### BMEG 3105: Data analytics for personalized genomics and precision medicine

Lecturer: Yu LI (CSE)

# Email:liyu@cse.cuhk.edu.hk

Scriber: HAN Mingyu(1155213539)

#### Lecture 1: Course Introduction

# Agenda

- Review the pre-course survey results
- Course logistics
- Brief overview of DATA in personalized genomics and precision medicine

# Pre-course Survey Results:

Response: 28/46 (Mostly BME UG students from U2-U4)

Drive of taking the course:	Help maybe needed for:
1. For the biological/ genomics/ health	1. Implementation the methods with
applications	programming
2. For the data analytics/ machine learning	2. Understand the mathematical background of
techniques	the data analytical methods
3. For research project experience	3. Understand concepts in data analytics
4. For the three credits and degree requirement	4. Apply the methods learned from the course
5. Just exploring a new field	to new data
	5. Finding additional materials and resources
	to learn more about the concepts
Interested topics:	What would we provide?
1. Neural networks	1. basic Python programming tutorials
2. Cancer genomics	2. not touch too much math
3.Data visualisation	3. fundamental concept
4. Protein-protein / RNA interaction	4. learn to apply the techniques into the data
5. Convolution neural networks	5. practical usage and project topics
	6. Additional resources and materials will be

# provided

ESTR3605: ESTR mirror of BMEG3105. Additional 45-min lecture each week

Minimum number of students to open the course:5.

#### Course Logistics

Lectures: Wed 9:30am-11:15am (11:05am), SC L4 Fri 9:30am-10:15am, MMW703

Tutorial: Fri 10:30-11:15am, MMW703

- Slides will be available the day before the lecture day
- there is no video recording, as required by the school
- TA will have several sessions to help you on Python programming
- All the materials will be available on https://lim-zq.github.io/BMEG3105-Fall-2025/

#### Teaching team:

Yu Li (Instructor)

liyu@cse.cuhk.edu.hk

Office hour: 3-5pm, Friday

Location: SHB-106

#### ZiqianLin(TA)

linziqian@link.cuhk.edu.hk

Office hour:3pm-5pm,Tuesday

Location: SHB-904

Xinyuan Liu(TA)

1155246738@link.cuhk.edu.hk

Office hour:3pm-5pm,Thursday

Location: SHB-904

#### Software and communications:

Blackboard	Piazza
	(https://piazza.com/cuhk.edu.hk/fall2025/bmeg
	3105)
The main software to manage the course	You can ask questions through Piazza
Grading will be through Blackboard	For a personal matter, please use the private
	post to the instructor and the TA

#### **Grading:**

- ♦ Homework (20%): Three grading homework(5%+5%+5%) and one non-grading programming assignment (5%, make sure you learn something from it)
- Scribing (10%): Grading scribing. Summarize one of the lectures. Submit it one week after the course. Each student should do at least one lecture. Notice that your note and scribing will be posted online, for others reference. You can choose to remove your name or not. You can sign for at most two, for additional 1%
- ♦ In-class quiz (10%): Two in-class quizzes. The questions will be simple. Mainly for checking the participation. The exact dates are on the website: Oct 15, and Nov26
- ♦ Midterm (20%): A grading midterm exam. One bonus question (2%)
- ❖ Project (20%): A grading project. You can give us your project and seek our help, or we will predefine some projects for you to choose (You should submit a mid-term report (5%), a final report (7%) + presentation (3%) together with the implementation (5%).)
- ❖ Final (20%): A grading final

#### All the exams or quizzes will be open book!

Programming: Non-grading assignment (5%), Grading programming included in the project (5%)

#### Bonus:

One bonus question in Midterm (2%)

One additional scribing: 1%

Pre-course survey + Post-lecture survey: 0.5% for each, and the maximum is 3%.

The bonus is sufficient to cover all the programming credit.

#### Attendance:

do not check the attendance except for the two participation quizzes as well as the mid-term and final

#### Programming:

Language: Python

#### Assignment:

	<u>Posted</u>	<u>Due</u>
Programming Assignment 0: Programming environment setup	<u>Sep 5</u>	<u>Sep 17</u>
Assignment 1: About the basic concept of data analytics-1	Sep 12	Sep 24
Assignment 2: About the basic concept of data analytics-2		Oct 15
Programming Assignment 1: About application of DA to the biology Oct		Nov 14
Assignment 3: DA in Personalized Genomics and Precision Medicine	Nov 12	Nov 21

Do be serious to the assignments. They can be very helpful for you to prepare for the mid-term and final.

