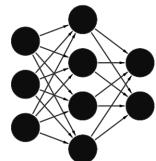


PREDICTION OF PATHOGENIC SNV

Prof. Giorgio Valentini
6 CFU

Luca Cappelletti

Course Project
Year 2017/2018



IT Master Degree
Universiy of Milan
Italy
29 giugno 2018

Indice

I Dataset	4
1 Data points	5
1.1 Retrieving the dataset	5
1.2 Composition	5
1.2.1 Training dataset	5
1.2.2 Testing dataset	5
2 Metrics	6
2.1 How the graphs are realized	6
2.1.1 Metric sample distribution	6
2.1.2 Plot graphs	6
2.1.3 Normalized plot graphs	6
2.2 CpGobsExp	7
2.2.1 Metric sample distribution	7
2.2.2 Metric values	7
2.3 CpGperCpG	8
2.3.1 Metric sample distribution	8
2.3.2 Metric values	8
2.4 CpGperGC	9
2.4.1 Metric sample distribution	9
2.4.2 Metric values	9
2.5 DGVCount	10
2.5.1 Metric sample distribution	10
2.5.2 Metric values	10
2.6 DnaseClusteredHyp	11
2.6.1 Metric sample distribution	11
2.6.2 Metric values	11
2.7 DnaseClusteredScore	12
2.7.1 Metric sample distribution	12
2.7.2 Metric values	12
2.8 EncH3K27Ac	13
2.8.1 Metric sample distribution	13
2.8.2 Metric values	13
2.9 EncH3K4Me1	14
2.9.1 Metric sample distribution	14
2.9.2 Metric values	14
2.10 EncH3K4Me3	15
2.10.1 Metric sample distribution	15
2.10.2 Metric values	15
2.11 GCContent	16
2.11.1 Metric sample distribution	16
2.11.2 Metric values	16
2.12 GerpRS	17
2.12.1 Metric sample distribution	17
2.12.2 Metric values	17
2.13 GerpRSpv	18
2.13.1 Metric sample distribution	18

2.13.2 Metric values	18
2.14 ISCApath	19
2.14.1 Metric sample distribution	19
2.14.2 Metric values	19
2.15 commonVar	20
2.15.1 Metric sample distribution	20
2.15.2 Metric values	20
2.16 dbVARCount	21
2.16.1 Metric sample distribution	21
2.16.2 Metric values	21
2.17 fantom5Perm	22
2.17.1 Metric sample distribution	22
2.17.2 Metric values	22
2.18 fantom5Robust	23
2.18.1 Metric sample distribution	23
2.18.2 Metric values	23
2.19 fracRareCommon	24
2.19.1 Metric sample distribution	24
2.19.2 Metric values	24
2.20 mamPhastCons46way	25
2.20.1 Metric sample distribution	25
2.20.2 Metric values	25
2.21 mamPhyloP46way	26
2.21.1 Metric sample distribution	26
2.21.2 Metric values	26
2.22 numTFBSConserved	27
2.22.1 Metric sample distribution	27
2.22.2 Metric values	27
2.23 priPhastCons46way	28
2.23.1 Metric sample distribution	28
2.23.2 Metric values	28
2.24 priPhyloP46way	29
2.24.1 Metric sample distribution	29
2.24.2 Metric values	29
2.25 rareVar	30
2.25.1 Metric sample distribution	30
2.25.2 Metric values	30
2.26 verPhastCons46way	31
2.26.1 Metric sample distribution	31
2.26.2 Metric values	31
2.27 verPhyloP46way	32
2.27.1 Metric sample distribution	32
2.27.2 Metric values	32
3 Metric distribution summary	33
4 Data correlation	34
4.1 Scatter plot	34
4.2 Correlation coefficient matrix	35
4.2.1 CpGobsExp and CpGperCpG	36
4.2.2 CpGobsExp and CpGperGC	36
4.2.3 CpGperCpG and CpGperGC	37
4.2.4 dbVARCount and DGVCount	37
4.2.5 mamPhyloP46way and verPhyloP46way	38
4.2.6 DnaseClusteredHyp and DnaseClusteredScore	38
4.2.7 mamPhastCons46way and verPhastCons46way	39
4.3 Identified data correlations	39
4.4 Correlation table after removing highly correlated data	40
5 Dataset visualization	41
5.1 PCA	41
5.1.1 Training dataset visualization	41

