

UNIVERSITÀ DEGLI STUDI DI TORINO
SCUOLA DI SCIENZE DELLA NATURA
Corso di Laurea Magistrale in Fisica dei Sistemi Complessi



Tesi di Laurea Magistrale

THE FANCY TITLE OF MY FANCY THESIS

Relatore:
Prof. Michele Caselle

Corelatore:
Dott. Matteo Osella

Candidato:
Filippo Valle

Controrelatore:
Dott. Matteo Cereda

Anno Accademico 2018/2019

Computers are incredibly fast, accurate, and stupid.
Human beings are incredibly slow, inaccurate, and brilliant.
Together they are powerful beyond imagination Albert Einstein,

Abstract

The interest in studying complex systems is increasingly spreading. Complex systems can be found anywhere and many common behaviours are observable, systems with different origins and purposes may share, for instance, some statistical laws.

An example can be the Zipf's law, well known in linguistics and texts analysis can be easily observed in the distribution of gene expressions in different samples of cancer tissues.

In recent years datasets with a large amount of cancer samples' data cancer are available, the most complete is The Cancer Genome Atlas (TCGA). From this dataset it is easy to get, for example, gene expression data from RNA sequencing experiments together with a lot of information about the sample itself.

If one studies the number of samples in which a gene is expressed above a certain threshold, the so called occurrence, it is easily verified that there are different kinds of genes. Some are present in the majority of samples, some others are present only in a subset of the whole dataset. This exact same behaviour can be found analysing words in a corpus of texts; some words, such as *the*, are present everywhere, other specific words are present only in texts regarding a certain subject. This suggests that there are similarities between a system of words and documents and a system of genes and samples.

Given a corpus of documents, they can be classified by their specific subject. In a similar way a set of samples can be classified, for example, by the tissue it comes from or by the type of the disease it is referred to.

The similarities between gene expression data and linguistics suggest the possibility to use topic modelling to classify data and separate samples and genes in different clusters. Topic modelling is a set of clustering algorithms in networks' theory. Given a set of words and documents, it describes documents as a mixture of topics. Topics are nothing but communities of words each with a given probability.

Purpose of this work is to build a bipartite network of genes and samples and use topic modelling to find communities. The goal is to separate samples depending on the site the tumour was and the disease type of the sample. Moreover it is possible to separate genes depending on their specific functions. In fact once a community structure of genes emerges, it is possible to run a hypergeometric test on the whole set in order to verify if they reveal some type of enrichment and to inspect their common properties.

The specific algorithm used in this work is particularly unique because it allows overlapping clusters; so it is possible to find genes belonging to different topics, this empowers a lot of new possibilities to investigate the network.

Furthermore a hierarchic approach can be used in finding topics, this let it possible to classify data at different layers. An ideal goal would be to separate healthy and diseased samples at the first layer, then separate by tissue, then by tumour type and so on.

Contents

Chapter 1

Introduction

In recent years the study of complex systems is becoming more interesting especially because some different systems can share interesting and fundamental properties. Network theory has proven to be a useful proxy to model and represent such complex systems.

This work wants to study and find universal statistical laws in different kind of biological systems. If one finds that two different systems share some important laws and data structure, therefore it is possible to use tool from different fields to study and gain information about each others. In particular two datasets containing information about some human healthy and diseased tissues will be analysed. This data come from biological experiments of RNA sequencing.

The ultimate goal of this work would be to study, develop and build machine learning's methods able to discriminate healthy and diseased tissues. Once diseased tissue are found, the next goal is to separate cancer types and ultimately sub-types, which is not always easy clinically.

The methods to gain this goal are derived firstly from linguistics, in particular a topic model approach will be widely described.

In chapter ?? I will describe the datasets used and introduce some basic biological properties of these datasets. In particular I'll use two datasets of gene expression data from diseased tissues and healthy tissues.

In chapter ?? I will describe the basics of component systems and give some basic mathematical definitions of quantities useful in general. This chapter refers in particular to the so called component systems.

In chapter ?? I will represent the gene expression from one sample from TCGA with respect to the genes' rank, one can easily obtain a Zipf's law. This law is well-known and in-depth studied in linguistics, demonstrating that different sources of data (genomic and linguistics) can share some statistical properties.

Demonstrated that linguistics and biological data share some laws in section ?? I will use topic modelling to perform network analysis on datasets. Using topic modelling one would find the inner structure of the samples. One would find clusters such that all samples in a cluster share the tissue and tumour type. As far as a document can contain a mixture of similar topics a single tumour can be very heterogeneous.

In chapter ?? I will discuss the results and the future developments

Many methods of the pipeline written in c++ using openMP and Boost [?] are available at <https://github.com/fvalle1/tacos>. During this work I used different python libraries such as pandas [?], scipy [?], numpy [?] and matplotlib [?]. Some analysis required tensorflow [?] and pySpark [?]. The topic modelling require graph-tool library [?]. The full work repository is on Github at http://github.com/fvalle1/master_thesis.

Chapter 2

Data presentation

2.1 Dataset

The goal of [?] [?] [?] [?]

2.2 RNA-Sequencing

Data come from a RNA-Sequencing [?] experiments.

RNA-Sequencing data provide a unique snapshot of the transcriptomic status of the disease and look at an unbiased population of transcripts that allows the identification of novel transcripts, fusion transcripts and non-coding RNAs that could be undetected with different technologies.

Briefly, long RNAs are first converted into a library of cDNA fragments through either RNA fragmentation or DNA fragmentation (see main text). Sequencing adaptors (blue) are subsequently added to each cDNA fragment and a short sequence is obtained from each cDNA using high-throughput sequencing technology. The resulting sequence reads are aligned with the reference genome or transcriptome, and classified as three types: exonic reads, junction reads and poly(A) end-reads. These three types are used to generate a base-resolution expression profile for each gene, as illustrated at the bottom; a yeast ORF with one intron is shown.

The general steps to prepare a complementary DNA (cDNA) library for sequencing are, in general:

- RNA Isolation: RNA is isolated from tissue and the amount of genomic DNA is reduced
- RNA selection/depletion: To analyze signals of interest, the isolated RNA can either be kept as is or filtered for RNA that binds specific sequences. The non-coding RNA is removed because it represents over 90% of the RNA in a cell, which if kept would drown out other data in the transcriptome
- cDNA synthesis: RNA is reverse transcribed to cDNA (DNA sequencing technology is more mature). Fragmentation and size selection are performed to purify

sequences that are the appropriate length for the sequencing machine. Fragmentation is followed by size selection, where either small sequences are removed or a tight range of sequence lengths are selected. Because small RNAs like miRNAs are lost, these are analyzed independently. The cDNA for each experiment can be indexed with a hexamer or octamer barcode, so that these experiments can be pooled into a single lane for multiplexed sequencing.

In order to collect Gene expression data is sufficient to count how many reads are mapped to a specific exon or gene.

Data was collected from TCGA¹ dataset at <https://portal.gdc.cancer.gov> [?].

The particular datatype considered was *Gene Expression Quantification*, with experimental strategy RNA-Sequencing. At the end the downloaded dataset consisted of:

- $N = 60483$ genes as **components**
- $R = 10672$ files as **realizations**

This type of data is just a small portion of all data available on the portal, this are the most useful data for this type of analysis.

As highlighted in [?] these data present a challenge to clustering tools, because of both the relatively large number of samples and the complex structure created by the inclusion of many different tissues

Attempts were made from GTEx [?] which is a similar source of data from healthy tissues. [?] tried to unify data from this two different sources and data are available from [?]. Anyway gene expression data were downloaded directly from GTEx v7²

2.2.1 normalization

Usually gene expression data can be normalized in different ways

- Counts
- RPK
- FPKM
- FPKM-UQ

Counts correspond to raw data. Anyway longer genes may have much reads than smaller gene, so it can be useful to normalize this data.

RPK³ solves this by dividing counts by gene length L ,

$$RPK = \frac{counts}{L}$$

¹The Cancer Genome Atlas

²<https://gtexportal.org/home/datasets>

³Reads Per Kilobase of transcript

FPKM⁴ are provided. FPKM calculation normalizes read count by dividing it by the gene length and the total number of reads mapped to protein-coding genes.

$$FPKM = \frac{RC_g * 10^9}{RC_{pc} * L} \quad (2.1)$$

where

- RC_g : Number of reads mapped to the gene
- RC_{pc} : Number of reads mapped to all protein-coding genes
- L : Length of the gene in base pairs

FPKM can be normalized to the 75th percentile read count value for the sample, in this case it is called FPKM-UQ. FPKM-UQ is obtained by:

$$FPKM - UQ = \frac{RC_g * 10^9}{RC_{g75} * L} \quad (2.2)$$

where

- RC_{g75} : 75th percentile read count value for genes in the sample

2.3 Clean data

2.3.1 Protein coding

The whole dataset contains infos on approximately 60000 elements with different *ENSG* identifier. Only $\simeq 20000$ of this are protein coding genes, using Ensemble⁵ protein coding genes are selected.

2.3.2 Thresholds

In order to filter out noise, it is useful to put a threshold on the data. Considering data in *FPKM* format, it is common opinion that values below 0.1 can be considered noise. Moreover data above 10^5 are trashed out, because they are a symptom of some kind of errors during experiment.

Given this thresholds ?? becomes

$$o_i = \frac{\sum_{j=1}^R \theta(n_{ij} - 0.1) * \theta(10^5 - n_{ij})}{R} \quad (2.3)$$

⁴Fragments Per Kilobase of transcript per Million mapped reads

⁵<https://ensemble.org>

Chapter 3

Data structure

The data studied in this work can be represented as component systems. These component systems can be represented by a two dimensional matrix in which rows represent components and columns are the possible realizations buildable given subset of the components. The entries of this matrix are the number of the components on the row needed during the realization of the column. In figure ?? an example of this kind of matrices.

metti citazioni [?] [?] [?]

3.1 Component systems

$$\begin{array}{c} \text{Components} \end{array} \left(\begin{array}{ccccc} & \text{Realizations} & & & \\ n_{11} & n_{12} & n_{13} & \dots & n_{1R} \\ n_{21} & n_{22} & n_{23} & \dots & n_{2R} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ n_{N1} & n_{N2} & n_{N3} & \dots & n_{NR} \end{array} \right)$$

Figure 3.1: Structure of a matrix representing component systems with $i = 0 \dots N$ rows and $j = 0 \dots R$ columns

The most common example of such systems is a set of books, in this case one puts on the rows the words in the whole vocabulary and the books' titles on the columns; the entry that correspond to row i and column j is the number of times the word i appears in book j . The same happens if one considers Wikipedia's pages. Other examples are Lego[®] sets where components are the Lego[®] bricks and realizations Lego[®] packages and protein domains; all these were described and well studied in [?].

Given a matrix with N components on the rows and R realizations on the columns and relative abundances n_{ij} as the entries, it is interesting to study some quantities that are universal and general characteristics of component systems.

First of all, the **occurrence** of a component, defined as

$$o_i = \frac{\sum_{j=1}^R (1 - \delta_{n_{ij},0})}{R}, \quad (3.1)$$

is the number of realizations in which the component's abundance is not null. A component that is present in all the realizations has got $O_i = 1$. Components with high ($\simeq 1$) occurrence are present in mostly all realisations of the datasets, in linguistics this components are articles. Components with low occurrence $\simeq 0$ are present only in a few realisations and are the most specific ones.

