



COMP20008 Elements of Data Processing Phase 4 Presentation

Semester 1 2018

**Name: Xulin Yang
Student id: 904904**

Structure of the talk:

- 1 research questions
- 2 why questions are valuable?
- 3 data sets used & why
- 4 data preprocessing methods used
- 5 findings
- 6 challenges



- **Questions:**
- (1) Identify the correlation between community citizens' chronic disease risks and citizens' health risks.
- (2) What are features in communities that causes citizens to have high risks?
- **outcomes**



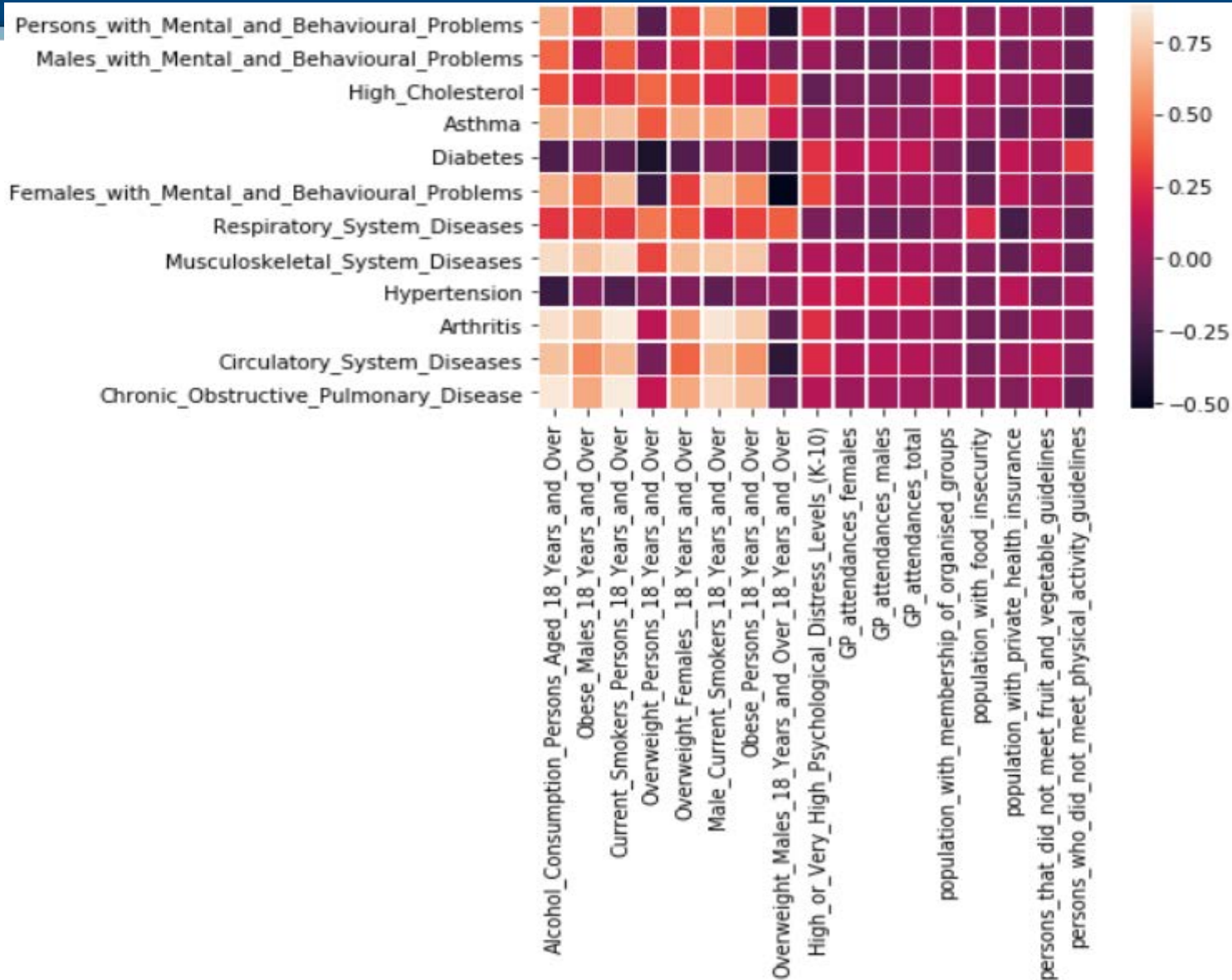
- Purpose
- Who benefits & interests?
- Innovation