5.1.2	Testing dataset visualization	41
5.1.3	Mixed dataset visualization	42
6	Dataset issues	43
6.1	Possible dataset errors	43
6.2	Biased testing dataset	43
II	Network implementation	44
7	Model architecture	45
7.1	Input	45
7.2	Output	45
7.3	Weight distribution based on input distribution	45
7.4	Weight distribution based on activation functions and regularization layers	45
7.5	Batch Size	45
7.6	Hidden layers	45
7.7	Possibility: Locally connected dense layers	45
7.8	Activation function	46
7.9	Regularization	46
7.10	Drop out	47
7.11	Loss function	47
7.12	Update policy	47
7.13	Network model representation	47
8	References	48

Parte I

Dataset

Data points

First we begin looking at the dataset, the distributions of the given metrics and the statistical analysis of these data points.

1.1 Retrieving the dataset

The dataset can be downloaded from <https://homes.di.unimi.it/valentini/ProgettoBioinformatica1718/data/>.

1.2 Composition

1.2.1 Training dataset

In the training dataset there are 981388 data points, each one comprised of 26 metrics. The first 356 are pathogenic and all the others are negative.

1.2.2 Testing dataset

In the test dataset there are 19018 data points, still each one comprised of 26 metrics. The first 40 are pathogenic and the following are negative.

2

Metrics

2.1 How the graphs are realized

All the graphs are in triples: positives, negatives and mixed.

The normalization is done, as usual, in the following way:

$$m' = \frac{\text{metric} - \mathbb{E}(\text{metric values})}{\max\{\text{metric values}\} - \min\{\text{metric values}\}}$$

Figura 2.1: Input normalization

2.1.1 Metric sample distribution

Are realized by calculating the frequencies and estimating the density distributions parameters via MLE.

2.1.2 Plot graphs

Plot graphs are realized by sorting the values of the single metrics.

2.1.3 Normalized plot graphs

Are realized by sorting the values of the metric, with the domain and codomain normalized.

2.2 CpGobsExp

2.2.1 Metric sample distribution

The data points seem to follow a **Beta** distribution.

