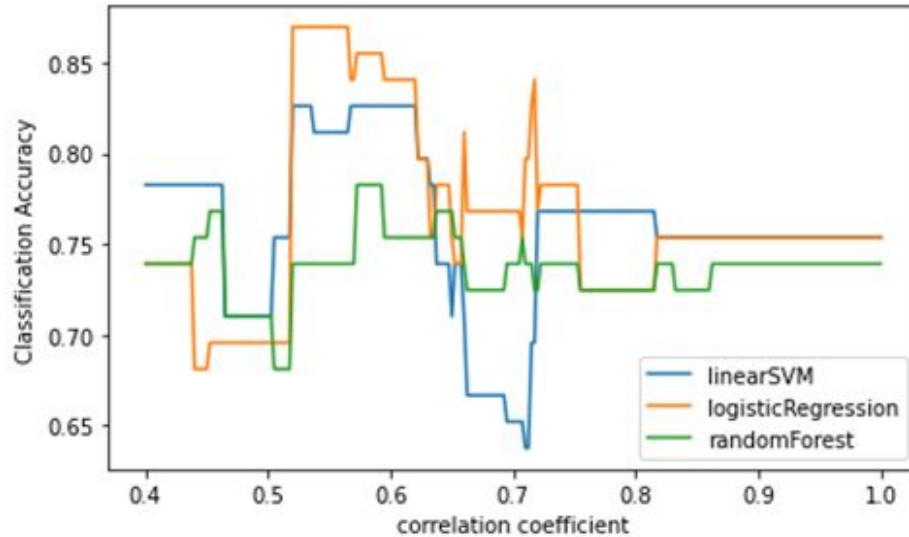
Random Forest



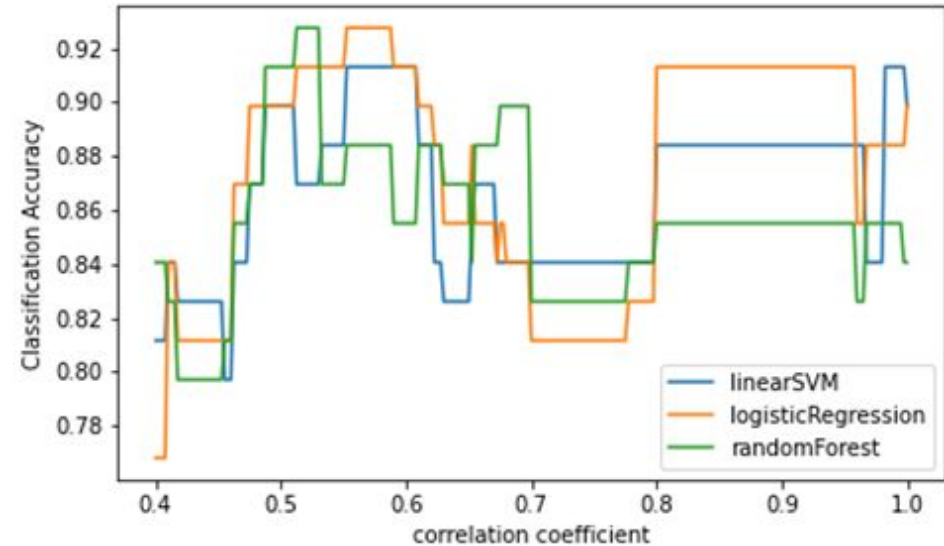
- Accuracy with correlation coefficients



PCA



Random Forest





- Selected miRNAs using correlations

- 11 miRNAs

Correlation
hsa-mir-4781:hsa-miR-4781-3p
brain-mir-112:brain-mir-112
hsa-let-7a-3:hsa-let-7a-5p
hsa-mir-148b:hsa-miR-148b-5p
hsa-mir-29b-2:hsa-miR-29b-3p
brain-mir-431:brain-mir-431
hsa-mir-378a:hsa-miR-378a-5p
hsa-mir-548h-5:hsa-miR-548h-5p
hsa-mir-3909:hsa-miR-3909
hsa-mir-625:hsa-miR-625-5p
hsa-mir-24-1:hsa-miR-24-3p