The sum across all columns is called **abundance** of a component and is defined as:

$$a_i = \sum_{j=1}^R n_{ij}; \quad (3.2)$$

dividing this by the global abundance

$$a = \sum_{i=1}^N a_i \quad (3.3)$$

naturally brings to the **frequency of a component** in the whole corpus

$$f_i = \frac{a_i}{\sum_{c=1}^N a_c}. \quad (3.4)$$

Dividing the abundance of a component by the sum of all abundances in a realisation gives the frequency of the component in that specific realization

$$f_{ij} = \frac{n_{ij}}{\sum_{c=1}^N n_{cj}}. \quad (3.5)$$

The sum of all abundances in the same realization

$$M_j = \sum_{c=1}^N n_{cj} \quad (3.6)$$

is called **size**.

It is expected that frequencies distribute according to the so called Zipf's law

$$f_i \propto r_i^{-\alpha} \quad (3.7)$$

where r is the rank: the position of a component when sorting all components by their frequencies in the whole dataset.

3.2 Universal laws in RNA-Seq

3.2.1 TCGA

Analysing TCGA dataset [?] the first interesting analysis is to plot the sorted abundance, this gives the so called Zipf's law. The analysis were made considering *Gene Expression Quantification* as data type, *Transcriptome Profiling* as data, *RNA-Seq* as experimental strategy, *HTSeq - Counts* or *HTSeq - FPKM* as workflow type. 5000 samples were downloaded and analysed. In figure ?? it is shown the frequency ranked

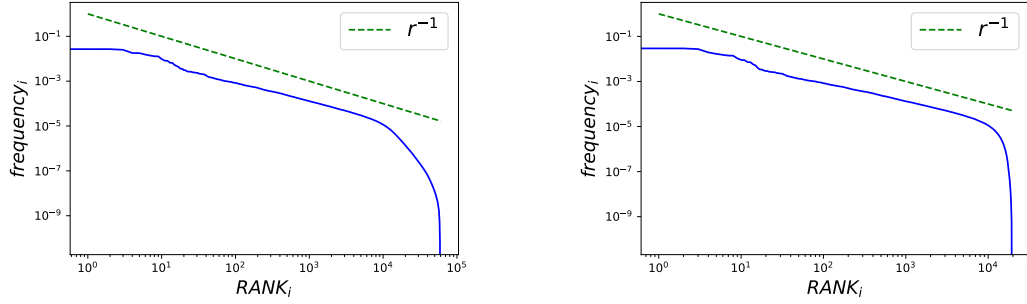


Figure 3.2: Zipf's law from FPKM normalised data. On the right considering only protein coding genes

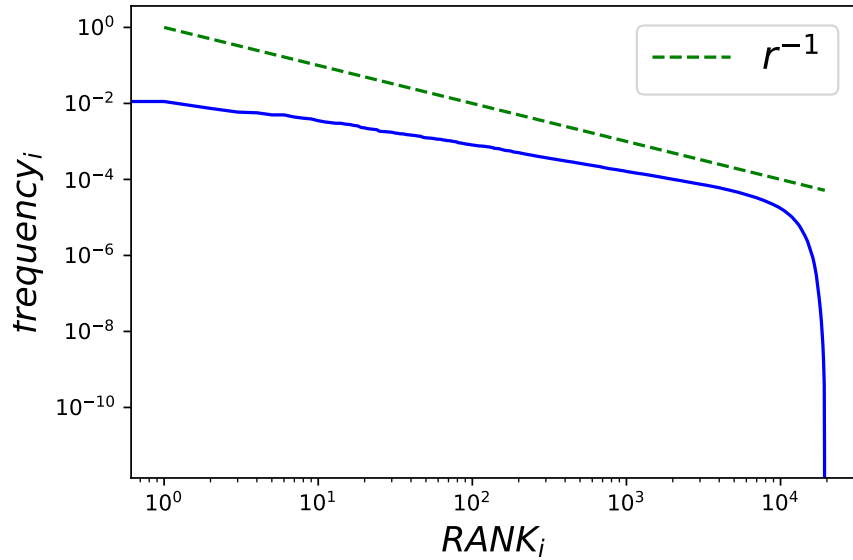


Figure 3.3: Zipf's law of protein coding genes considering counts

plot. It is interesting that this kind of data distribute according a power law with exponent close to 1, this same behaviour can be found in completely different systems such as linguistics' ones [?]. Another interesting fact is that considering in the analysis also non-coding genes gives a double-scaled power law. This is due to the fact that non coding genes are also more specific and rare, so their frequencies are quite small compared to protein coding genes.

Changing normalisation and considering counts instead of FPKM, the result is quite similar. The power law is more flat, meaning that genes have more similar abundances in the whole dataset.

3.2.2 GTEx

A pretty similar analysis can be made on GTEx's [?] healthy samples. RNA sequencing raw counts data were download from file version *2016-01-15 v7 RNASEQCv1.1.8*. All

~ 11000 samples available were considered at this time.

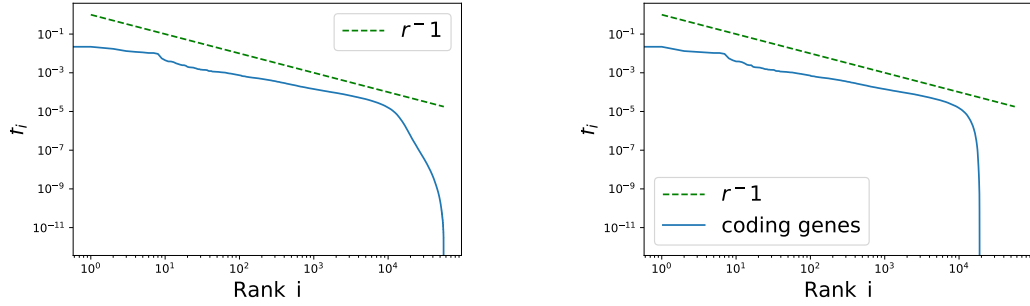


Figure 3.4: Zipf's law from GTEx count data. On the left all genes considered, on the right only protein coding ones

Not surprisingly in the GTEx dataset it is retrieved the same behaviour at this time. The power law with exponent $\simeq 1$ is found and considering non coding genes gives to a knee in the power law.

Going further in the analysis it is possible make an histogram of occurrences defined by O_i , also known as U_i .

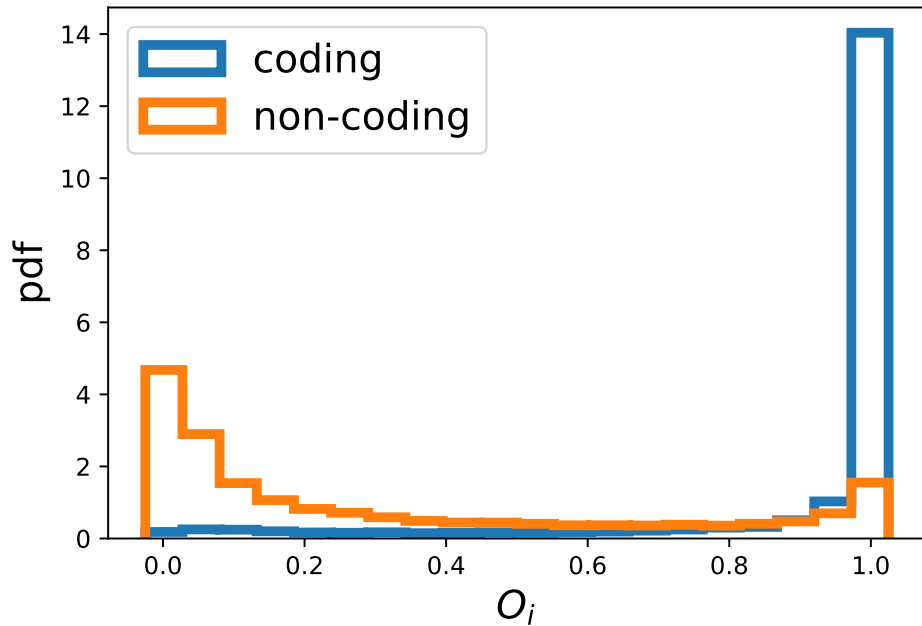


Figure 3.5: The histogram of the occurrences O_i

Also in this kind of analysis it is possible to see the different behaviour of coding and not coding genes. The protein coding genes express almost in every sample, so their occurrence is near to 1, non coding genes are more specific, so they are present only in a subset of the dataset and many of the have small occurrence. The same behaviour can be observed considering just all samples of a given tissue. In this case

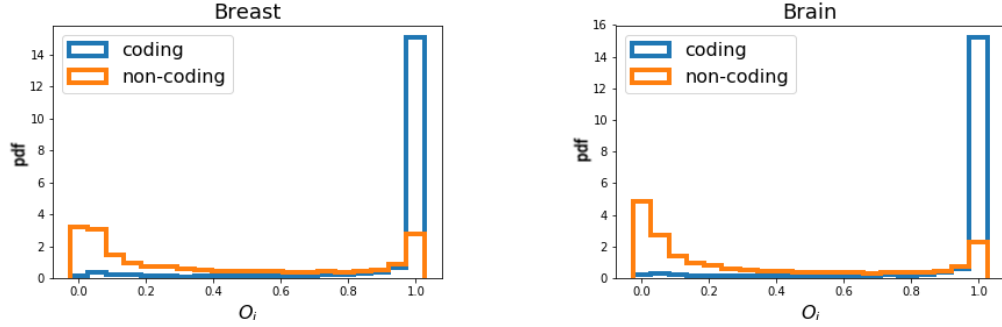


Figure 3.6: Same behaviour is observed looking at one tissue a time.

$O_i = 0$ means that the genes has a non zero expression in just one of the samples of the tissue considered; in other words if a gene never express in a tissue it is not considered when constructing these tissue specific U distributions.

From now on except were explicitly declared analysis will be made considering protein coding genes and counts with no normalisation.

3.3 Null model construction

The kind of data considered in this work comes from RNA Sequencing experiments. This experiments use wet biology methods to extract information from samples. If one imagines it exists an unknown function that describes the gene expression across the samples considered, what experimenters people do is to sample this function, picking up some genes.

In this section it is described a null model of sampling, this is useful to verify if the data distributions seen are just an effect of this experimental sample or if they carry some useful and interesting information.

As described in [?] a random matrix has to be created. This matrix is a collection of components and realizations exactly as [?]. The values of abundances of each component in each realization n_{ij} are randomly assigned with a probability determined by the global abundance in the whole dataset [?]. Values of each column are extracted until the size [?] is reached. Strictly speaking it is a multinomial process

$$P(n_i; M) = \frac{M!}{\prod_{i=1}^N n_i} \prod_{i=1}^N f_i^{n_i} \quad (3.8)$$

where n_i is the number of components with frequency f_i , being $f_i = \frac{a_i}{\sum_{i=1}^N a_i}$ as defined in [?].

Figure ?? shows an example of this, M components are picked up with respect to their frequency in the dataset. The most abundant components, which are also the ones with higher frequency (frequency is nothing but the normalised abundance), have a greater probability to be picked up.