#### Scribing:

Summarize one of the lectures, Submit it within one week after the lecture

Grades will be deducted by 25% for each additional late day

Lecture date: 5 Sep. Deadline: 12 Sep, 11:59pm

First two lectures: additional one week for submission

#### Midterm Date:

Oct 17 (Fri) In class, open-book

#### <u>Project</u>

Project milestone report(Proposal): 1-page report (Due: Nov 7)

- 1. Title, author
- 2. What problem do you want to do? Why is the problem interesting? (1%)
- 3. What data are you going to process? The source, the size, the sample of the data(1%)
- 4. What's the output of your method? (1%)
- 5. How are you going to do it? Describe the method step by step, from input to output(1%)
- 6. What are the expected results? How are you going to evaluate the results? (1%)
- 7. What have you done?

Project report: No length requirement, submit the report together with codes (5%, whether it is correct or not) (Due: Dec2)

- 1. Title, author
- 2. What problem do you want to do? Why is the problem interesting and important?(0.5%)
- 3. What data have you processed? The source, the size, the sample of the data (0.5%)
- 4. What have you done to resolve the problem? Describe the method step by step, from input to output (2%)
- 5. What are the results? (1.5%)
- 6. Result evaluation (1.5%)
- 7. Any idea for further improvement? (1%)

Project presentation (Date: Nov 21, Nov 28):7 mins for each student

Will be graded in the following way:

1.Logic (1%)

What is the problem?

Why is it important?

How do you resolve it?

The overview of your idea

The overview of the results

2.Clarity (1%)

Whether the audience can understand and follow the presentation Slides

3.Preparation (1%)

Clear illustration No typos, no grammar errors

#### Late days

Can be used on Assignment 1,2,3 and programming assignment, project mid-term report

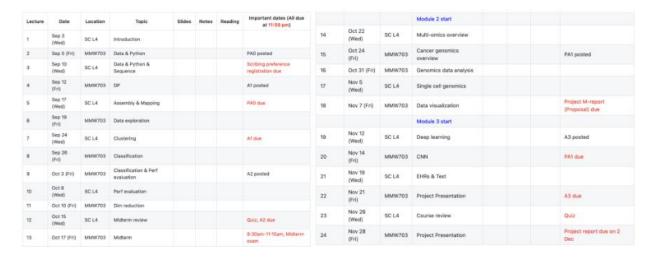
Cannot be used on final project report and scribing report

6 late days total, 2 max for assignments

Grade will be deducted by 25% for each additional late day

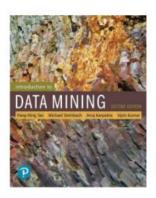
Let TA know when it is necessary to use the late dates

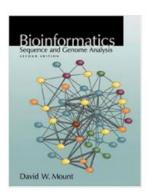
#### Course schedule

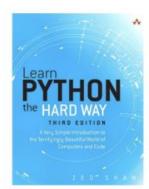


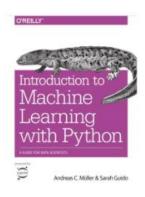
#### Reference book

#### Module 1









Module 2 + Module 3

Cancer genomics, Single cell, Multi-omics, EHRs, Protein-RNA interaction are all cutting-edge research topics. There is no such a book covering all of them.

provide the reading materials and reference books when we are there.

#### Overview of DATA in personalized genomics and precision medicine

#### 1. Why data analytics?

1. Massive of data is being collected and warehoused

1.1 Web data: Meta, Google, Amazon, X, TikTok

1.2 Biological data: DNA sequences, protein sequences

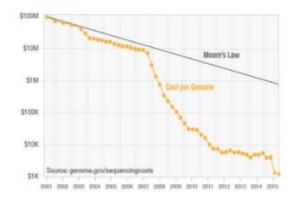
1.3 Bank/credit card transaction data: Alipay, PayPal

- 1.4 Mobile data: China Mobile, CSL
- 2. Computers have become cheaper and more powerful
- 3. Data analytics are useful
- 3.1 Aggregate data
- 3.2 Generate hypothesis
- 3.3 Support the conclusion

