

## SECB3203 PROGRAMMING FOR BIOINFORMATICS – GROUP PROJECT (25%)

### Project Description

This project provides you with hands-on experience in utilizing Python, Microsoft Azure, or AWS cloud computing. The project is designed to align with current trends in bioinformatics and apply various bioinformatics analyses to the collected dataset.

To accomplish this, you are going to form a three-person team for the project and select a specific bioinformatics trend as the project's focal point. You will be responsible for collecting sample data for analysis. You'll have the flexibility to choose between Python, Microsoft Azure, or AWS cloud computing for data analysis, with a comprehensive report documenting your findings as a mandatory component.

Throughout the project's timeline, you are required to update your project's progress and share your work on GitHub regularly. Also, we recommend looking for a mentor or client (who can be your academic advisor or any other related professional) to enhance the project's depth and guidance.

Note:

- Create a folder named your project title at [https://github.com/sean-seah/SECB3203\\_25261](https://github.com/sean-seah/SECB3203_25261) then the README.md of the folder must include showing your name and your member's name.
- Don't rely on your member to finish everything alone.

### Examples of Trends in Bioinformatics

- Focus
  - Discovery of biomarkers: Schizophrenia disease, cancer (breast/lung/kidney/stomach/skin)
  - Classification of diseases: Alzheimer's disease, Parkinson's disease, cancer (breast/lung/kidney/stomach/skin)
  - Prediction of cancer prognosis: cancer (breast/lung/kidney/stomach/skin)
  - Protein structural class classification
  - Protein structure prediction
  - Gene-disease predictions
- Bioinformatics analysis/approach
  - Gene selection/feature selection/variable selection
  - Deep learning: Convolutional neural network
  - Network analysis/Pathway analysis/Gene expression analysis
- Data
  - Image/Video: GitHub/KAGGLE/research papers
  - Gene expression data: GEO
  - Pathways: KEGG
  - Protein-protein interaction network: STRING

**Project Activities and Due Dates**

<b>Task</b>		<b>Marks (%)</b>	<b>Start Date</b>	<b>Due Date</b>	<b>Duration</b>
<b>1</b>	<b>Project Proposal and Group Formation.</b>  - Students in groups of 3 - Choose a trend in bioinformatics <u>Proposal</u> 1.0 Introduction 1.1 Problem Background 1.2 Problem Statement 1.3 Objectives 1.4 Scopes 1.5 Conclusion	2.5	Week 2 15 Oct	Week 4 29 Oct, 0000 MYT submit to the e-learning	7 days
<b>2</b>	<b>Project Progress 1.</b>  - Software and hardware requirements - Flowchart of the proposed approach	2.5	Week 5 5 Nov	Week 9 3 Dec, 0000 MYT submit to the e-learning	14 days
<b>3</b>	<b>Project Progress 2.</b>  - <u>Importing Dataset</u> - Understanding the data - Importing and exporting data in Python - Getting started analyzing data in Python - Python packages for Data Science  - <u>Data Wrangling (Pandas/Numpy)</u> - Identifying and handling missing values - Data formatting - Data normalization (centering/scaling) - Binning - Indicator variables	5.0			
<b>4</b>	<b>Project Progress 3</b>  - <u>Exploratory Data Analysis (Pandas/Numpy)</u>	2.5	Week 7 19 Nov	Week 11 17 Dec, 0000 MYT	10 days

## PROJECT

2023

Task		Marks (%)	Start Date	Due Date	Duration
	<ul style="list-style-type: none"> <li>- Descriptive statistics</li> <li>- Basic of grouping</li> <li>- ANOVA</li> <li>- Correlation</li> </ul>			submit to the e-learning	
5	<p><b>Project Progress 4</b></p> <ul style="list-style-type: none"> <li>- <u>Model Development (Pandas/Keras/TensorFlow/Matplotlib/...)</u></li> <li>- Simple and multiple linear regression</li> <li>- Model evaluation using visualization</li> <li>- Polynomial regression and pipelines</li> <li>- R-squared and MSE for In-Sample Evaluation</li> <li>- Prediction and decision making</li> </ul> <p><b>Note:</b> make sure to submit the codes and flowchart for model development</p>	5.0	Week 9 3 Dec	Week 12 24 Dec, 0000 MYT submit to the e-learning	20 days
6	<p><b>Project Progress 5</b></p> <ul style="list-style-type: none"> <li>- <u>Model evaluation (sklearn)</u></li> <li>- Model evaluation</li> <li>- Over-fitting, under-fitting, and model selection</li> <li>- Ridge regression</li> <li>- Grid search</li> <li>- Model refinement</li> </ul>	5.0	Week 11 17 Dec	Week 14 7 Jan, 0000 MYT submit to the e-learning	20 days
7	<p><b>Project Assessment and Demo</b></p> <ul style="list-style-type: none"> <li>- <u>Live presentation/Recorded video (5 minutes)</u> <ul style="list-style-type: none"> <li>- Individual and group assessment</li> <li>- shared the recorded video and your GitHub link with your client and put your client's comments in the report (section 4.0)</li> </ul> </li> <li>- <u>Report (follow the contents below)</u> <ul style="list-style-type: none"> <li>1.0 Introduction</li> <li>1.1 Problem Background</li> <li>1.2 Problem Statement</li> <li>1.3 Objectives</li> </ul> </li> </ul>	7.5	Week 12 24 Dec	Week 14 7 Jan, 0000 MYT submit to the e-learning submit to the github submit to the youtube	

## PROJECT

2023

Task	Marks (%)	Start Date	Due Date	Duration
1.4 Scopes 2.0 Data collection and pre-processing 2.1 Importing dataset 2.2 Data wrangling 2.3 Software and hardware requirements 3.0 Flowchart of the proposed approach 3.1 Exploratory data analysis 3.2 Model development 3.3 Model evaluation 4.0 Testing and validation 5.0 Conclusions				
<p>Note:</p> <ul style="list-style-type: none"><li>• Take photos of your meetings (including meetings with your client) and put them in GitHub to monitor your project progress.</li><li>• All progress must be tallied between GitHub and what you submitted to the e-learning.</li><li>• During e-learning submission, you just copy-paste or print in pdf what you updated in the GitHub, and make sure the submitted file is pdf.</li></ul>				