Using this construction on data of counts on both dataset, by definition the Zipf's law sampled are identical to the data's one. By construction the distribution of the

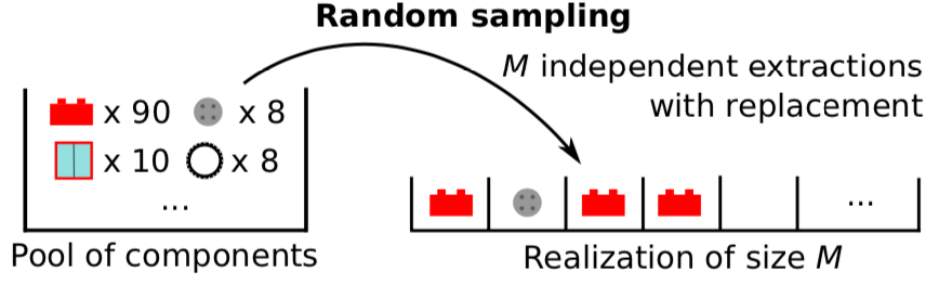
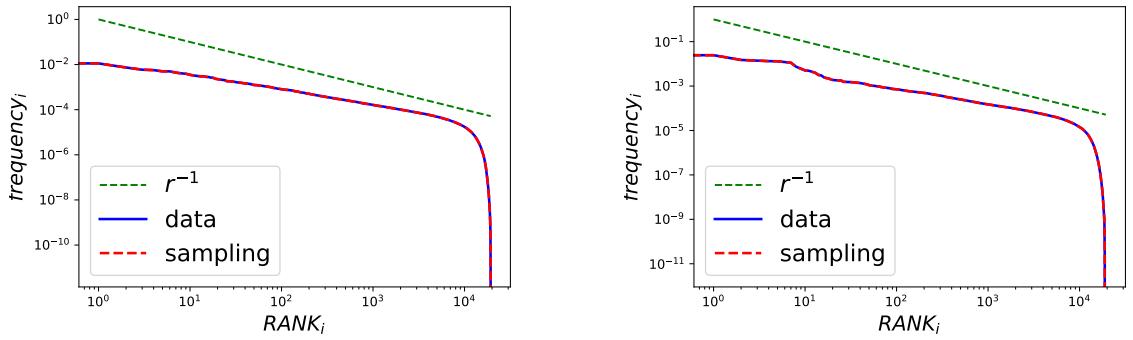
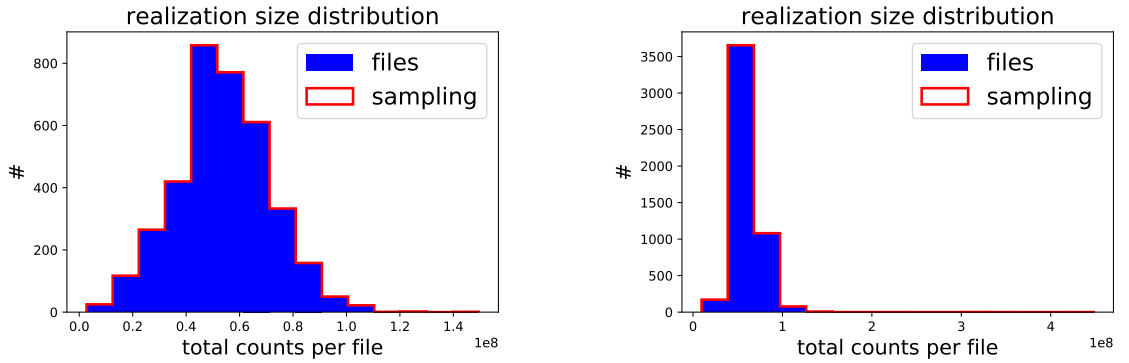

 Figure 3.7: Random sampling of components to build a realization of size M


Figure 3.8: Zipf's law sampled; TCGA(left) and GTEx (right)


 Figure 3.9: Distribution of sizes M ; TCGA(left) and GTEx (right)

sizes of the sampling and of the data are identical.

Looking at the U s, it is evident that data is different from sampling. This is a signal that the null model is not enough to explain the data matrices. In particular from figure ?? it is evident that the null model generate the matrices in a manner such that more components have high occurrence with respect to the original data. This can be easily explained, in fact in real world there are some genes that are highly expressed but only in a subset of the whole dataset; these genes are specific for certain type of samples. The null model gets the information the such genes are highly expressed from the abundance and so samples these quite often (components with high abundance have a greater chance to be picked up by the null model sampling).

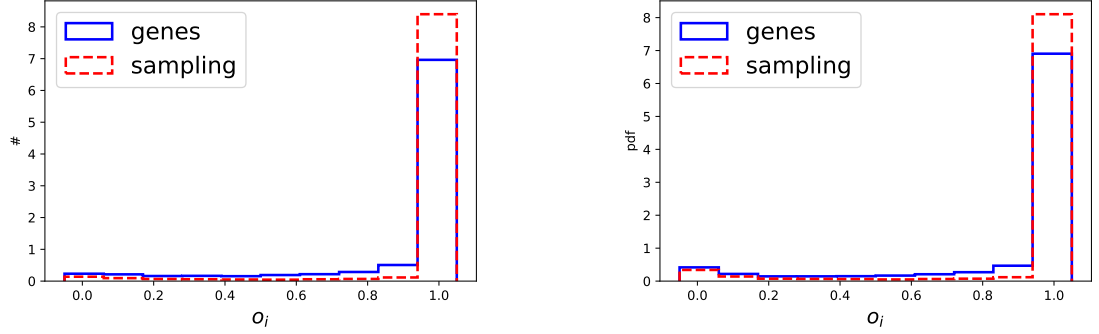


Figure 3.10: Occurrence distributions; TCGA(left) and GTEx (right)

Looking at the Heaps's law [?], again the curves differ and the null model is not enough complete to explain the trend. In figure ?? the Heaps's law is presented compared to the one obtained by sampling, note that each data point share the abscissa with a sampling one (figures ?? are nothing but the histograms of the abscissas of ??). It happens that the sampling curve is above the data's one. This means that to build a sample of size M just by sampling it is necessary to use a greater number of different genes than the number of different genes actually expressed in nature. In other words in real world are expressed only the genes that are really useful in the sample, and this is not describable just by sampling. This fact is coherent with the fact that the U s differ. Another way to see this is looking at the histograms of the number of different

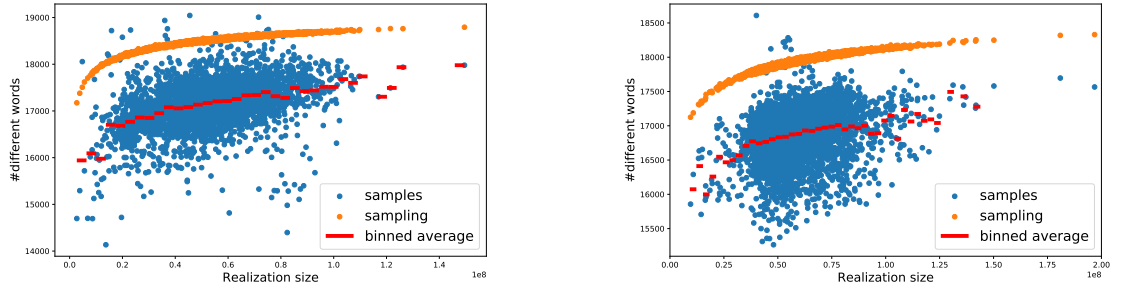


Figure 3.11: Heaps' law; TCGA(left) and GTEx (right)

genes expressed, actually the distribution of the ?? y axis. Figure ?? shows that these distributions are completely different if one looks at the data and at the samples.

3.4 Statistical laws differentiate by tissue

Observing the GTEx dataset of healthy samples it is possible to study how it is possible to see the tissue differentiation and how to study tissues' differences, [?] suggests the approach.

First of all could be interesting to study which is the fraction of transcript that can be explained by a certain number of genes. One can reduce the realisations to the ones

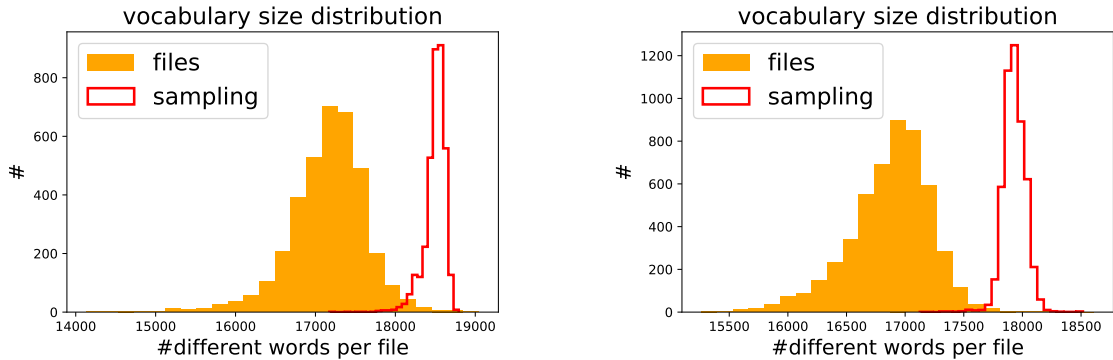


Figure 3.12: Occurrence distributions; TCGA(left) and GTEx (right)

that share the tissue. Then one estimates the average per each component (gene), at this point one has the average abundance of each gene in a tissue, dividing by the sum of all the components it is possible to obtain the fraction of the total counts in the tissue due to each gene. Sorting from greater to smaller and integrating (cumulative summing) one have the fraction of transcript due to 1, 2, 3... genes. This is plot in ?? . Here, if a curve is steep it means that a few genes' counts represent a great fraction

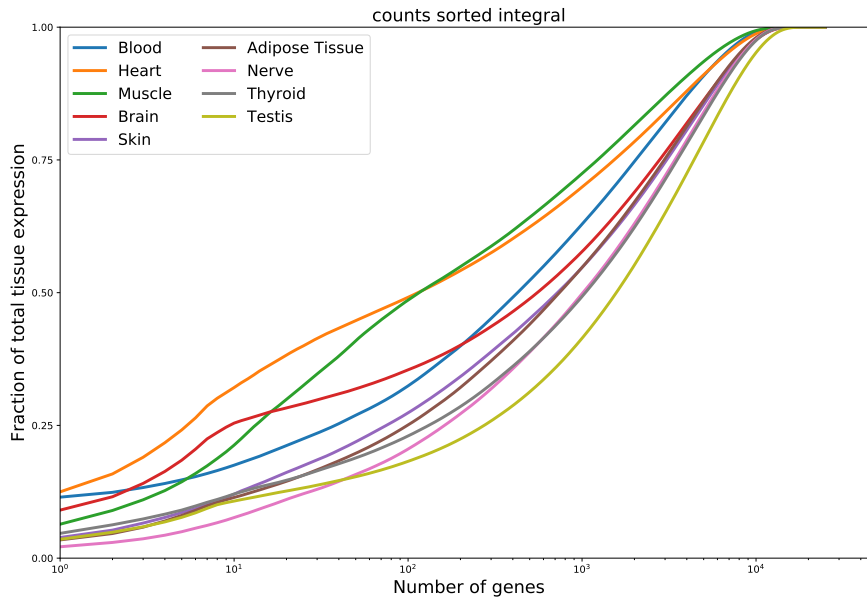


Figure 3.13: The integral of the sorted abundances for each tissue

of the total. If a curve is smooth it means that many genes are necessary to describe the whole transcriptome for that particular tissue. This analysis shows that different tissues have a different complexity in terms of the number of genes necessary to build the transcriptome (in average). In figure ?? the same analysis is done for the sub-tissues of Brain, also this sub-type separate by tissue.