# 2. Why data analytics in personalized genomics and precision medicine?

Tons of sequencing and health data available and waiting to be analyzed

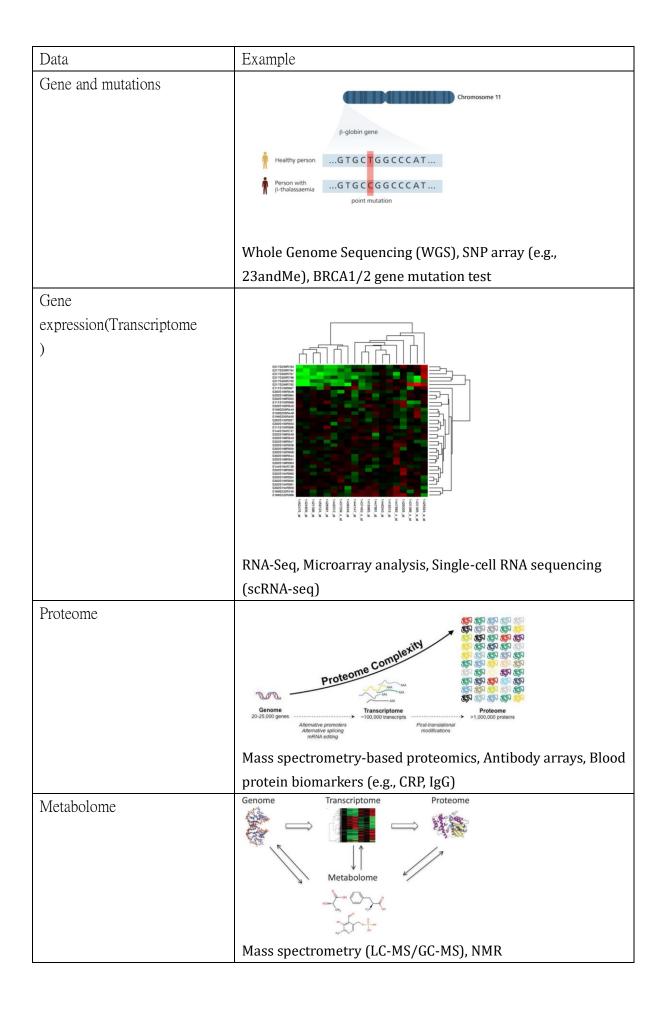
#### Sequencing cost decreasing dramatically



#### Global efforts in building biobank

Location	Biobank	N (goal)
Canada	CARTaGENE biobank <sup>118</sup>	43,000
USA	All of Us <sup>23</sup> Million Veteran Program <sup>49</sup>	1,000,000 > 600,000
Mexico	The Mexico City Prospective Study <sup>50</sup>	150,000
iceland	deCODE Genetics	500,000
UK	UK Biobank <sup>36</sup> Avon Longitudinal Study of Parents and Children (ALSPAC) <sup>30</sup>	500,000 > 15,000
Netherlands	Lifelines Biobank <sup>100</sup>	> 167.000
Denmark	Danish National Biobank 121	
Norway	HUNT - Nord-Trendelag Health Study 122	125,000
Sweden	Biobank Sweden	
Finland	FinnGen	500,000
stonia Estonian Biobank <sup>123</sup>		52,000
Israel	Project 10K	10,000
Saudi Arabia	Saudi Biobank	200,000
Qatar	Qatar Blobank <sup>124</sup>	
China China Kadoorie Biobank <sup>51</sup> Guangzhou Biobank <sup>101</sup>		> 500,000 30,000
Japan	BioBank Japan <sup>136</sup>	200,000
Korea	National Biobank of Korea 127	500,000
Talwan	Taiwan Biobank <sup>128</sup>	200,000

3. What data can we have to measure a person?



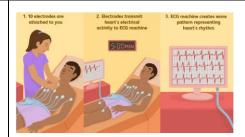
	spectroscopy, Blood glucose, cholesterol, amino acids levels
Molecular network & cellular network	Metabolite Gene/Protein Metabolism Protein/gene interaction Signaling Regulation
	Protein-protein interaction maps, Metabolic pathways (e.g.,
	KEGG), Gene regulatory networks
Microbiome(Oral and gut )	
	16S rRNA sequencing (to identify gut bacteria composition),
Organ(Biomedical imaging )	Stool test for pathogens  MRI, CT Scan, X-Ray, Ultrasound, Echocardiogram

Hospital test(Blood test and so
on)

Blood Test	Result	Normal Value
WBCs (billion/L)	8.00	3.5 to 10.5
Neutrophils (%)	62	40 to 70
Lymphocytes (%)	28	25 to 45
Monocytes (%)	10	2 to 8
Eosinophils (%)	1	1 to 5
Basophils (%)	0	0 to 1
RBCs (trillion/L)	3.84	4.3 to 5.7
Hb (g/dL)	11.7	13 to 17
Hematocrit (%)	37	37 to 52
Platelets (billion/L)	262	150 to 450

