

Introduction to statistical learning: supervised and unsupervised learning

prof. Janez Povh, Ph.D.

University in Ljubljana, Faculty of mechanical engineering

21st February 2017

① Supervised and unsupervised learning

What is SL?

- ① Statistical learning (SL) is about (**deeper**) understanding of the data!
- ② SL tries to answer:

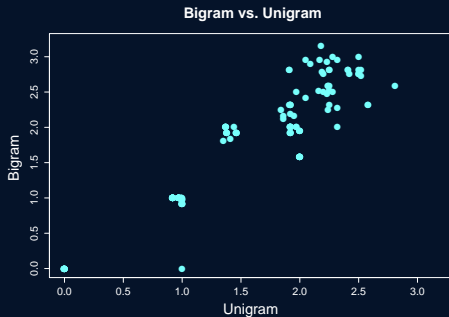
Question: SL

What are the hidden relations between the data instances and/or the data variables?

Example: regression

Example: SL

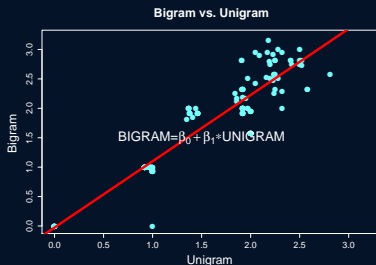
- dataset of songs;
- song's complexity: **UNIGRAM** and **BIGRAM** entropy.



Example: regression

Question: SL

What is the **relation** between BIGRAM and UNIGRAM?

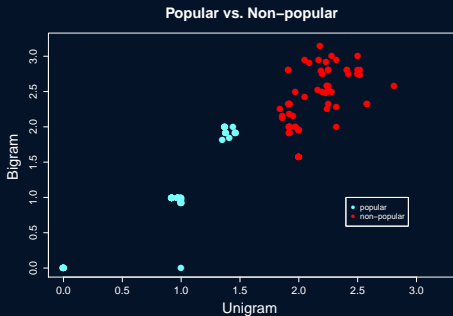


Example: **classification**

Example: SL

We have a dataset of songs. For each song we know:

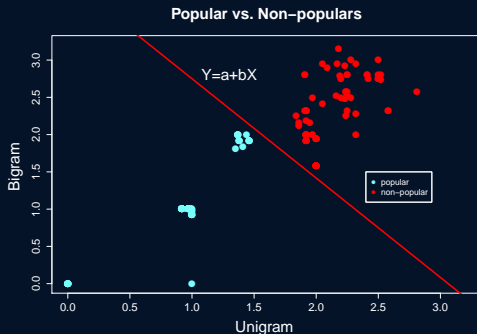
- song's complexity: **UNIGRAM** and **BIGRAM**.
- song's popularity: **popular** and **non-popular**.



Example: classification

Question: SL

Can we **predict** popularity based on BI-GRAM and UNI-GRAM entropy?



Supervised learning - definition

Definition: Supervised learning

- We **consider** variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ (predictors) and \mathbf{Y} (response);
- We **measured** $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ and \mathbf{Y} on given (**training**) set;
- We **look for** function \mathbf{f} such that

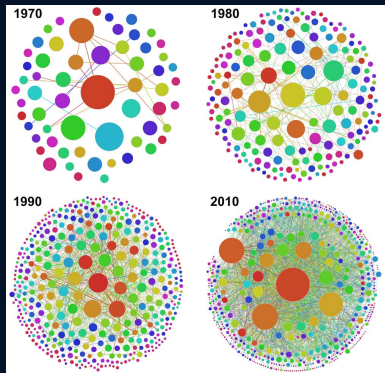
$$\mathbf{Y} = \mathbf{f}(\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p).$$

- **Examples:** classification, linear regression.

Example: clustering

Example

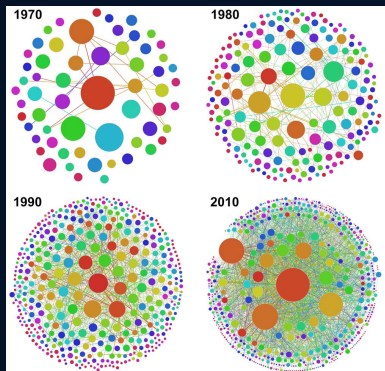
- Dataset of Slovenian scientists that published at least 1 paper in 1970-2015;
- Their collaboration in 1970, 1980, 1990 and 2010:



Example: clustering

Question

Can we **explain** the groups (communities) from the figure?

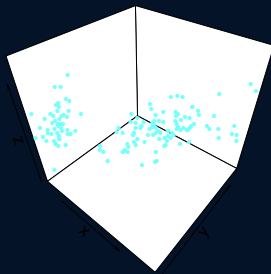


Example: dimension reduction

Example

- Dataset of cancer patients;
- p variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ (**features**): sickness history, current status etc.
- How to **best** visualize these data in 2-dimensions?

3D scatter plots



Example: dimension reduction

Example

- Dataset of cancer patients;
- p variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ (**features**): sickness history, current status etc.
- How to **best** visualize these data in 2-dimensions?

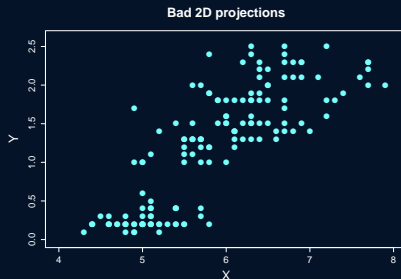


Figure: Scatter plot for 2D data

Example: dimension reduction

Example

- Dataset of cancer patients;
- p variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ (**features**): sickness history, current status etc.
- How to **best** visualize these data in 2-dimensions?

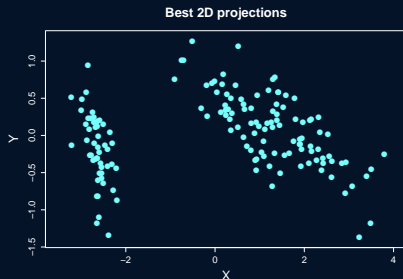


Figure: Scatter plot for 2D data

Definition: Unsupervised learning

- We **consider** variables $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$;
- We **measured** $\mathbf{X}_1, \mathbf{X}_2, \dots, \mathbf{X}_p$ on given set;
- We **look for**
 - hidden groups of data instances (clusters) (\implies **CLUSTERING**) OR
 - **few** new variables enabling better visualization: (\implies **PRINCIPAL COMPONENT ANALYSIS**) OR
 - **few** new variables enabling better interpretation: (\implies **FACTOR ANALYSIS**).

Supervised vs. unsupervised learning - summary

	SUPERVISED LEARNING	UNSUPERVISED LEARNING
AVAILABLE DATA	X_1, X_1, \dots, X_p, Y	X_1, X_1, \dots, X_p
GROUND TRUTH	<ul style="list-style-type: none">• We know the answer on available data• Results can be evaluated• Results can be compared	<ul style="list-style-type: none">• Several solutions possible• No universal measure to compare solutions• Subjectivity