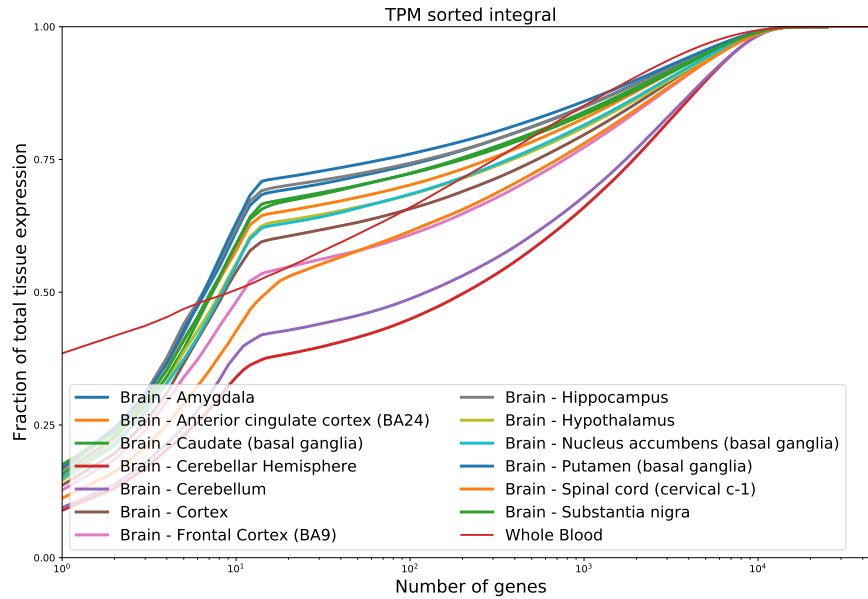


Figure 3.14: The integral of the sorted abundances for sub-types of Brain. This is done using TPM to avoid biases due to gene lengths. Blood is plotted for reference.

Coming back to the Zipf's law ??, it is now obvious that ?? represents nothing but the integral of the Zipf's law. So estimating the Zipf looking at a tissue a time, it is evident that each tissue has its particular slope. The steeper the Zipf the simplest is the tissue: the transcript can be described with a few genes. In figure ?? the tissue with an extreme behaviour.

The point where the ?? reaches 1 corresponds to the total number of genes expressed, the remaining ones have a 0 expression and do not contribute to the transcript. This can be visualised again with the Heaps' law. In figure ?? it is evident that there is some kind of tissue differentiation even when looking at the Heaps' law.

All these analysis suggest that there must be a sort of hidden structure in the data that is somehow related with the tissue each sample comes from. In particular there are many different Zipf's laws hidden behind the data and each sample is build looking at one of these a time. Also given two samples with a similar size, it happens that the number of genes necessary to build that realisation is not always the same (shown by Heaps' law) and it is somehow related to the tissue of the sample.

In conclusion, some interesting laws were found that some statistical.

...

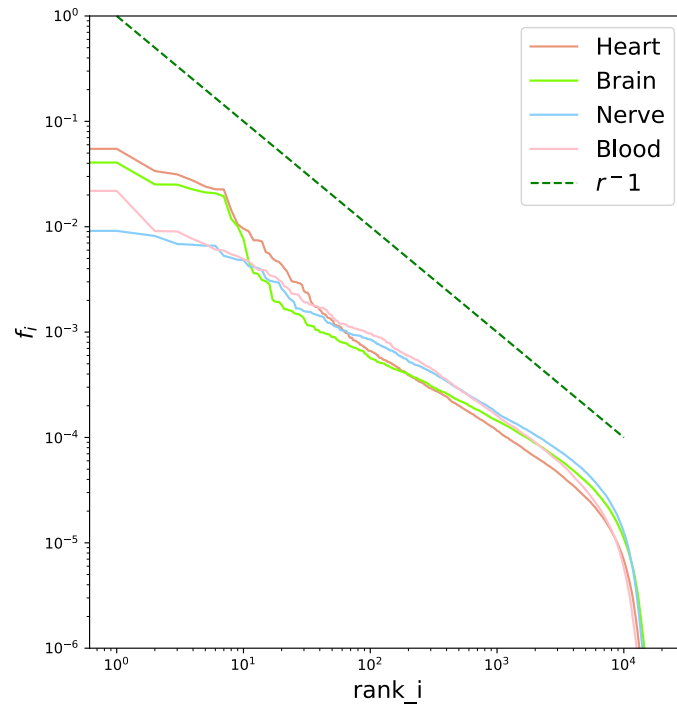


Figure 3.15: The integral of the sorted abundances for each tissue

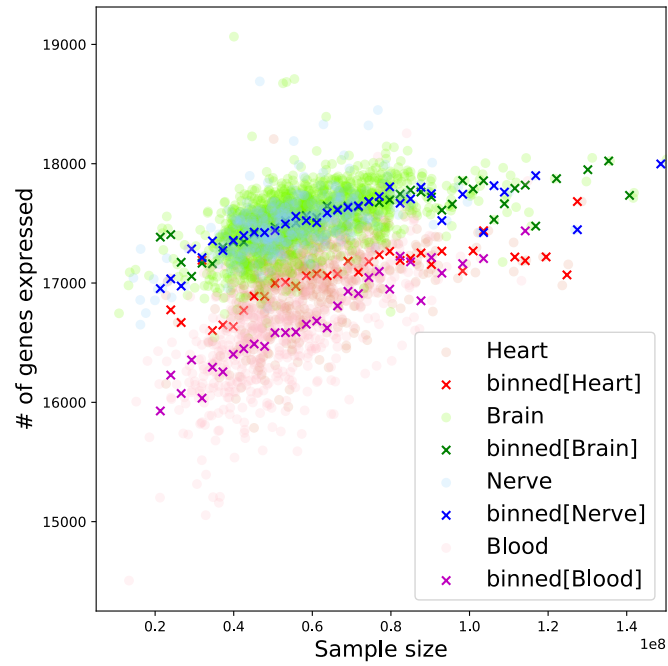


Figure 3.16: The integral of the sorted abundances for each tissue

Chapter 4

Scale laws

One of the goals of this work is to search, reveal, study and use universal laws in bulk gene expression data. As in chapter ?? approaches from different field of science are considered at this point.

In can be interesting to study the behaviour of the gene expression across samples.
taylor [?]

4.1 Scaling

gene expression across samples? gamma?

Given a matrix of components and realisations as ?? with expression entries n_{ij} it is possible to estimate the mean of a row $m_i = \langle n_{ij} \rangle_j$ and its variance $\sigma_i^2 = \langle n_{ij}^2 \rangle_j - \langle n_{ij} \rangle_j^2$.

First of all it could be interesting to represent the variance of expression σ^2 versus the average m across tissues.

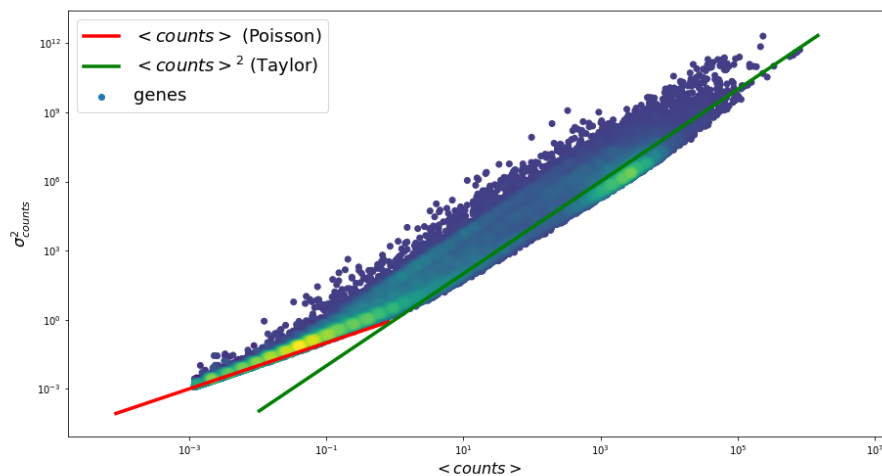


Figure 4.1: Variance versus occurrence

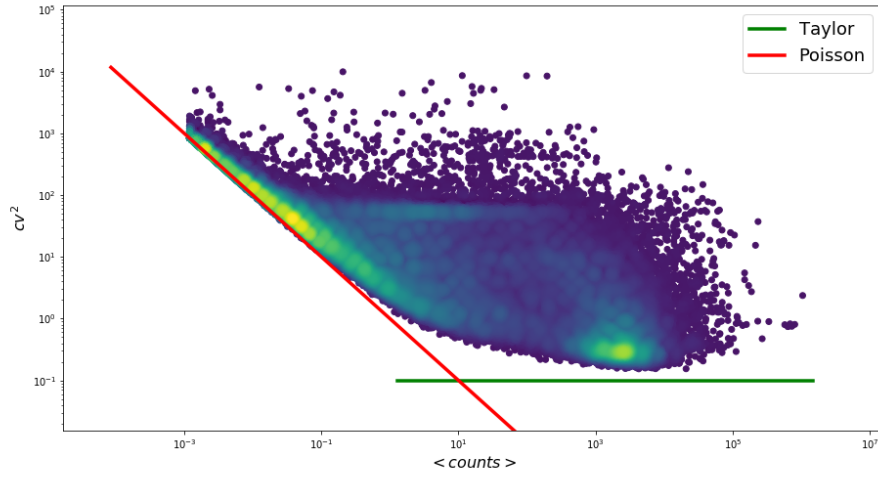


Figure 4.2: Variance versus occurrence

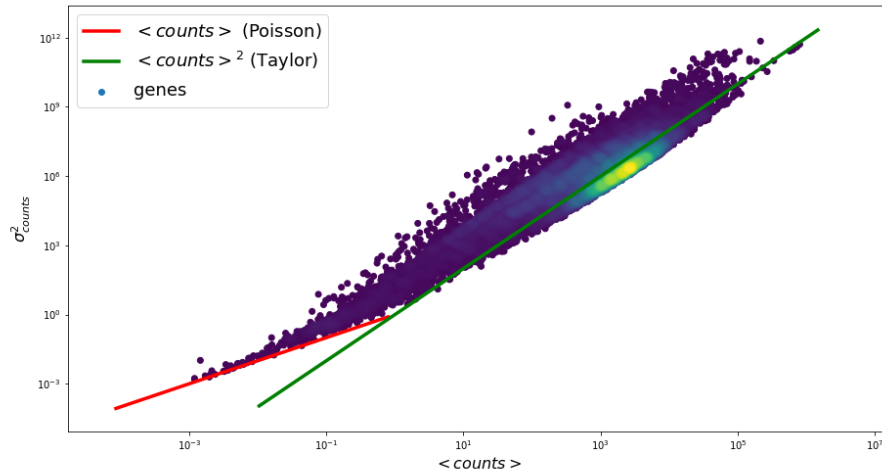


Figure 4.3: Variance versus occurrence

< FPKM > versus occurrence One can be interested in finding genes that are expressed often, and what is the average expression of them. To manage this it is plotted the average expression $\langle FPKM \rangle$ versus the number of samples in which that gene is expressed that is, considering the thresholds θ , $\sum_j \theta(FPKM_{ij} - 0, 1) \theta(10^5 - FPKM_{ij})$