● Selected miRNAs using Overlap

● 14 miRNAs

overlap

hsa-mir-186:hsa-miR-186-5p

hsa-mir-144:hsa-miR-144-3p

hsa-mir-151a:hsa-miR-151a-3p

hsa-mir-99b:hsa-miR-99b-5p

hsa-mir-98:hsa-miR-98

hsa-mir-148a:hsa-miR-148a-3p

hsa-let-7g:hsa-let-7g-5p

hsa-let-7f-2:hsa-let-7f-5p

hsa-let-7a-1:hsa-let-7a-5p

hsa-mir-30d:hsa-miR-30d-5p

hsa-mir-15a:hsa-miR-15a-5p

hsa-mir-589:hsa-miR-589-5p

hsa-mir-144:hsa-miR-144-5p

hsa-let-7f-1:hsa-let-7f-5p

- Accuracies



	Both PCA & RF (overlap)	Overlap + Correlation
Linear SVM	80.95	85.71
Logistic Regression	95.24	95.24
Random Forest	90.48	95.24

Validation of the results



- HMDD v3.2 was used with the final set of 25 miRNAs
 - 10 miRNAs were identified
 - 2 from less correlated set
 - 8 from overlapped data from both PCA & RF

```
corr_validate
```

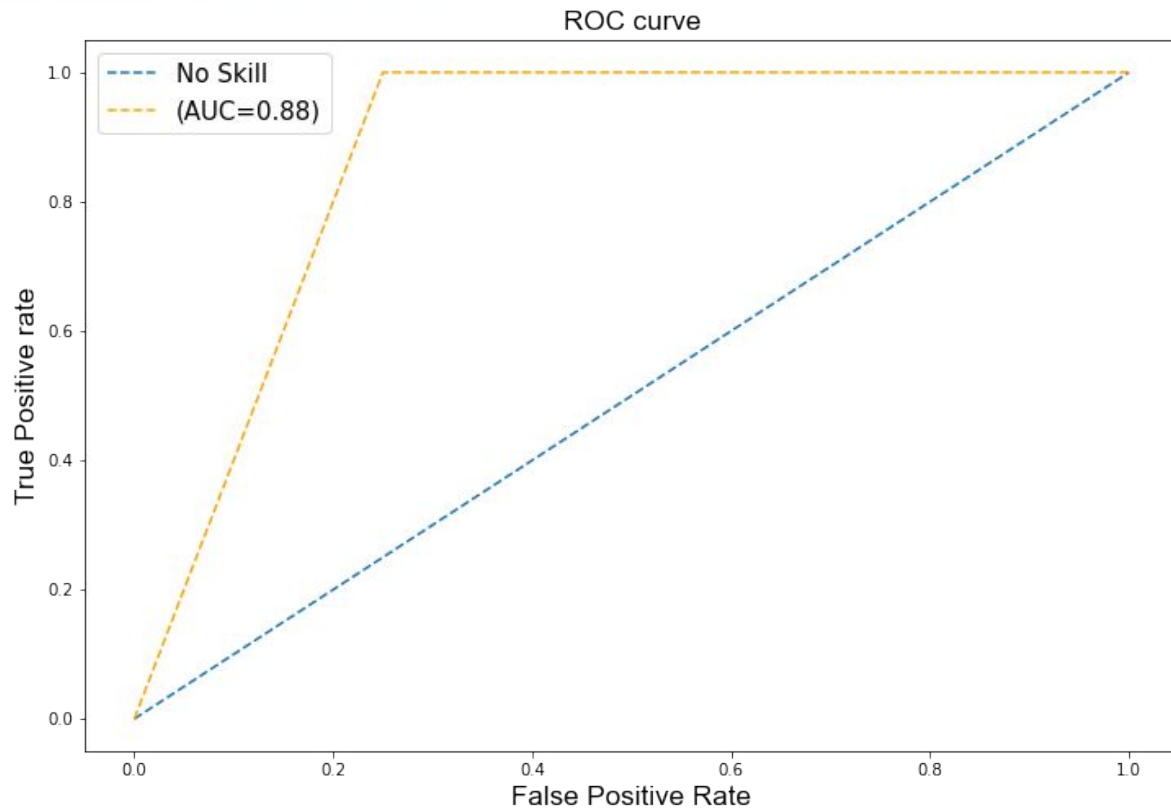
```
['hsa-mir-148b', 'hsa-mir-29b-2']
```

```
overlap_validate
```