- LGA11 Chronic Disease Modeled Estimate
- LGA11 Health Risk Factors Modeled Estimate
- LGA11 Psychological Distress Modeled Estimate
- Local Government Area LGA profiles data 2011

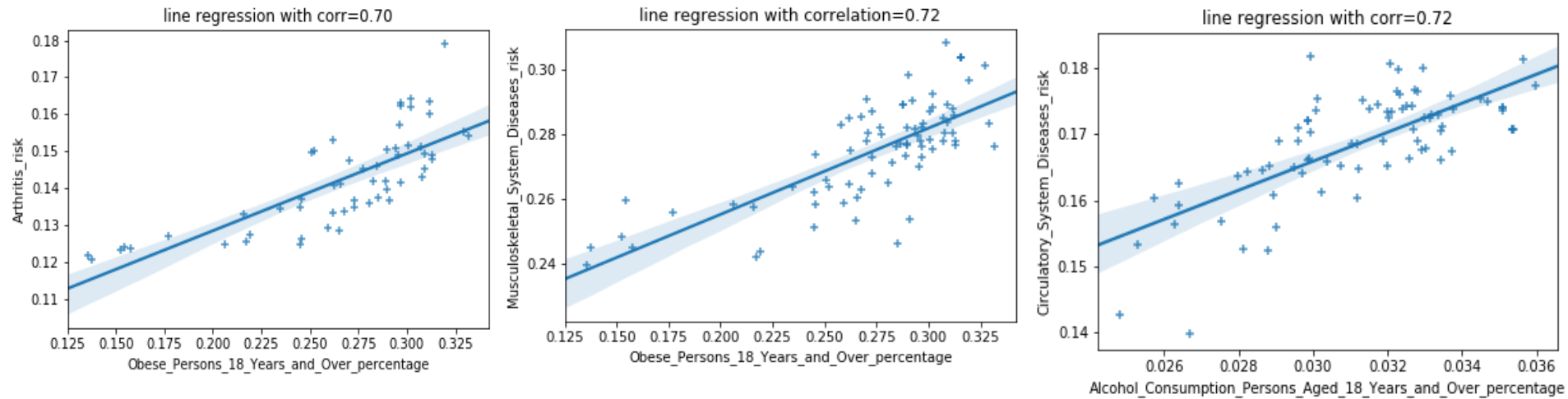
	Data_set	#Rows	3 example columns	Column data type	Reliable data%
1	LGA11 Chronic Disease Modeled Estimate	80	"Arthritis"	Discrete	70%
2			"Arthritis_RRMSE"	Qualitative	
3			"Local_Government_Area_Code"	Categorical	
4	LGA11 Health Risk Factors Modeled Estimate	80	"Current_Smokers"	Discrete	84%
5			"Current_Smokers_RRMSE"	Qualitative	
6			"Local_Government_Area_Code"	Categorical	
7	LGA11 Psychological Distress Modeled Estimate	80	"People_felt_high_mental_pressure"	Discrete	97.50%
8			"People_felt_high_mental_pressure_RRMSE"	Qualitative	
9			"Local_Government_Area_Code"	Categorical	
10	Local Government Area LGA profiles data 2011	79	"Unemployment_rate"	Discrete	70%
11			"People_lack_physical_activity"	Discrete	
12			"Local_Government_Area_Code"	Categorical	
13					

- Wrong data corrected
- Unreliable data dropped
- Missing data
- Columns replaced by schema
- Data normalization
- Deriving data from raw data