4.1.1 Tissue differentiation

Per gene type scaling

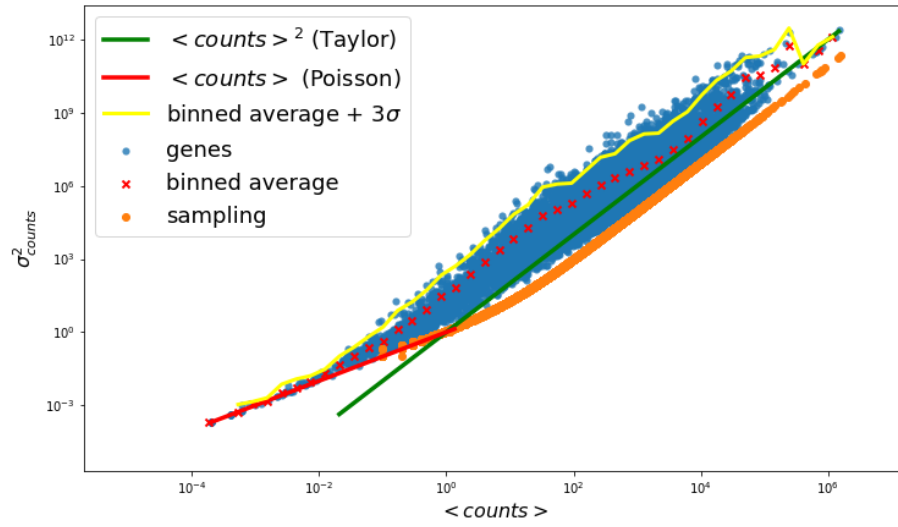


Figure 4.4: Caption

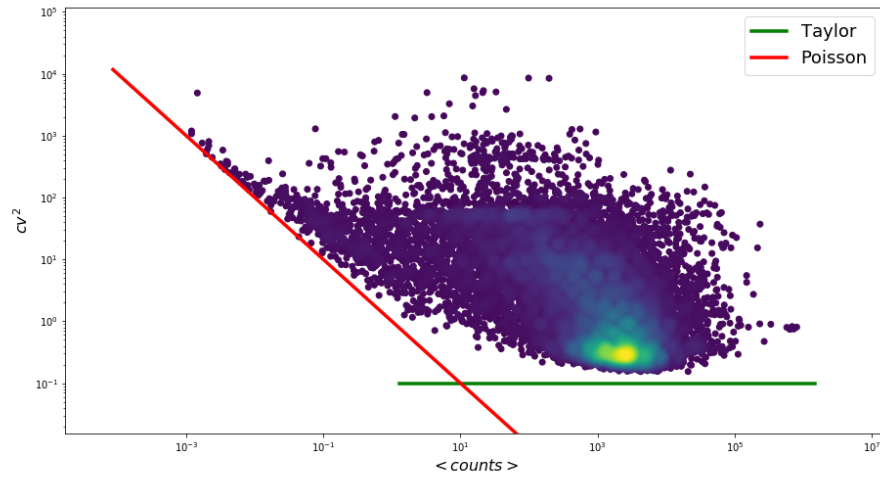


Figure 4.5: Variance versus occurrence

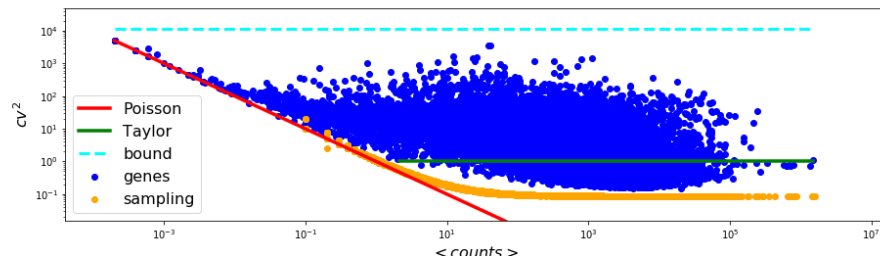


Figure 4.6: Caption

Chapter 5

Topic modelling

Once extensively analysed the structure of the dataset, the goal becomes to develop a machine learning method to learn the hidden structure of the data.

Remembering that in chapter ?? it emerged some kind of structure behind data, where each tissue seemed to be sampled by a different power law, a topic modelling approach it is here proposed.

The idea is that behind data there are hidden variables that describe the relation between the genes and the samples. Let's call these variables topics. Firstly it is necessary to build a bipartite network of genes and samples, than nodes are linked considering the gene expression value in the dataset.

The output of this kind of model are set of genes, the topics, with a probability distribution $P(gene|topic)$ and probability distributions of these topics inside each sample $P(topic|sample)$, both gives the relation between a *sample* and a *gene*.

In this work it is used an innovative and recent approach to topic modelling, the algorithm was presented by [?] and extends the so called stochastic block models [?]. Topic modelling is being developed and studied to approach linguistics problems, so this algorithm was developed considering words and books as input, links represents the abundance of a word in a book. In chapter ?? was evident that there are many similarities between data considered in this work and linguistics' dataset. Referring to data used in this project documents will be **samples**, words will be **genes**. It is expected that topics represent some properties of samples due to the gene expression distribution in samples.

The ultimate goal would be to be able to separate healthy and diseased samples, than separate and find well known tumour types, than extend the actual knowledge and retrieve the tumour sub-types.

One of the advantages of this particular algorithm is that it is hierarchical, so it apply community detection at different layers. So the output has got different resolution, the extreme one is the separation, by definition, between samples and genes, than it is possible to have few big clusters until the other extreme were the number of clusters is comparable with the number of nodes.

What the algorithm does is to run a sort of Montecarlo simulation and find the best partition of the data. The probability that the hidden variables θ describe the data G

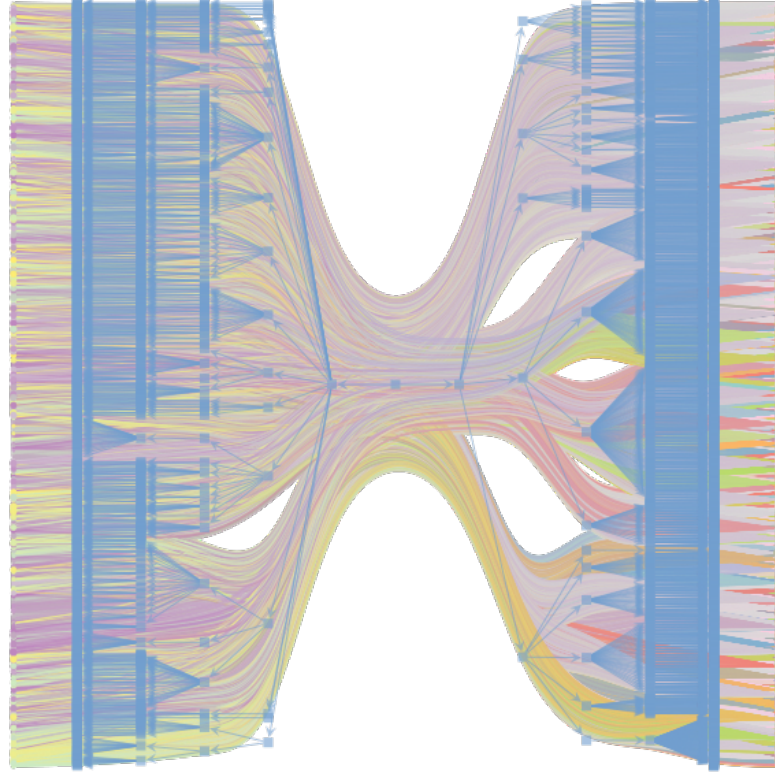


Figure 5.1: An example of a bipartite network. Samples are on the left, genes are on the right. Each link is weighted by gene expression value. On the left side all nodes of the same colour are clusters of samples. On the right side all nodes with same colour are set of genes, also known as topics.

Blue lines represent the cluster structure, each blue square is a set of nodes, lines delineate the hierarchical structure.

It is clear in the middle the network separation in genes and samples.

$P(\theta|G)$ can be written as a likelihood times a prior as

$$P(\theta|G) = \frac{\underbrace{P(G|\theta)}_{\sum_{\theta} P(G|\theta)P(\theta)} \overbrace{P(\theta)}^{\text{prior}}}{P(G)}.$$

It is possible to define a description length

$$\Sigma = -\ln P(G|\theta) - \ln P(\theta),$$

so that $P(\theta|G) \propto \exp -\Sigma$. Moreover the likelihood $P(G|\theta)$, can be written as $\frac{1}{\Omega}$ where Ω is the number of possible realisations given θ . This can be represented as a microcanonical ensemble with entropy $S = \ln(\Omega)$. Note that $\Sigma = S - \ln P(\theta)$ According to [?] entropy S can be written as

$$S = \frac{1}{2} \sum_{r,s} n_r n_s H \left(\frac{e_{rs}}{n_r n_s} \right),$$

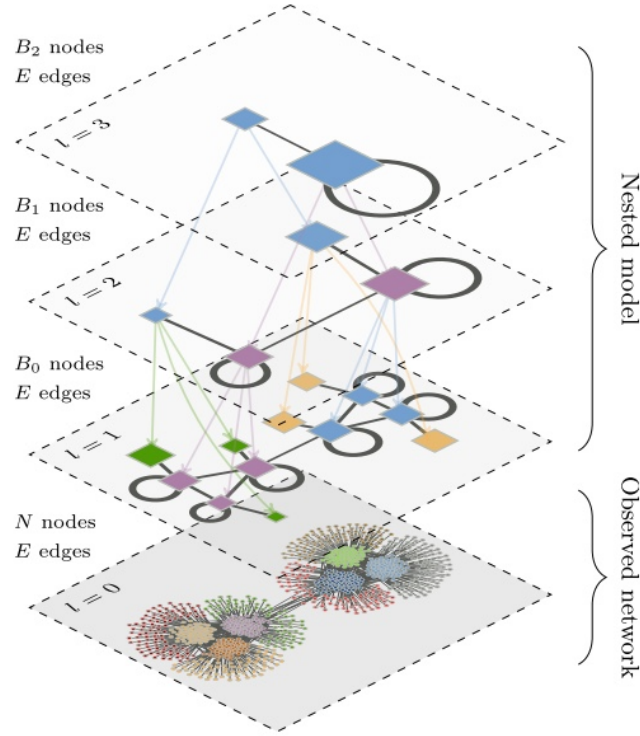


Figure 5.2: Hierarchical structure

where n_r is the number of nodes in block r , e_{rs} the number of links between nodes of group r and group s and H is the Shannon entropy $H(X) = x \log_2(x) + (1-x) \log_2(1-x)$. Note that S is minimal if $\frac{e_{rs}}{n_r n_s}$ is close to zero, r and s are two completely separated blocks or if it is close to 1, r and s are groups with many connections; this allows to find groups with nodes very disconnected or topic and clusters with a lot of connections. The algorithm tries to minimise S , so that Σ is minimised, so $\exp -\Sigma$ is maximised, but this is $P(\theta|G)$ that is the required probability to maximise.