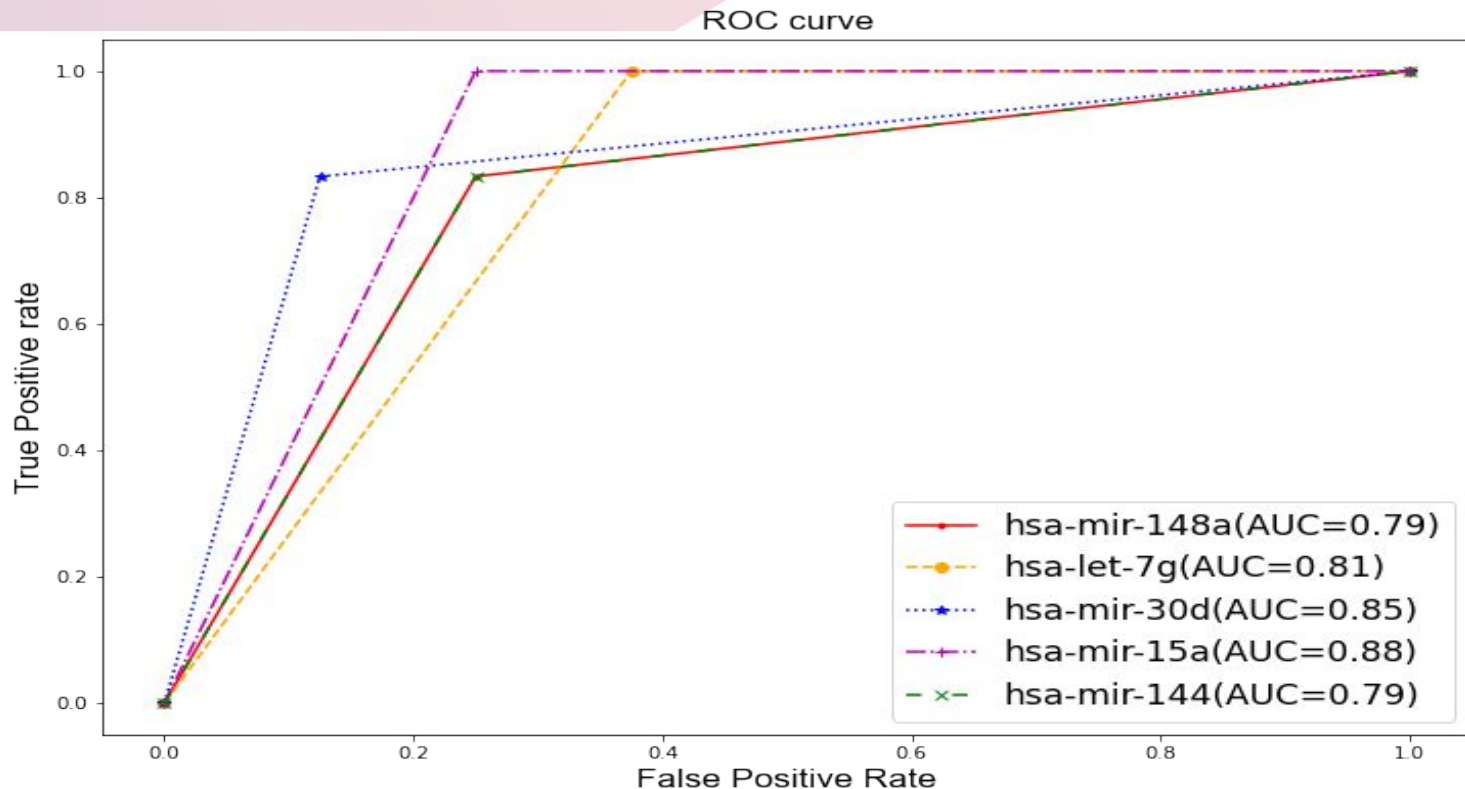
```
['hsa-mir-186',  
'hsa-mir-144',  
'hsa-mir-151a',  
'hsa-mir-98',  
'hsa-mir-148a',  
'hsa-let-7g',  
'hsa-mir-15a',  
'hsa-mir-144']
```


