# A Network Tour Of Data Science Final Project

# Course Suggester

Cherqui Alexandre

Gusset Frédérick

El Hamamsy Laila

Triquet Thomas

25.01.2018

## What is needed ?

Enrolments*

- List of all courses taken by all students between 2007 and 2016

- Standardized data, minimal cleaning required

Course Descriptions*

- Enrich the data with specific course attributes

- Manually input, fastidious to standardize.

*Courtesy of Francesco Pinto, EPFL Data Scientist

# Baseline : Enrolments

- Weight between two courses proportional to number of students who took both


Unweighted Degree Distribution


Weight Matrix


Spectral Clustering

# Baseline : Enrolments

- Weight between two courses proportional to number of students who took both



Normalized Laplacian Eigenvectors
Architects are not engineers

# Enriching : Study Plans

- Redundant with the enrolments as students tend to stick to their study plans



AR

Normalized Laplacian Eigenvectors – All Sections
Projection on the 2nd and 3rd eigenvectors

# Enriching : Study Plans

- Redundant with the enrolments as students tend to stick to their study plans



Normalized Laplacian Eigenvectors – All Sections
Spring View

# Baseline : Enrolments

- Removing Architects

- Spectral clustering on first eigenvectors of the normalized Laplacian

- Distinction between faculties

- Distinction between sections

# Baseline : Enrolments



- Decision to focus on the STI faculty for Master courses

- Objective : propose a diverse ensemble of courses by enriching this graph



Normalized Laplacian Eigenvectors – STI Faculty Isolating Sections

# Enriching : Professors

- Sparse



Unweighted degree Distribution



Distribution of Component Sizes



Weight Matrix

- Largest giant component 16 due to courses with multiple professors
- Highest node degree 8

# Enriching : Assistants

- Even more sparse

- Assistants are not necessarily specified on the course descriptions



Unweighted degree Distribution



Distribution of Component Sizes

- Largest giant component (4) : due to courses with multiple assistants

- Second largest component : 3

# Enriching : Topics

- Keywords, contents, summaries, etc...

- Topic detection to create links between courses

- Matrix factorization to get a Course x Topic matrix

- Graph constructed based on distance between courses



```
Topics in LDA model:
Topic #0: direct diode dimensional digital different differential difference diagnostic
 device development determine develop detection detail described design description dept
h describe depend
Topic #1: direct diode dimensional digital different differential difference diagnostic
 device development determine develop detection detail described design description dept
h describe depend
Topic #2: direct diode dimensional digital different differential difference diagnostic
 device development determine develop detection detail described design description dept
h describe depend
Topic #3: 6ghz aberration ability able abstract acoustic acquaint acquire acquires acqui
sition active act activity actuator add addition address adaptive advance afm
```

# Enriching : Requirements



- 3 graphs:
  - Link a course to its requirement.
  - Link the requirements of the same course.
  - Link the courses that have the same requirements.



Unweighted degree distribution



weight matrix

- Courses with same requirements : not enough connections, courses recommended are irrelevant although in the same section.

- Requirements of same course : immediate neighbours are pertinent

- Course to its requirements : similar results as the requirements of the same course but with less information.