The MonteCarlo works in few steps:

- a node i is chosen
- the group of i is called r
- a node j is chosen from i 's neighbours, the group of j is called t
- a random group s is selected
- move of node i to group s is accepted with probability $P(r \rightarrow s|t) = \frac{e_{ts} + \epsilon}{e_t + \epsilon B}$
- if s is not accepted, a random edge e is chosen from group t and node i is assigned to the endpoint of e which is not in t

in figure ?? an example of these steps.

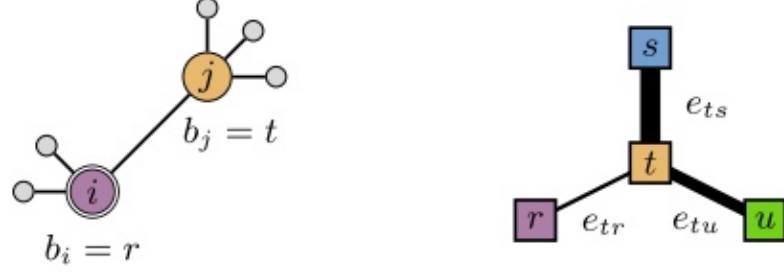


Figure 5.3: Left: Local neighbourhood of node i belonging to block r , and a randomly chosen neighbour j belonging to block t . Right: Block multigraph, indicating the number of edges between blocks, represented as the edge thickness. In this example, the attempted move $bi \rightarrow s$ is made with a larger probability than either $bi \rightarrow u$ or $bi \rightarrow r$ (no movement), since $e_{ts} > e_{tu}$ and $e_{ts} > e_{tr}$.

Once the model run it is possible to estimate the probability distribution of words inside a topic

$$P(w|t_w) = \frac{\# \text{ of edges on } w \text{ to } t_w}{\# \text{ of edges on } t_w}$$

and the topic distribution inside a document

$$P(t_w|d) = \frac{\# \text{ of edges on } d \text{ from } t_w}{\# \text{ of edges on } d}$$

In case of overlapping partitions the presence of a word in a topic is not trivial and can be estimated as

$$P(t_w|w) = \frac{\# \text{ of edges on } w \text{ to } t_w}{\# \text{ of edges on } w}$$

or the presence of a document in a cluster

$$P(t_d|d) = \frac{\# \text{ of edges on } d \text{ to } t_d}{\# \text{ of edges on } d}$$

See appendix ?? for detailed analysis of the math behind the algorithm and <https://cloud.docker.com/repository/docker/fvalle01/hsbm> for the extension of [?] to non linguistics component systems datasets.

5.1 Metrics and benchmarks

Before put the hands on topic modelling, it is useful to define some metrics to test and benchmark the model.

Looking at the cluster side of the network, the outputs are sets of samples, the clusters. One can state the the model works if all, or at least the majority, of samples in the same cluster share some property. Here the tissue is considered as property.

Note that this work's model is a non supervised one, but a ground truth is available from metadata. So every sample has a certain probability to have a certain property (the true tissue label), let's call this $P(C)$ and a certain probability of being in a cluster (model's output), let's call this $P(K)$.

It is possible to define some quantities, the homogeneity

$$h = 1 - \frac{h(C|K)}{H(C)} \quad (5.1)$$

defining the entropy as

$$H(C|K) = \sum_{c \in \text{tissues}, k \in \text{clusters}} \frac{n_{ck}}{N} \text{Log} \left(\frac{n_{ck}}{n_k} \right) \quad (5.2)$$

where n_{ck} is the number of nodes of type c in cluster k , N the number of nodes and n_k the number of nodes in cluster k . It is evident that if all nodes inside cluster k are of the same type c $n_{ck} = n_k$, $H(C|K) = 0$ and $h = 1$, it is actually a complete homogeneous situation.

Another quantity can be defined and it is completeness:

$$c = 1 - \frac{h(K|C)}{H(K)}, \quad (5.3)$$

$H(K|C)$ is defined in the same way as ???. Completeness measures if all nodes of the same type are in the same cluster.

Ideally one wants a method which output is both homogeneous and complete. So it is possible to define the V-measure, which is the harmonic average of the two:

$$MI = 2 \frac{hc}{h+c}, \quad (5.4)$$

which is actually the mutual information between $P(C)$ and $P(K)$ [?].

In the next sections will be studied also the maximum fraction of label in the same cluster $\max_{c \in k} \frac{n_{ck}}{n_k}$. Also the number of different labels in the same cluster will be studied.

5.1.1 LDA

commenti vari

As in [?]

$$P(w, z, \beta, \theta | \alpha, \eta) = \prod_n^{N_d} P(w|z, \beta) P(z|\theta) \prod_k^K P(\beta|\eta) \prod_d^N P(\theta|\alpha) \quad (5.5)$$

where

- N number of documents

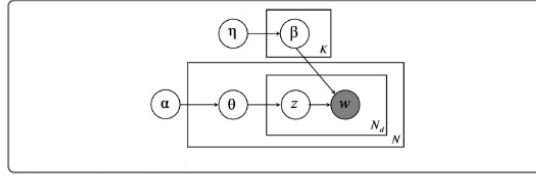


Figure 5.4: LAD scheme

- K number of topics
- w words
- N_d number of words in document d
- η and α are parameters of the model

in ?? $P(\theta|\alpha)$ and $P(\beta|\eta)$ are Dirichlet distributions

the outputs are the topic distribution in documents $P(z|d)$ and the word distribution in topics $P(w|z)$

5.1.2 Hierarchical clustering

da scrivere

`sklearn.hi`

5.2 Pre-process

To make the algorithm faster, could be interesting to do a pre-processing of the data. Vary approaches were tested, all of them involving the quantities defined in ??. The goal is to identify components which are able to best separate the realisations.

Low occurrence genes were selected firstly to make topic modelling. A 0.5 threshold was set. This method select genes that appears (have expression greater than zero) only in less than half samples. This doesn't consider genes that appear everywhere (whith occurrence $\simeq 1$), but changes their behaviour across realisations.

tf-idf (term frequency–inverse document frequency) should help. This approach doesn't consider values, but a transformed version

$$n_{ij}^{new} = \frac{n_{ij}}{M_j} \times (1 - \text{Log}(o_i))$$

which increases the importance of components with small occurrence o_i . This approach doesn't actually select components, which is still an issue.

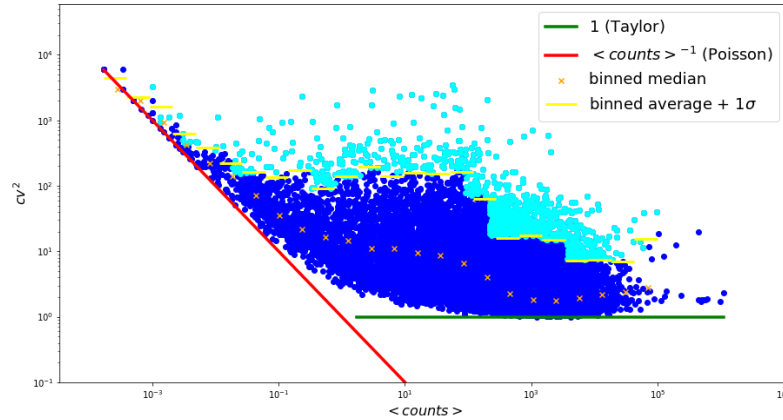


Figure 5.5: Highly variable genes

Highly variable genes can be selected. This is done using the CV^2 analysis from [metti referenza giusta](#). Plotting the coefficient of variation versus the mean for each component reveals which components have higher variance with respect to components which, in average, have a similar behaviour. Binned averages and variances were estimated, and only genes with a CV^2 over a σ greater than the bin's mean were considered.

Distance from boundaries can be a similar and alternative method to select highly variable genes. [metti figura giusta!!](#) The distribution as discussed in [disctuti e linka](#)

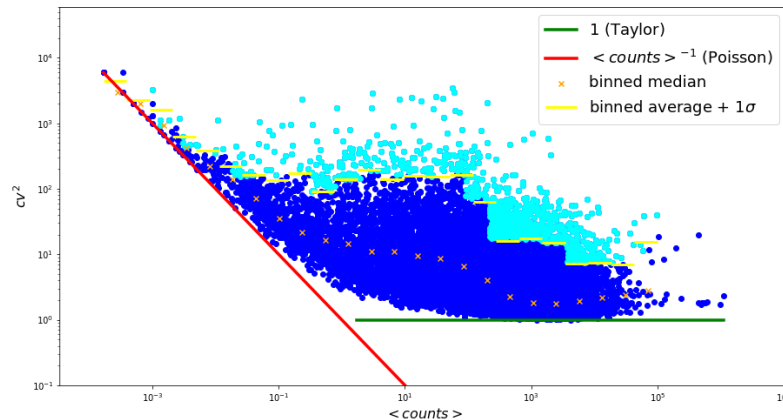


Figure 5.6: Genes distant from the boundaries

have a Poisson-like and a Taylor-like boundaries. So can be considered only components that are the most distant from this boundaries.

Using the last two approaches got the point and actually help topic modelling to succeed.

5.3 Run

5.3.1 Run on GTEx

Firstly the algorithm is run on a subset of 5 tissues of GTEx mettì 5 tissues?