Finally have three DataFrame: disease_risks,
health_risks, lga_features

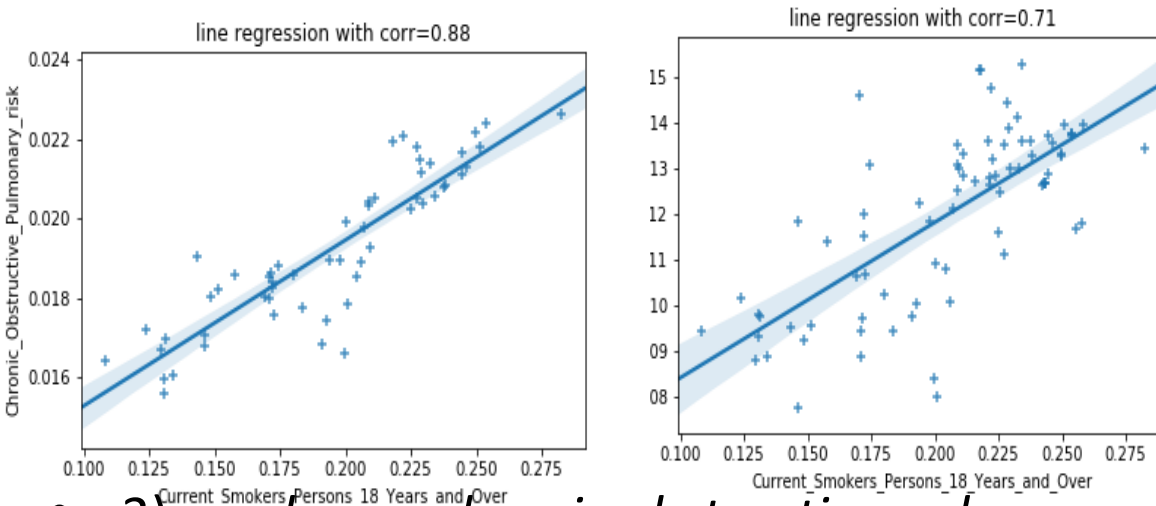


heat map against disease risk and health risk



1) obese v.s. arthritis risk 2) obese v.s. musculoskeletal system diseases risk

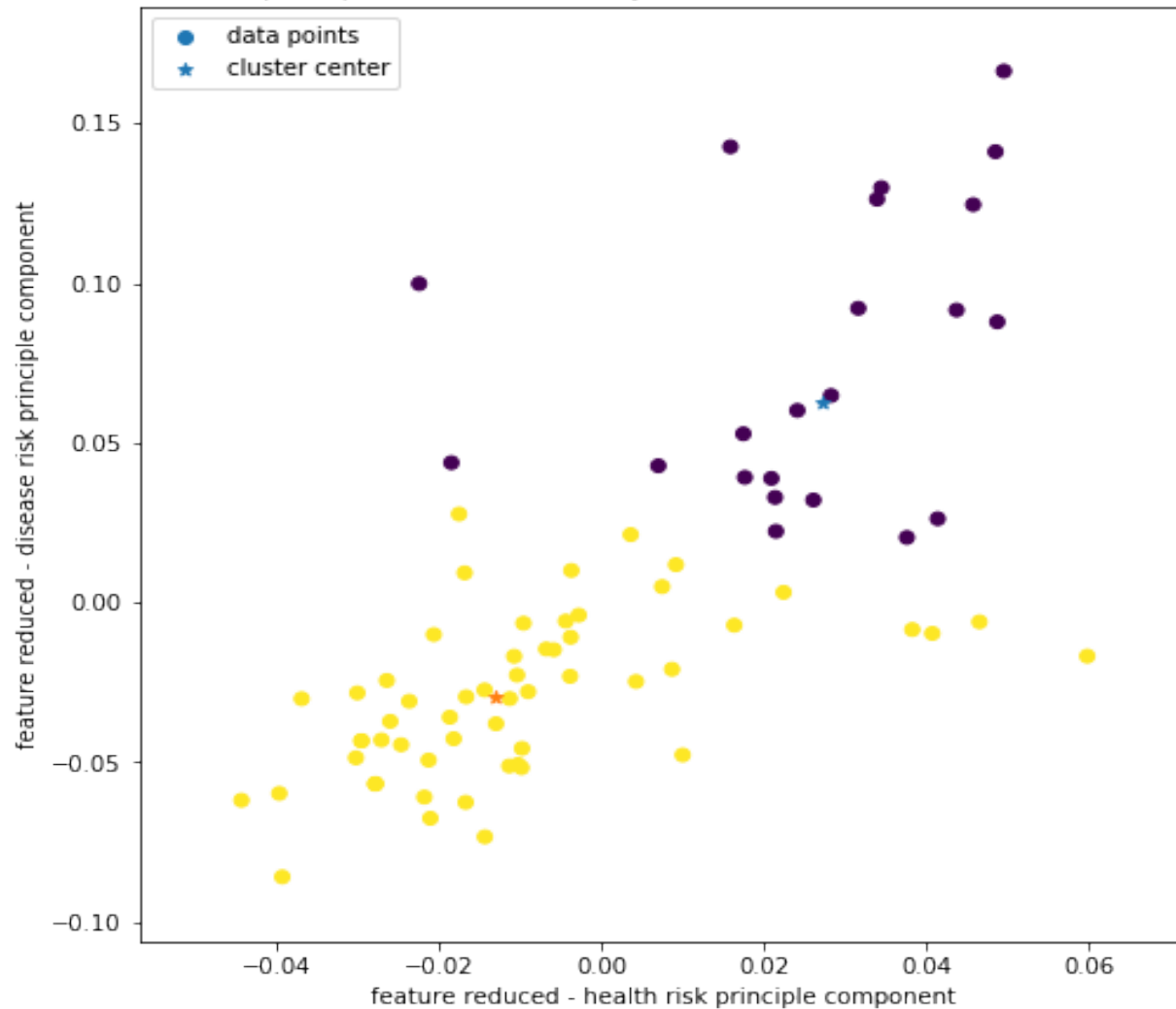
5) alcohol v.s. circulatory system disease