# Enriching : Final Graph

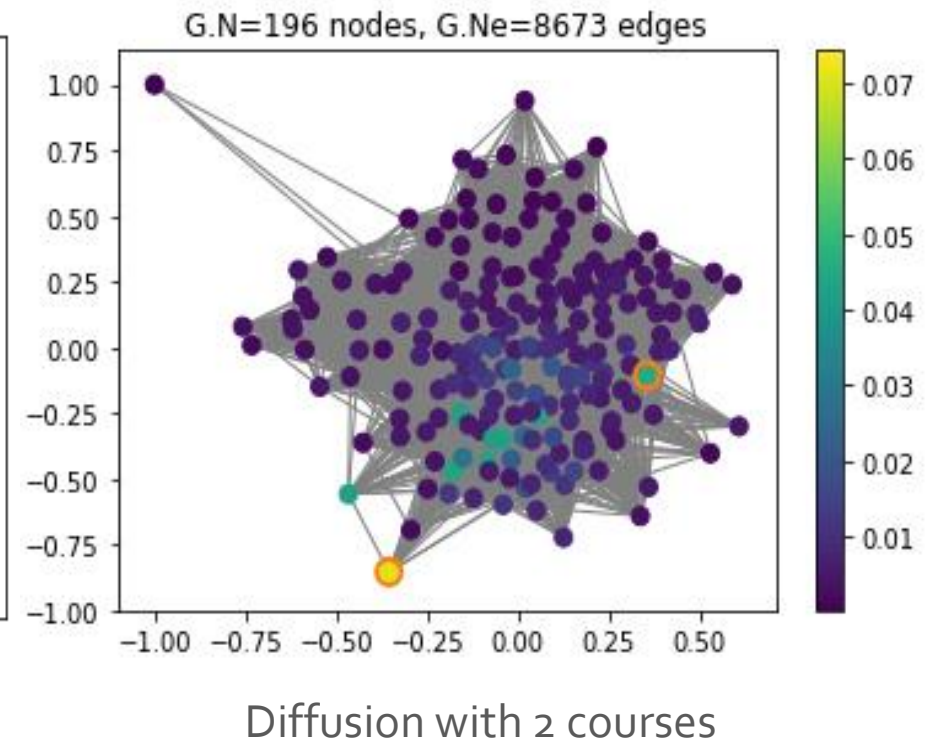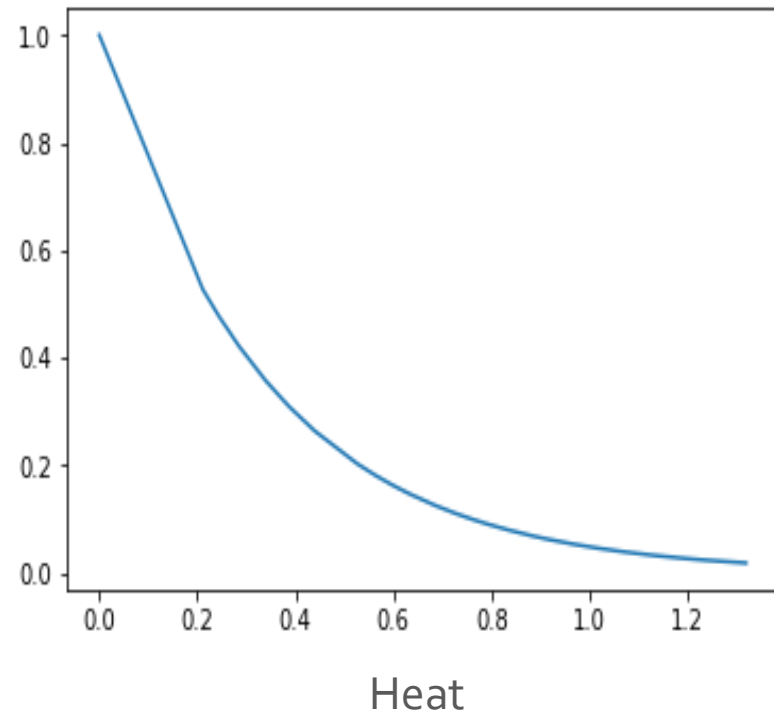| Data | Status | Reason |
| --- | --- | --- |
| Students | Kept | Baseline |
| Study Plans | Discarded | Same information as the enrolments |
| Requirements | Kept | Links based on common topics of the courses |
| Keywords | Discarded | Interesting but needed more information |
| Professors | Discarded | Too sparse |
| Assistants | Discarded | Too sparse |

- Weighted sum of the requirements and the baseline to obtain the final graph
- Empirical weights : not possible to do supervised learning
  - Using the probability of taking two courses would overfit the baseline graph
  - Did not have the data required to do a good metric

# Graph Diffusion

- Low pass filter.