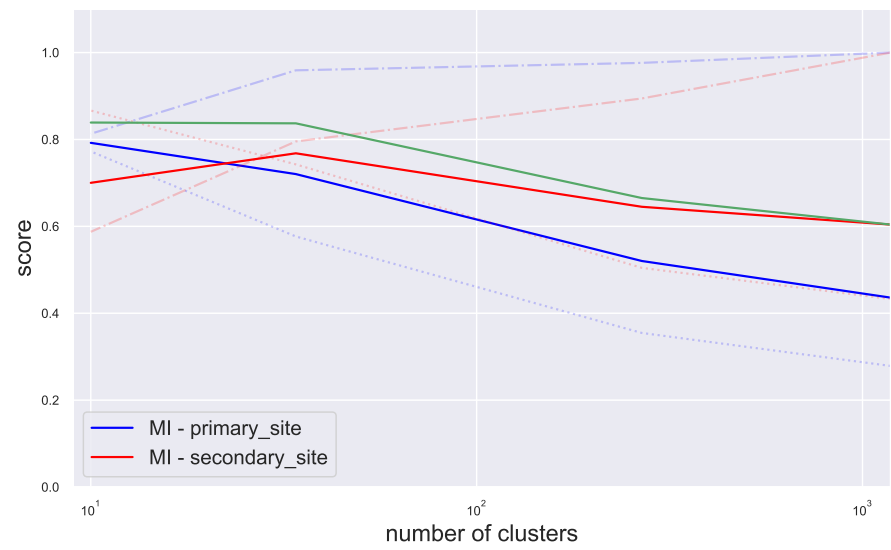


Figure 5.7: Scores across hierarchy

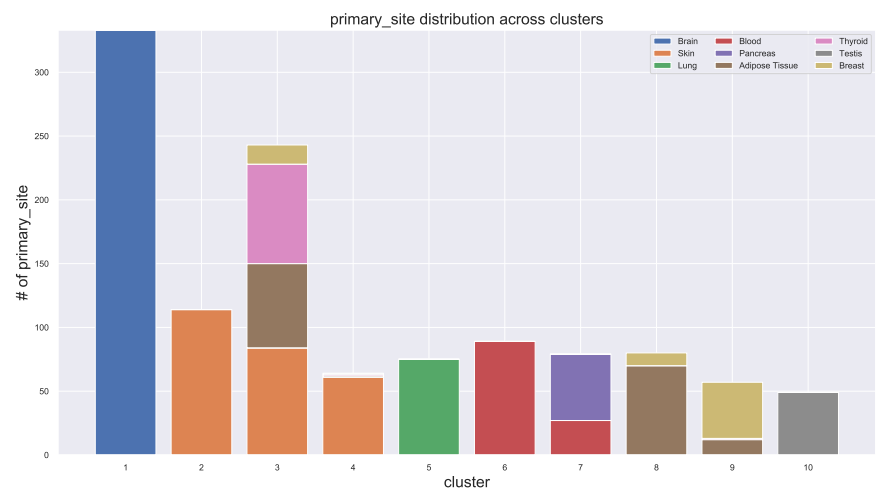


Figure 5.8: Caption

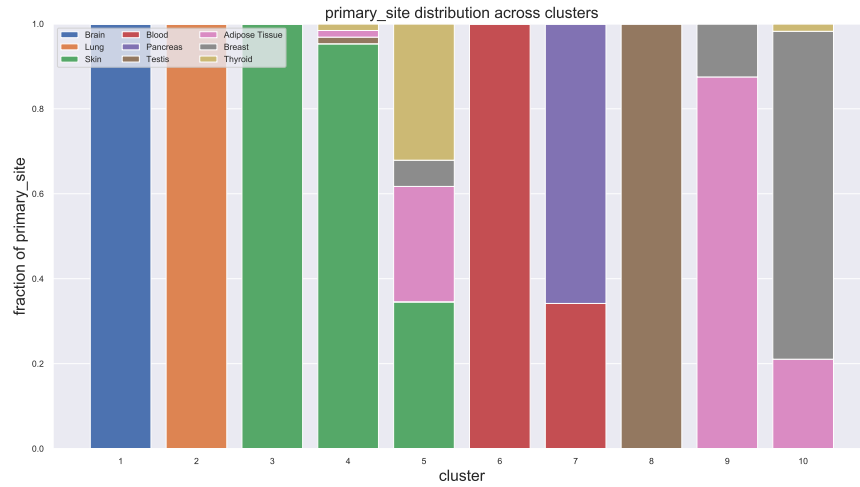


Figure 5.9: Caption

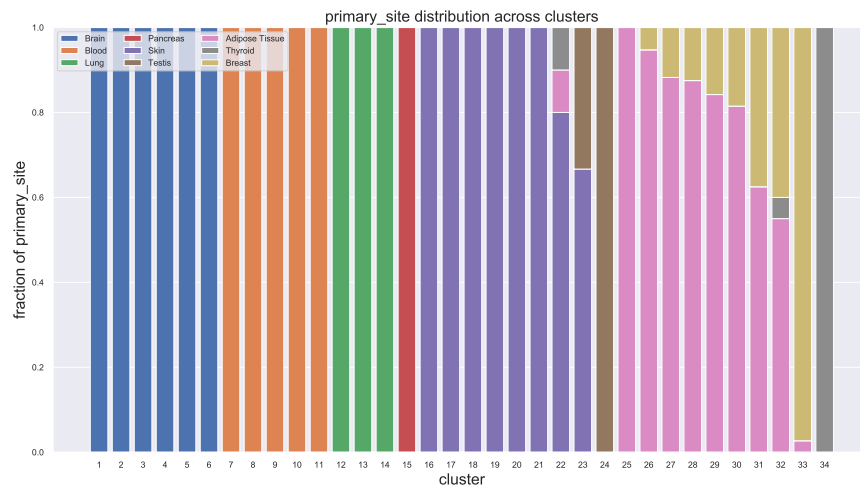


Figure 5.10: Caption

5.3.2 Run on TCGA

5.3.3 Mixed runs

hierarchical clustering This is better than LDA because...

Enrichment test are made once for each topic, starting from the layer with more genes per single topic. Test are made across multiple categories.

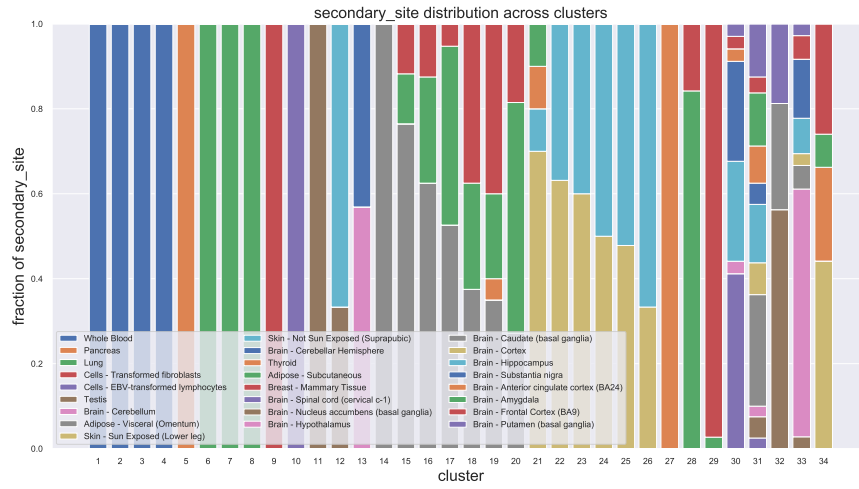


Figure 5.11: Caption

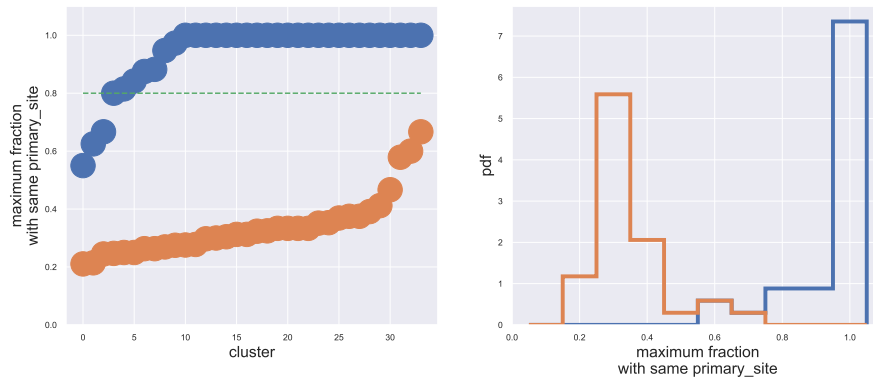


Figure 5.12: Caption

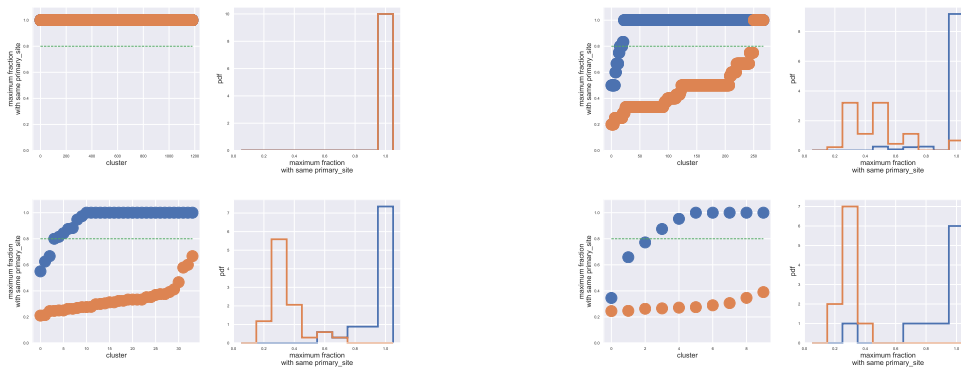
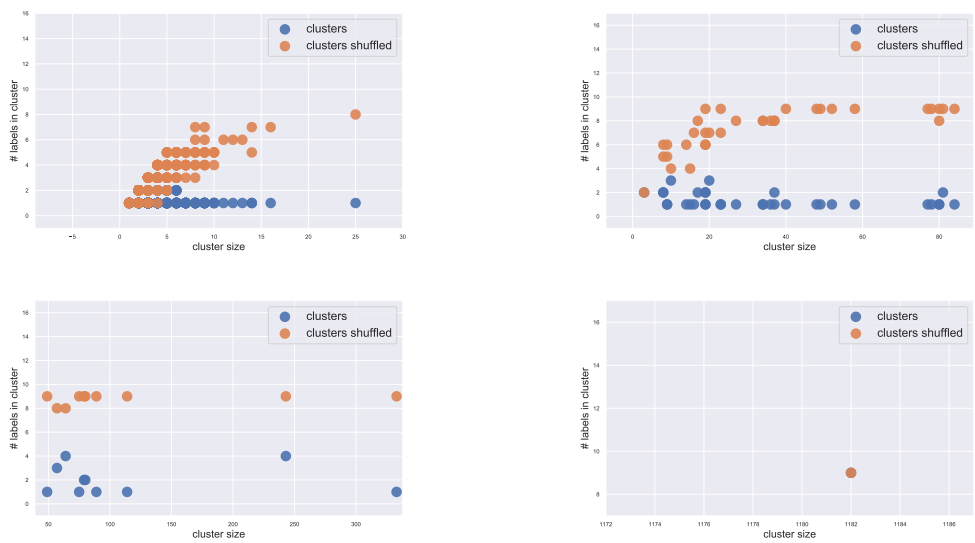




Figure 5.13: Caption



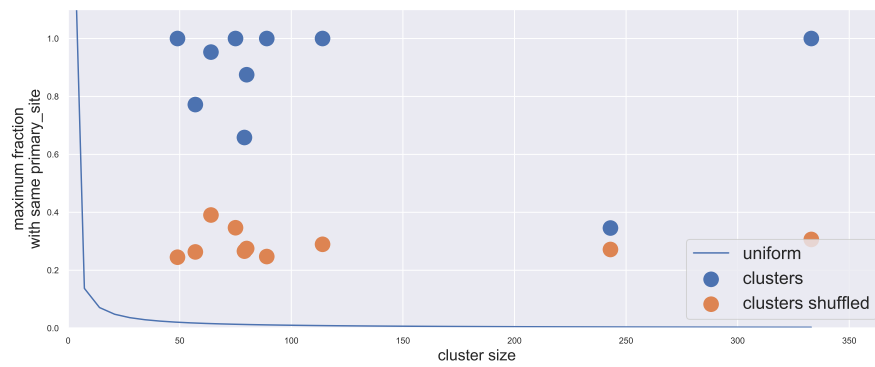
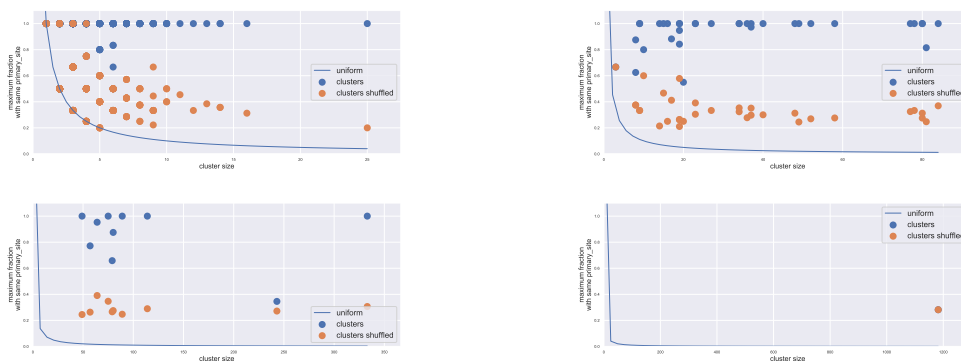


Figure 5.14: Caption



Chapter 6

Results

So discovered that
future: Loredana, Jacopo

Chapter 7

Conclusions

Finally...