# MRI, CT Scan, X-Ray, Ultrasound, Echocardiogram

# Electrocardiography



# Resting ECG

Demographic
information(Age,
gender, location and so on)

Demographic Categories	Frequency	Valid Percentage	U.S. National Census Data (2012), 5
Gender			
Female	150	49.5	51.1
Male	153	50.5	48.9
Age			
18-24	24	9.5	11.2
25-34	53	20.9	13.4
35-44	36	14.2	12.9
45-54	67	26.5	14.2
55-64	61	24.1	12.3
65+	12	4.7	13.5
Not specified	50	_	_
Ethnicity			
Euro-American/Caucasian	241	79.5	63.0
African American	26	8.6	13.1
Hispanic/Latino(a)	17	5.6	16.9
Asian American	11	3.6	5.1
Other	8	2.6	1.9
Marital status			
Married	168	55.4	56.4
Never married	70	23.1	26.9
Divorced/separated/widowed	68	21.5	16.7
Employment			
Full time	144	47.5	64.7
Part time	45	14.9	25.7
Unemployed/retired	114	37.6	9.6
Highest level of education			
4-year college or graduate degree	140	46.2	28.5
High school degree or other	163	53.8	57.2
Annual household income			
Less than US\$20,000	39	12.9	
US\$20,000-US\$39,999	84	27.7	
US\$40,000-US\$59,999	63	20.8	
US\$60,000-US\$79,999	43	14.2	N/A*
US\$80,000-US\$99,999	33	10.9	100,700,0
US\$100,000 or up	41	13.5	
Total	303	100	

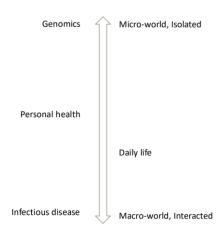
Age, Gender, Postal code, Ethnicity, Education level, Occupation

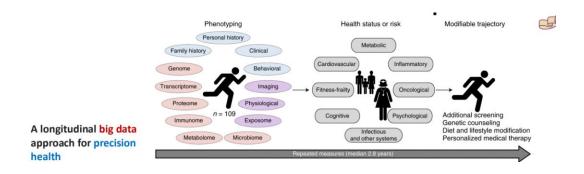
# Personal statement and doctor diagnosis



	diagnoses (e.g., ICD-10 codes: I10 for hypertension)
Living habit(Exercise)	Wearable data (steps, active minutes, heart rate from Fitbit/Apple Watch), Self-reported workout logs
Diet	Food diary apps (MyFitnessPal), Nutritional intake analysis (calories, macros, micronutrients)
Family history	History of heart disease, cancer, or diabetes in first-degree

	relatives
Communication and social media data	Language sentiment analysis, social network size, frequency of interactions (for mental health research)
Environment(Pollution)	PM2.5 exposure level (from local air quality monitors), Lead level in drinking water
Travel history(Global pandemic)	Travel itinerary during a global pandemic



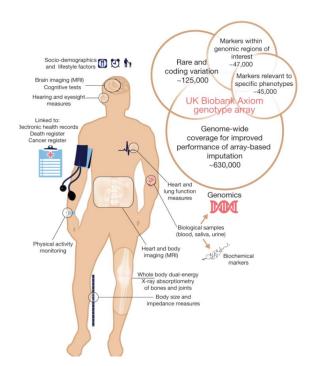


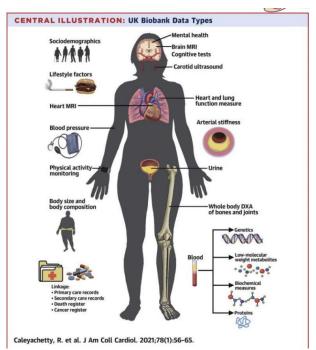
Goal: Use long-term, frequent data tracking for personalized health.

Method: Tracked 109 people for 4-8 years, collecting a massive amount of data.

Data Collected: Everything from genes and proteins to blood tests, wearables, and lifestyle questionnaires.

Purpose: To build a complete health profile and assess risks for diseases like diabetes, heart disease, and cancer. Enable personalized actions like diet plans, genetic counseling, and custom medical treatment.





The UK Biobank is a major biomedical database with deep genetic and health data from 500,000 people. It combines genomics, imaging (MRI, DEXA), physical measures, and long-term health records to study disease causes and prevention.

#### Learning outcome

- 1.Learn thefundamental conceptof data analytics
- 2.Know the various data ingenomicsandmedicine
- 3. Apply the data analytics techniques toprocess the data and resolveproblems in biology