- Diffuse from multiple courses.
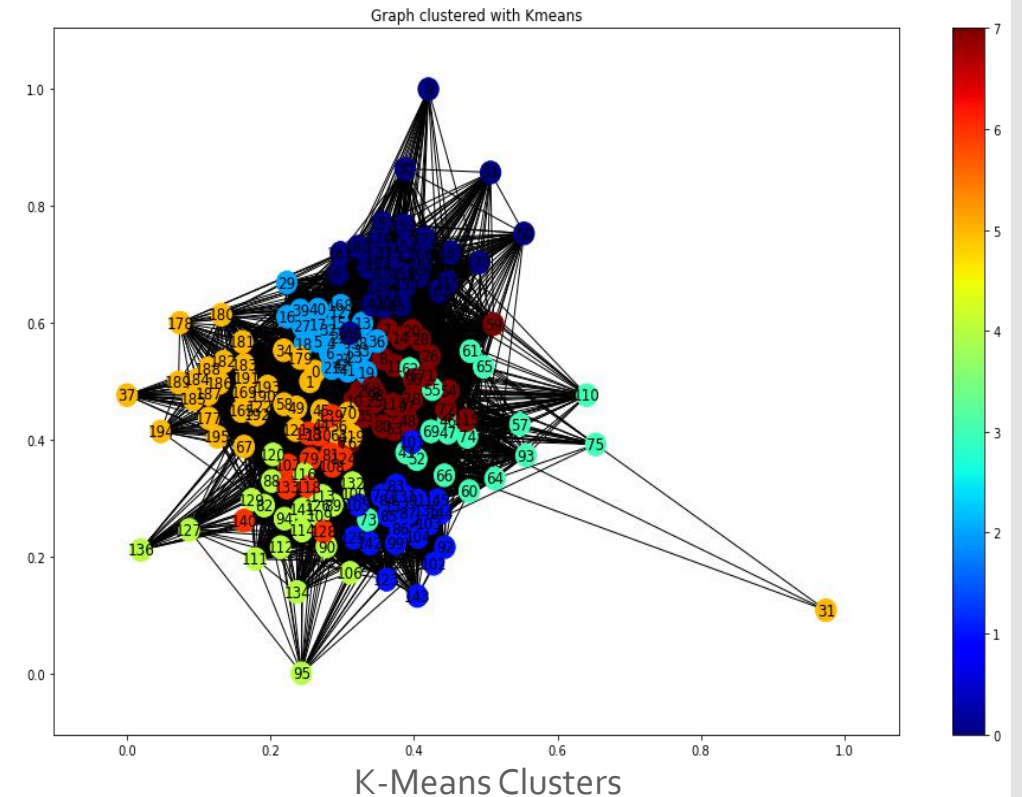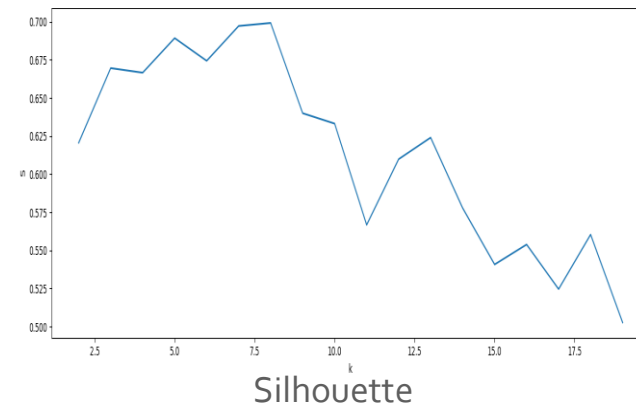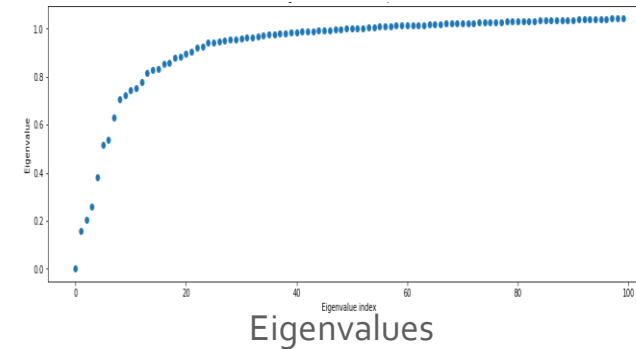
- Recommendation: sort the courses in descending order