Evaluation of the results

- ROC curve for obtained 25 miRNAs



- ROC curves for 5 miRNAs with the highest AUC values



Evaluation of methodology



- Analyzed another dataset using the same methodology
- Used the data set under the accession number GSE147218 in NCBI database
- Obtained results were,
 - Overlap
 - Correlation Coefficient

```
['hsa-miR-193b-5p']
```

```
hsa-miR-1287-3p
```

```
hsa-let-7f-1-3p
```

```
hsa-miR-1283
```

```
hsa-miR-4703-3p
```

```
hsa-miR-155-3p
```

```
hsa-miR-1254
```

```
hsa-miR-3131
```

```
hsa-miR-26b-5p
```

Significance of our methodology



- Both statistical and machine learning approaches were used
- Feature selection was done using more than one method
 - Random Forest
 - PCA
 - Correlation coefficient values
- Higher accuracy compared to the study done using the same data set
 - Study by Leidinger et al. : 93.3%
 - Our Study : 95.24%

Deliverables Addressed in Phase 1



Milestone 01 : Background study

: Dataset selection

- Data set from NCBI database under the accession number GSE46579

Milestone 02: Preprocessing dataset (Galaxy tool)

: Data visualization and normalization

Milestone 03: Statistical analysis

- Using significance values, fold change and AUC values

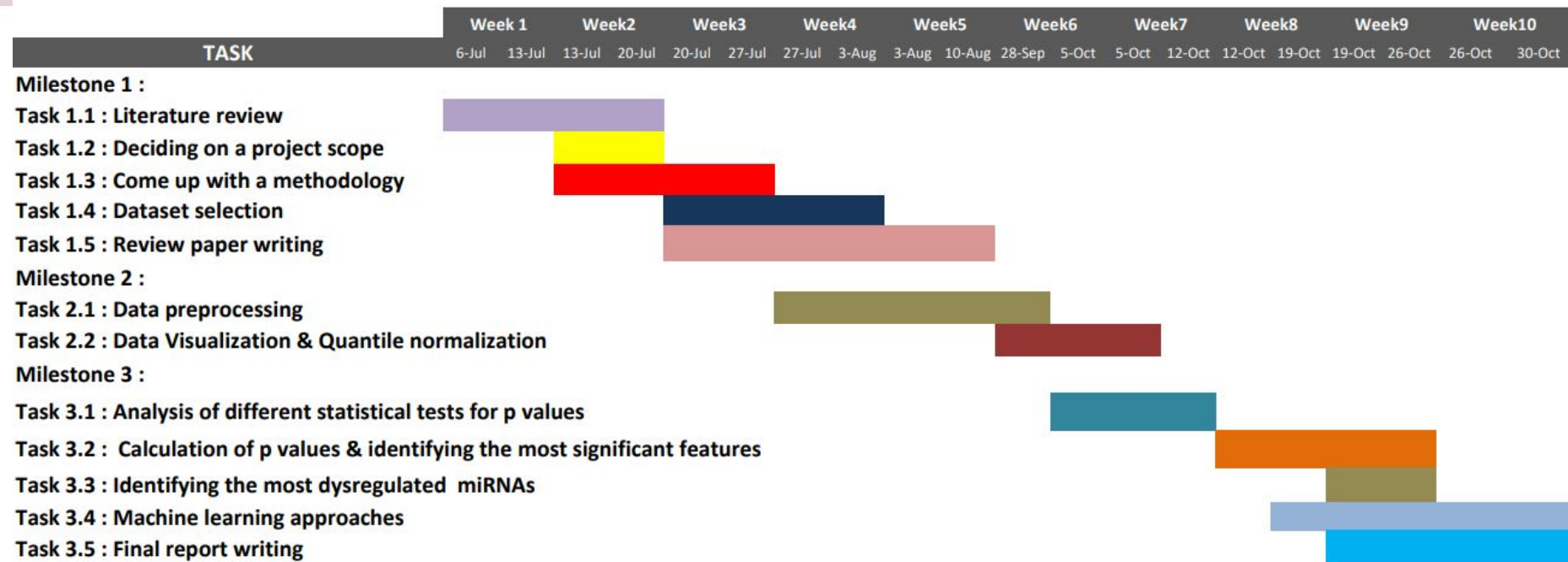
Deliverables Addressed in Phase 2



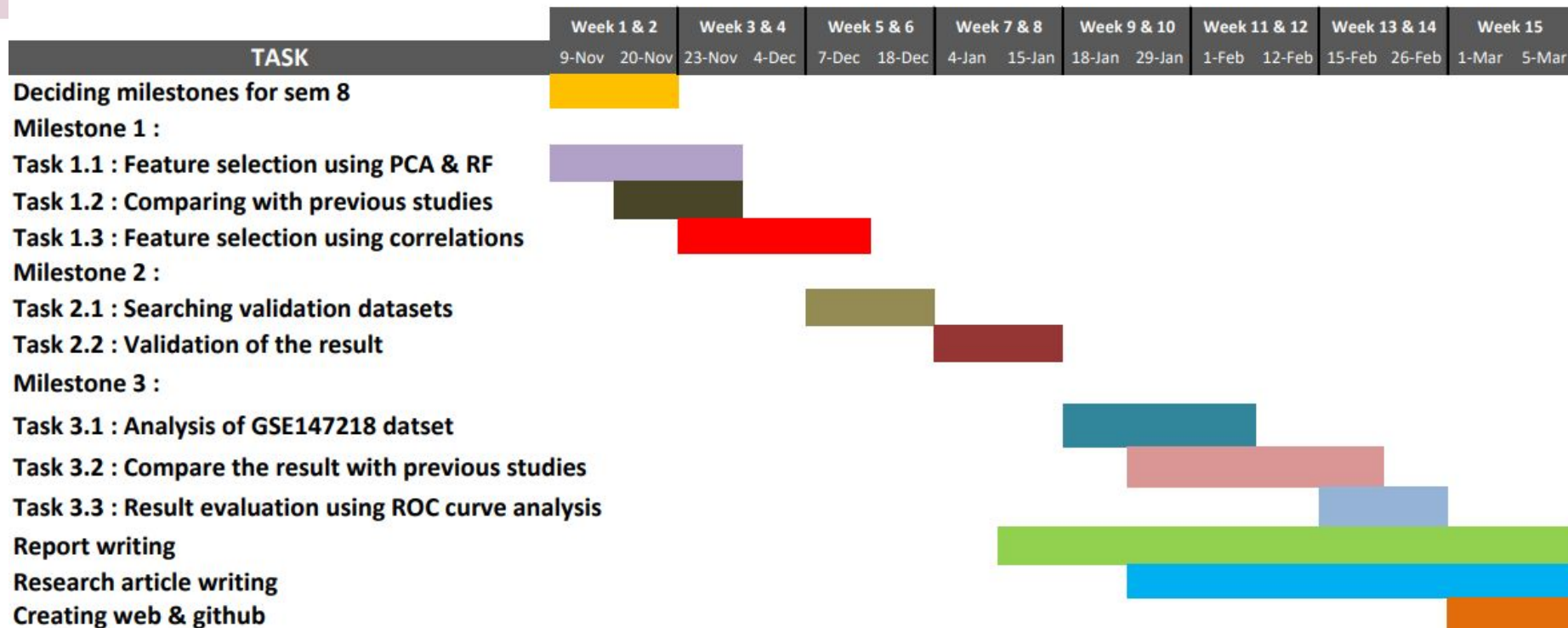
- **Milestone 01** : Identify Biomarker miRNA set
 - Using,
 - PCA
 - Random Forests
 - Correlation coefficient
- **Milestone 02** : Result Validation
 - Using HMDD v3.2
- **Milestone 03** : Analyzing another data set using the same methodology
 - Use GSE147218 dataset

: Evaluating the obtained results

Phase 1 Gantt Chart



Phase 2 Gantt Chart



Future plan



- Developing a web based tool
 - Visualizing the distributions of datasets, preprocessing, statistical analysis and classification effortlessly.



Thank you



Q & A