Appendix A

hierarchical stochastic block model

The algorithm is called hierarchic Stochastic Block Model.

The first step of hSBM, as discussed in [?], is to create a bipartite network G with two kind of nodes: **words** and **documents**. Every time a word w is present in a document d an edge e_{wd} is created. If a word count in the entire corpus is under a certain threshold, that word is ignored. The aim is to find a partition $b \in \{b_i\}$ with $B = |\{b_i\}|$ blocks.

These kind of models are called *generative models*: given the data the model should generate a network G with probability $P(G|\theta, b)$, where b is the partition and θ any additional parameter of the model.

Using well-known Bayes theorem one could estimate the probability that an observed network is generated by partition b

$$P(b, \theta|G) = \frac{P(G|b, \theta) \overbrace{P(b, \theta)}^{prior}}{\underbrace{P(G)}_{\sum_{\theta} P(G|\theta, b) P(\theta, b)}} \quad (\text{A.1})$$

defining the amount of information needed to describe the data as the description length

$$\Sigma = -\ln P(G|b, \theta) - \ln P(b, \theta) \quad (\text{A.2})$$

the ?? can be written as $\frac{e^{-\Sigma}}{P(G)}$, so maximising that is equivalent to minimise the description length ???. The probability of obtaining a Graph from a set of parameters is $P(G|b, \theta) = \frac{1}{\Omega(A, \{n_r\})}$, where $\Omega(A, \{n_r\})$ is the number of graph that is possible to generate with audience matrix A and n_r the counts of block partition $\{b_i\}$

In case of a weighted network the likelihood becomes $P(G, x|b, \theta)$, where x are the weights.

A.0.1 Algorithm

First of all a $B \times B$ matrix is created. The entry e_{rs} of this matrix represents the number of links between nodes of group r and nodes of group s , with $r, s \in \{b_i\}$. At the beginning B groups are formed at random and the initial B is a hyper-parameter of the model.

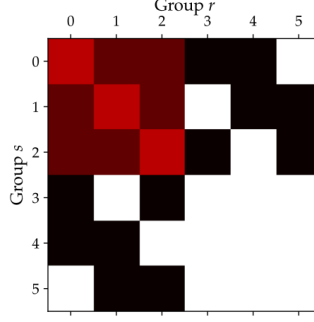


Figure A.1: Example of a edge's matrix from [?]

It is useful to define a traditional entropy:

$$S_t = \frac{1}{2} \sum_{r,s} n_r n_s H \left(\frac{e_{rs}}{n_r n_s} \right) \quad (\text{A.3})$$

where n_r is the number of nodes in groups r , e_{rs} is the number of edges between nodes of group r and nodes of group s , and $H(x) = -x \ln(x) - (1-x) \ln(1-x)$. This entropy is equivalent to the micro-canonical entropy of a system with $\Omega(A, \{n_r\})$ accessible states $S_t = L \ln \Omega$.

The algorithm uses a Markov Chain Monte Carlo to minimise this entropy. At each step a node changes block and the new configuration is accepted if S is decreased.

Note that ?? can be corrected taking care of degree distribution obtaining corrected entropy S_c

$$S_c = -\sum_{r,s} \frac{e_{rs}}{2} - \sum_k N_k \ln(k!) - \frac{1}{2} \sum_{r,s} e_{rs} \ln \left(\frac{e_{rs}}{e_r e_s} \right) \quad (\text{A.4})$$

How to change group of a node? At each step according to [?] node i can change group from r to s with a probability

$$P(r \rightarrow s|t) = \frac{e_{ts} + \epsilon}{e_t + \epsilon B} \quad (\text{A.5})$$

where j is a random neighbour of i : $j \in N_i$, $t \in \{b_j\}$ its block as defined in [?]. ϵ is a parameter that according to [?] has no significant impact in the algorithm, provided it is sufficiently small.

?? can be rewritten as

$$P(r \rightarrow s|t) = (1 - R_t) \frac{e_{ts}}{e_t} + \frac{R_t}{B}$$

defining $R_t = \frac{\epsilon B}{e_t + \epsilon B}$

This is done in four steps for each node i :

- a node j is chosen from i 's neighbours, the group of j is called t
- a random group s is selected

- move of node i to group s is accepted with probability R_t
- if s is not accepted, a random edge e is chosen from group t and node i is assigned to the endpoint of e which is not in t

This steps mime probability ??; note that for $\epsilon \rightarrow \infty$ this gives a uniform probability.

To enchant the probability to go into a minimum, a bounce of these moves is made, only the set of moves with the minimum S is accepted.

How many blocks B ? Note that the number of blocks B is a free parameter and must be inferred as described in [?]. This implies a slight modification of the algorithm such that it became possible to admit that a new group is created. When a group s is chosen, the algorithm can now accept a **new group** and ?? became

$$P(r \rightarrow s) = \Sigma_t P(t|i) \frac{e_{ts} + \epsilon}{e_t + \epsilon(B + 1)} \quad (\text{A.6})$$

being $P(t|i) = \Sigma_j \frac{A_{ij} \delta(b_j, t)}{k_i}$ the fraction of neighbours of i belonging to group t , e_t the number of edges in group t , k_i the degree, and b_j groups.

Using this modification it is now possible to add new groups and B is no longer a parameter.

How to find hierarchic layers? After the algorithm is run, one may would to add a new hierarchic level, this is done considering the B groups as nodes and repeating the process. As done before a matrix of edges like ?? is created, where edges are considered between groups of the previous layer.

The posterior probability became

$$P(\{b_l\}|A) = \frac{P(A|\{b_l\})P(\{b_l\})}{P(A)} = \prod_l^L P(b_l|e_l, b_{l-1}) \quad (\text{A.7})$$

where $l = 0 \dots L$ is the layer, A the audience matrix, b_i blocks. Note that $e_0 = A$ and $B_L = 1$. Maximising ?? gives the correct number of layers.

Adding a layer is done in 3 steps described in [?]:

Resize find $B_l \in [B_{l-1}, B_{l+1}]$ by bisection

Insert a layer l

Delete l and linking nodes from layer $l - 1$ directly to groups of layer $l + 1$

One marks initially all levels as not done and starts at the top level $l = L$ [?]. For the current level l , if it is marked done it is skipped and one moves to the level $l - 1$. Otherwise, all three moves are attempted. If any of the moves succeeds in decreasing the description length Σ ??, one marks the levels $l - 1$ and $l + 1$ (if they exist) as not done, the level l as done, and one proceeds (if possible) to the upper level $l + 1$, and repeats the procedure. If no improvement is possible, the level l is marked as done and one proceeds to the lower level $l - 1$. If the lowest level $l = 0$ is reached and cannot be improved, the algorithm ends.

Overlapping partitions As described in [?] one of the advantages of this approach is that it is possible to let a node belonging to multiple groups. In this case b_i becomes \vec{b}_i , with component $b_{ir} = 1$ if node i is in group r , 0 otherwise. The number of 1s in vector \vec{b}_i is called $d_i = |\vec{b}_i|$.

The probability of having a graph G being generated from an audience matrix A and a partition $\{\vec{b}_i\}$ is

$$P(G|A, \{\vec{b}_i\}) = \frac{1}{\Omega}$$

if Ω is the number of possible graphs. Entropy ?? is $S_t = Ln\Omega$. This corresponds to an augmented graph generated via a non overlapping block model with $N' = \sum_r n_r > N$ nodes and the same audience matrix A .

First of all, it is necessary to sample the distribution of mixture sizes $P(\{n_d\})$ where n_d is the number of nodes which mixture has got size d , $n_d \in [0, N]$ and $d \in [0, D]$ (typically $D = B$ and in the non-overlapping case $D = 1$), this is done by sampling uniformly from

$$P(\{n_d\}|B) = \left(\binom{D}{N} \right)^{-1}$$

which is probability of having n nodes whose mixture has size d . $\binom{B}{N}$ is the number of histograms with area N and B distinguishable bins. $B - 1$ can be used instead of B to avoid node with no group, in this case $d \in [1, B]$.

Given the mixture sizes, the distribution of node membership is sampled from

$$P(\{d_i\}|\{n_d\}) = \frac{\prod_d n_d!}{N!}$$

.

At this point for each set of nodes with $d_i = d$ it is necessary to sample $n_{\vec{b}}$; the number of nodes with a particular mixture \vec{b} . It is sampled from

$$P(\{n_{\vec{b}}\}_d|n_d) = \left(\binom{D}{n_d} \right)^{-1}, \quad (\text{A.8})$$

next all mixtures \vec{b}_i of size d must be sampled, they are given by

$$P(\{\vec{b}_i\}_d|\{n_{\vec{b}}\}_d) = \frac{\prod_{|\vec{b}_i|=d} n_{\vec{b}}!}{n_d!} \quad (\text{A.9})$$

the global posterior as defined in [?] is

$$P(\{\vec{b}_i\}|B) = \left[\prod_{d=1}^B P(\{\vec{b}_i\}_d|\{n_{\vec{b}}\}_d) P(\{n_{\vec{b}}\}_d|n_d) \right] P(d_i|n_d) P(n_d|B) \quad (\text{A.10})$$

At this time it is necessary to obtain the distribution of the edges between mixtures. Defined $e_r = \sum_s e_{rs}$ the number of half-edges labelled r , $m_r = \sum_{\vec{b}} b_r$ the number of mixtures containing group r the algorithm samples the probability distribution of the edges count

$$P(\{e_{\vec{b}}\}|\{\vec{b}_i\}, A) = \prod_r \left(\binom{m_r}{e_r} \right)^{-1}$$

and the labelled degree sequence $\{\vec{k}_i\}$ from

$$P(\{\vec{k}_i\}_{\vec{b}}|\{e_{\vec{b}}\}, \{\vec{b}_i\}) = \frac{\prod_k n_k^{\vec{b}_k!}}{n_{\vec{b}}!}$$

Word documents separation Following what is done in [?], the probability of a group $P(b_l)$ at a certain level l is intended as the disjoint probability of group of words and group of documents.

$$P(b_l) = P_w(b_l^w)P_d(b_l^d) \tag{A.11}$$

Doing this let words and documents be separated by construction. Considering the process described above if two nodes are not connected at the beginning it is impossible that they end up in the same block. It is easily verified in [?] that this property is preserved and fully reflected in the final block structure.

List of Figures

List of Tables

Acknowledgements

Ringrazio bla bla bla...