Heat



Diffusion with 2 courses

# Unsupervised Spectral Clustering With K-Means



Eigenvalues

Silhouette

K-Means Clusters

- Selection of k based on the eigenvalues of the normalized Laplacian
- Selection of k based on the silhouette analysis to assess the distance between clusters
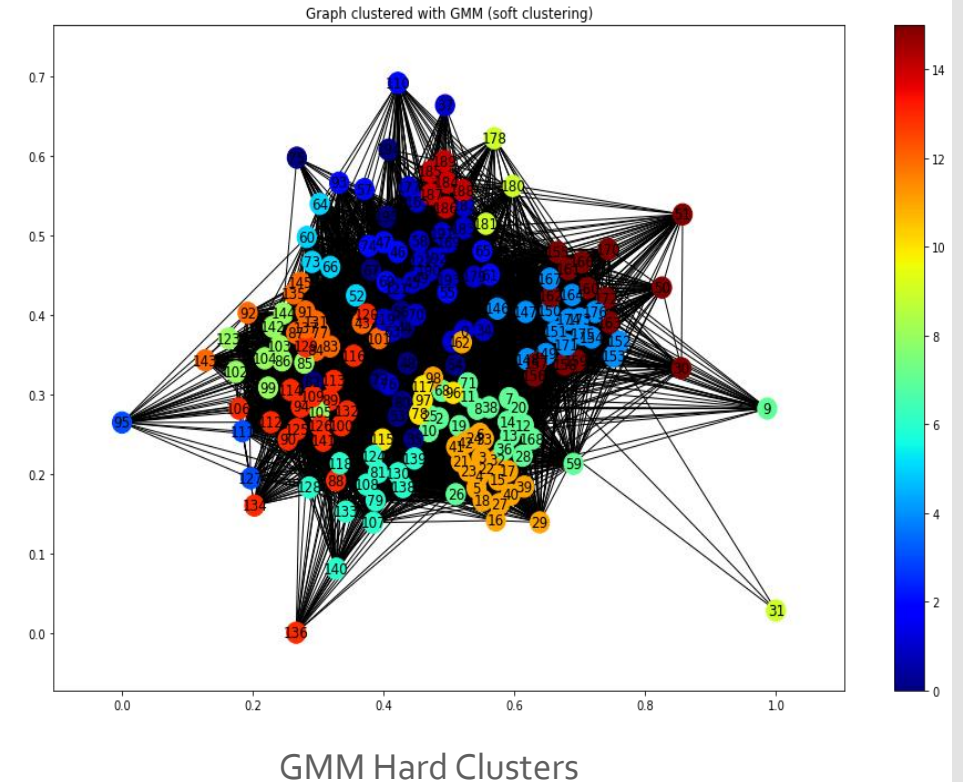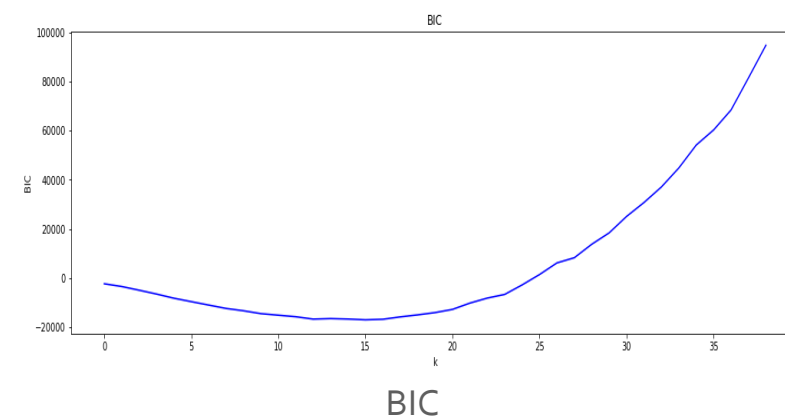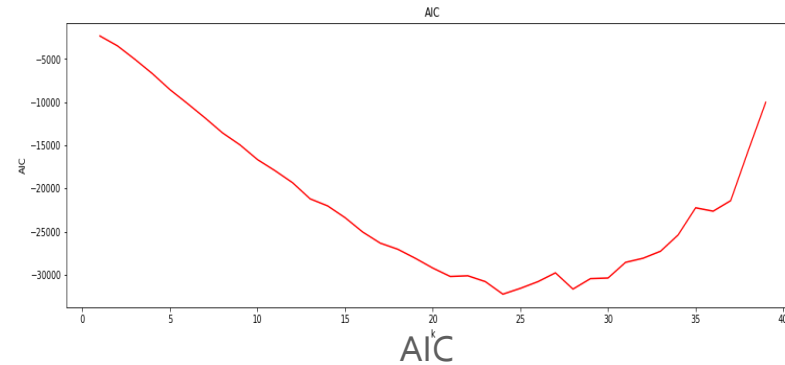
15

# Unsupervised Spectral Clustering With GMMs



AIC



BIC



GMM Hard Clusters

- Selection of k based on AIC and BIC citerions
- 15 clusters with BIC

# Results

- Recommendation From **Automatic Speech Processing**

| | Results |
|---|---|
| K-means | Biomedical signal processing<br>Microwaves<br>Systems and architectures for signal processing<br>Wireless receivers<br>Image analysis and pattern recognition<br>microwaves<br>Speech processing… |
| GMM soft clustering | Biomedical signal processing<br>Image analysis and pattern recognition<br>Systems and architectures for signal processing<br>Microwaves<br>Speech processing… |
| Diffusion | Advanced machine learning<br>Applied machine learning<br>Biomicroscopy I<br>Biomicroscopy II<br>Image optics<br>Optical communications… |