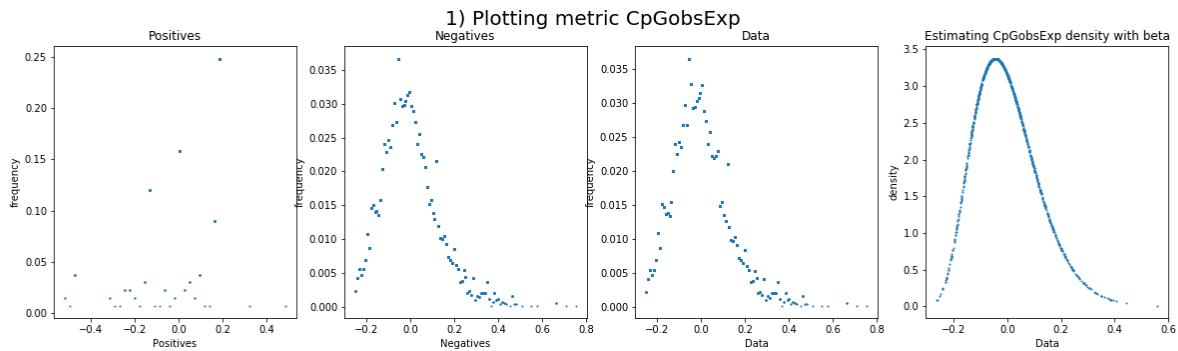


Figura 2.2: Sampling distribution of metric CpGobsExp

2.2.2 Metric values

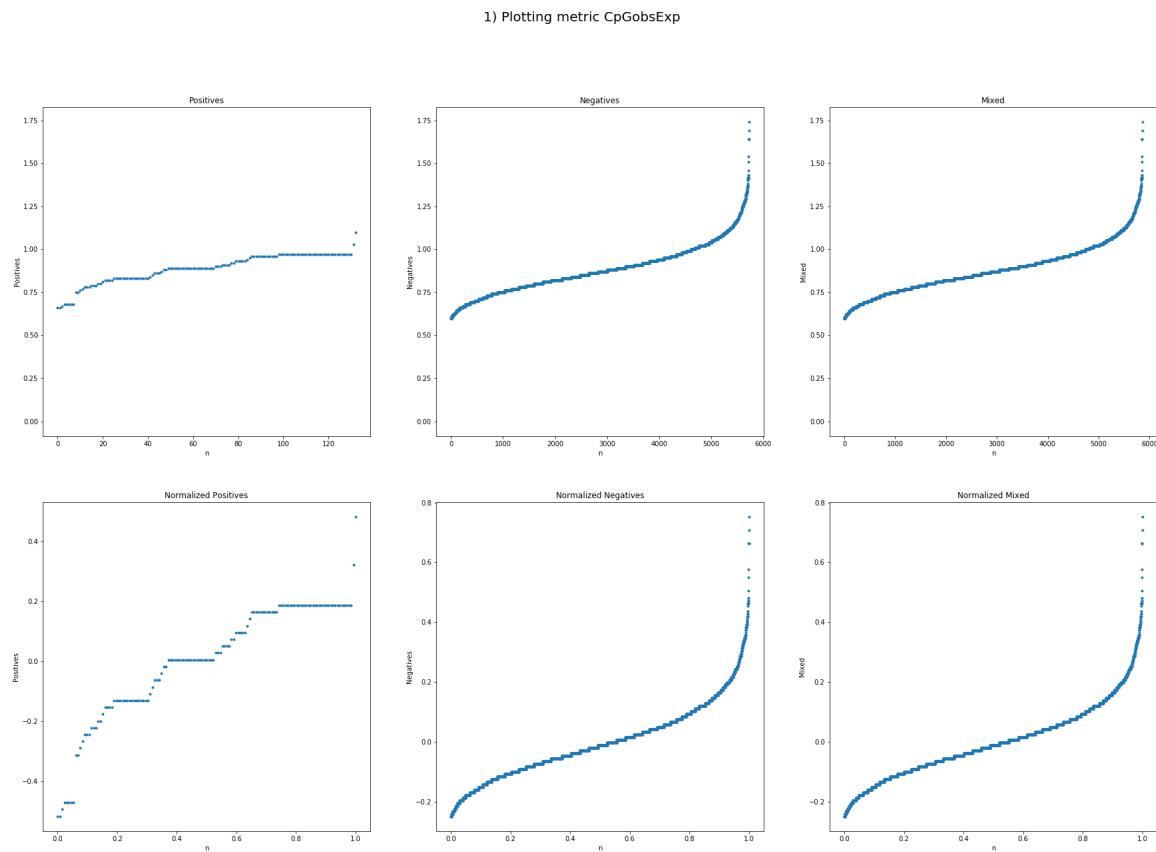


Figura 2.3: Values of metric CpGobsExp

2.3 CpGperCpG

2.3.1 Metric sample distribution

The data points seem to follow a **Beta** distribution.

