# README

# Document Tree

```
|   a_semi_model_final.py
|   a_train_final.py
|   a_train_final_sim.py                        # model running script
|   a_Transformer.py
|   a_utilize_semi.py
|   generate_Sigma.py                           # Sigma generating script
|   gen_Data_SIM.R                              # SimDat generating script
|   get_Folds.R
|   pack_prepare.R                              # package preparation for R
scripts
|   ReadMe_Simulation.txt
|   run_example.sh                              # Example for the whole
pipeline
|   Simdat_generation.pdf                       # Principle of simulation data
generation
|   simulation_data_functions_v4.R
|
├─.idea
|
├─data
|       embedding_selected.csv                  # embedding_selected for
simulation
|       embedding_selected.npy
|       Norms.csv
|       Sigma.csv                               # sigma generated for SimDat
genration
|       Sigma.npy
|
├─example                                       # example for embeddings and
input data
|       embedding_example.csv
|       Input_example.csv
|       input_example_sim.csv
|
├─Results
├─Results_example
|
├─Simulation
```

```
|    └─SimDat                           # Simdat generated by me
|        ├─SimDat.1                      # Simdat.1 is generated using
|        |  |  SimDat.1.csv              # All of SimDat.1
|        |  |  SimDat.1.Rds
|        |  |  SimDat.1_labeled.csv      # labeled part of SimDat.1
|        |  |  SimDat.1_label_patient_info.csv
|        |  |  SimDat.1_unlabeled.csv    # unlabeled part of SimDat.1
|        |  |  SimDat.1_unlabeled.csv6116_ALL.pkl
|        |  |  SimDat.1_unlabeled.csv9061_ALL.pkl
|        |  |  test_patients.csv         # test patients' numbers
|        |  |  train_patients_450.csv    # train patients' numbers
|        |  |  unlabeled_patients.csv    # unlabeled patients' numbers
|        |  |
|        |  ├─test
|        |  |      test_data.csv
|        |  |
|        |  └─train                      # Simdat.2 is generated using
|        |         train_data.csv
|        |
|        └─SimDat.2
```

# Input Data & Data Generation

## Input Structure

- ID: patient_num
- Y: the label 0 or 1
- T: integers that flags the dates
- other columns for the covariant(columns length is flexible)

```
"","ID","Y","T","S.1","S.2","S.3","S.4","S.5","S.6","S.7","S.8","S.9","S.10"
"100",6,0,0,0,0,0,0,0,0,0,0,0,0
"101",6,0,1,0,0,0,0,0,0,0,0,0,0
"102",6,0,2,0,0,0,0,0,0,0,0,0,0
"103",6,0,3,0,0,0,0,0,0,0,0,0,0
"104",6,0,4,0,0,0,0,0,0,0,0,0,0
"105",6,0,5,0,0,0,0,0,0,0,0,0,0
"106",6,0,6,0,0,0,0,0,0,0,0,0,0
"107",6,0,7,0,0,0,0,0,0,0,0,0,0
"108",6,0,8,0,0,0,0,0,0,0,0,0,0
"109",6,1,9,13,5,6,7,7,10,11,18,4,4
"110",6,1,10,8,4,3,4,7,12,14,11,2,4
"111",6,1,11,3,0,4,2,3,6,6,3,4,2
"112",6,1,12,0,3,2,0,3,2,10,7,2,4
"113",6,1,13,2,1,0,1,1,1,2,2,0,0
"114",6,1,14,1,1,1,2,0,0,4,2,0,1
"115",6,1,15,4,1,2,4,3,5,4,8,4,2
"116",6,1,16,3,0,4,1,2,9,9,9,4,1
"117",6,1,17,7,0,5,2,1,2,4,7,1,0
"118",6,1,18,0,0,0,0,0,0,0,0,0,0
```

see example input at: `/example/input_example_sim.csv`

# Data Generation

---

you need to download R and run `pack_prepare.R` before generating simulation data.

## Principle of simulation data generation.

See detailed description at `/Simdat_generation.pdf`

## Pipeline of data generation

---

1. **First generate Sigma for multivariant norm distribution.**

run `genrate_Sigma.py`

this python script **generate simulation embeddings** by cutting example_embedding.csv. **Then generate Sigma matrix** using dot production of different concepts' embeddings.

- selected embeddings saved at `/data/embeddings_selected.csv`
- Sigma saved at `` `/data/Sigma.csv ``

The **embedding dimension** and **covariant column number** is both set to 10 initially, change the `d` or `l` in the script if needed. Remember to **adjust the column names** at line.35 if changing the column numbers



Example of selected embeddings



Example of generated Sigma

---

2. **Generate simulation data**

run `Rscript gen_Data_SIM.R 2000 500 100 200` at **terminal**

```
Rscript gen_Data_SIM.R {1} {2} {3} {4}
# ARGS: 1.total number 2.labeled number 3.train(labeled) num 4. test(labeled) num
#Rscript gen_Data_SIM.R 2000 500 100 200
```

Generated data is saved at `/Simulation/SimDat`

- SimDat.1 & SimDat.2 are generated using different Hazard Models.
- ==Remember to adjust the line.213-line.222 of== `simulation_data_functions_v4.R` ==if changing the column numbers==

# Run Model

run `a_train_final_sim.py`

- Look into `Simulation_runmodel_descript.txt` for description of different ARGS.

# Result Structure

- Results are saved in `/Results/`, you can see example results at `/Results_Example/`
- Document Tree of Results directory:

```
C:.
|   Incident_epoch29_818__results_RETTAIN.csv      # important
|   Incident_epoch30_818__results_RETTAIN.csv      # important
|   Prevalence__818__results_RETTAIN.csv           # important
|
results_RETTAIN.csv1_Attenation_value_patient_visit_code_prediction_label_we
ight_test.pkl
|
results_RETTAIN.csv1__embedding_patient_ori_codeRW_codeVistRW_hiddenFCN_labe
l_test.pkl
|   results_RETTAIN.csv1__embedding_patient_ori_hiddenFCN_label_train.pkl
|   results_RETTAIN.csv_incident_evaluation.txt     # important
|   results_RETTAIN.csv_prevalence_evaluation.txt  # important
|   _results_RETTAIN.csv_code_weights.csv
|
└─results_RETTAIN.csv_model
    |   keras_metadata.pb
    |   saved_model.pb
    |
    ├─assets
    └─variables
            variables.data-00000-of-00001
            variables.index
```

  - Incident_epoch{}__results_RETTAIN.csv : **incident phenotyping(prediction on every time spot) results**
  - Prevalence__results_RETTAIN.csv: **Prevalence(EVER/NEVER) binary phenotyping results**
  - .txt: the calculated AUC, ACC, F1, Speci etc.