• 3) smoke v.s. chronic obstructive pulmonary risk 4) smoke v.s. asthma

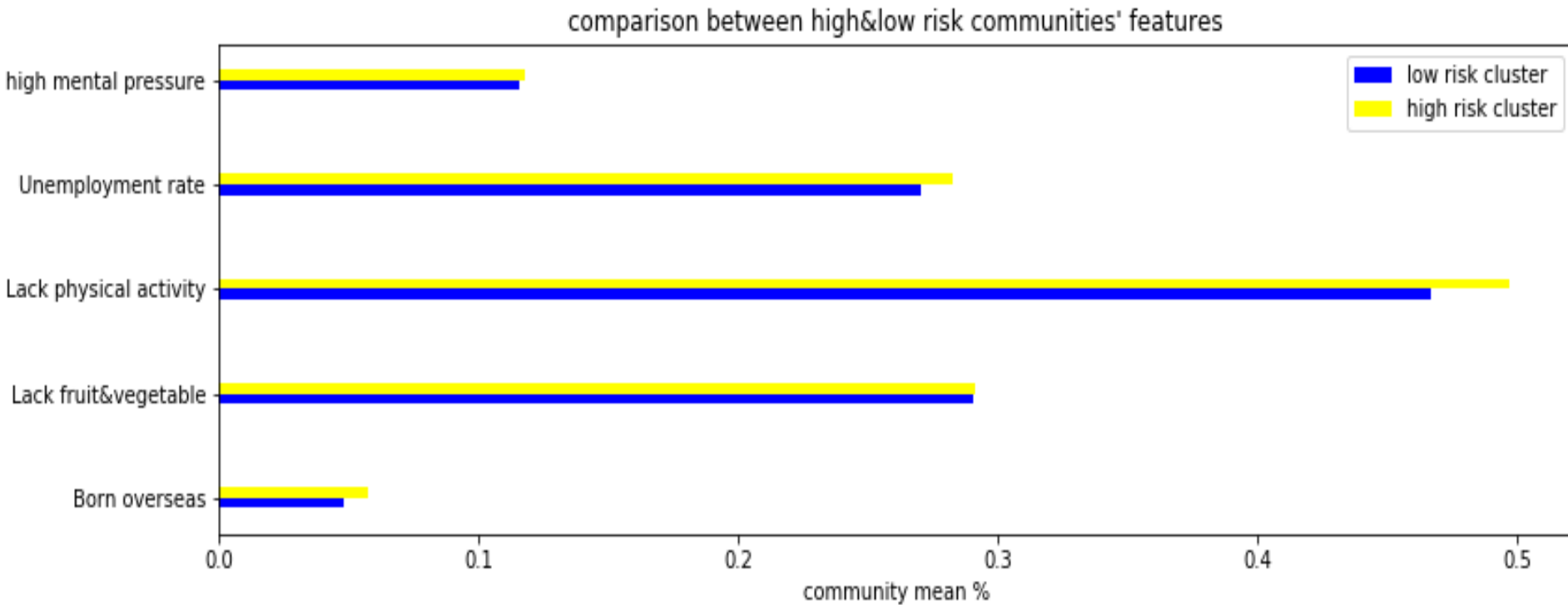


principle coordinates analysis with k-mean cluster (k=2)



The sixteen communities
with high risks:
Moreland,
Kingston,
etc

principle coordinated analysis with K-mean (k=2) clustering



community feature for clustered data



- circulatory system disease <> drinking alcohol
- musculoskeletal system diseases & arthritis <> obese adults
- Asthma & chronic obstructive pulmonary disease <> smokers
- high people felt high mental pressure,
- high unemployment rate
- High rate of lack of exercises
- High immigrants' population percentages

Features of
Community with

High risk

<> : correlated

- Appropriate data sets to use
- Selecting columns(100+ columns in lga census data set) also consumes a lot of time



Thanks for your attentions



- Any questions and queries?