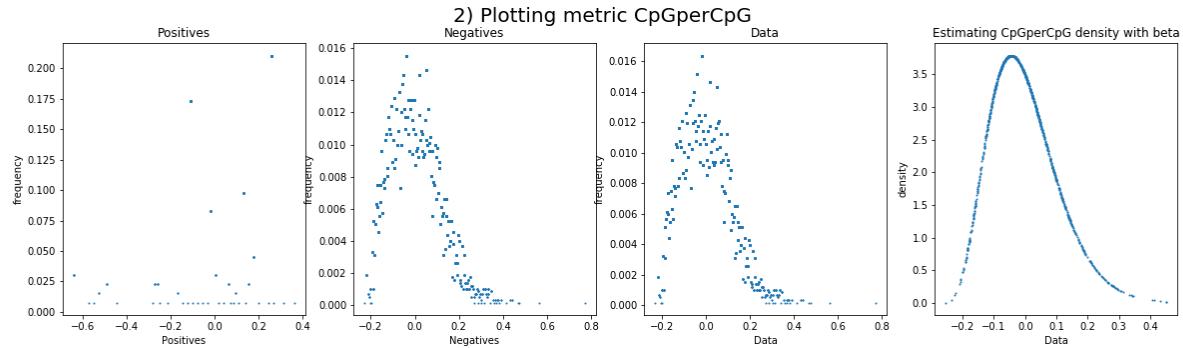


Figura 2.4: Sampling distribution of metric CpGperCpG

2.3.2 Metric values

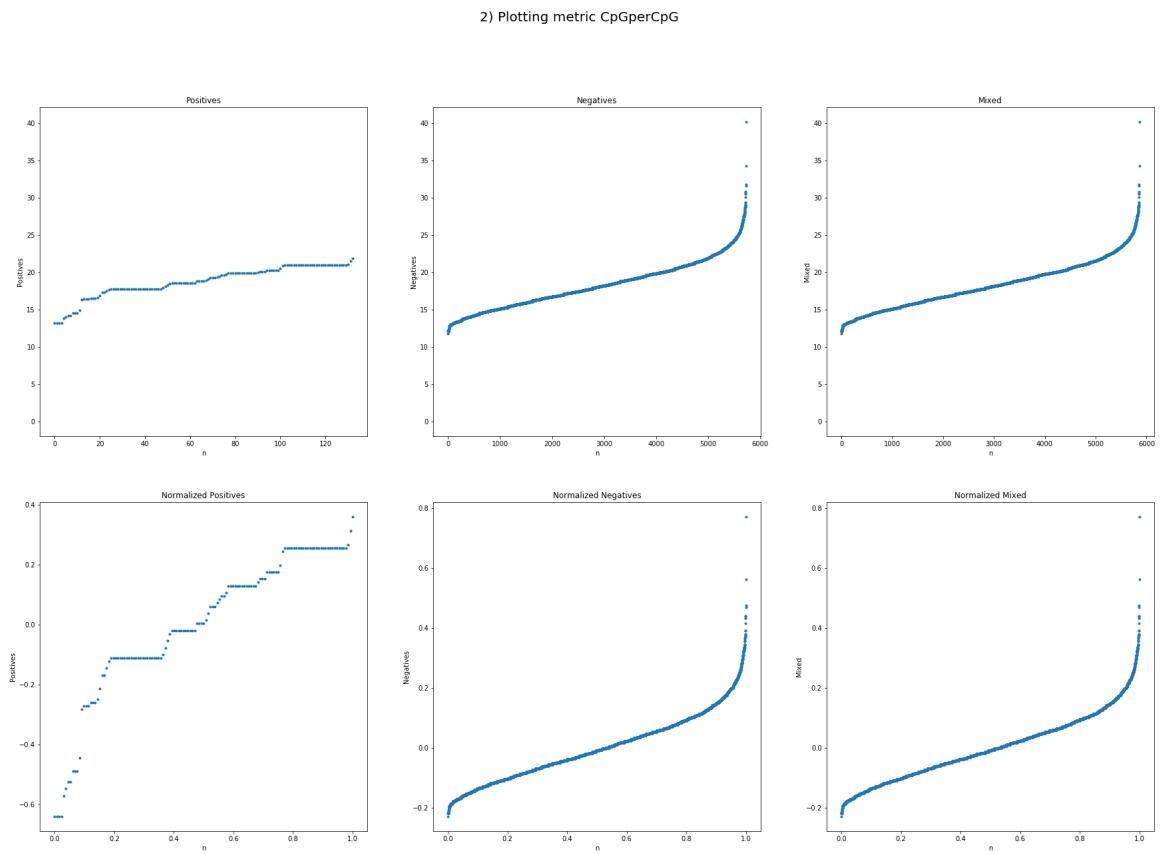


Figura 2.5: Values of metric CpGperCpG

2.4 CpGperGC

2.4.1 Metric sample distribution

The data points seem to follow a **Gaussian** distribution.

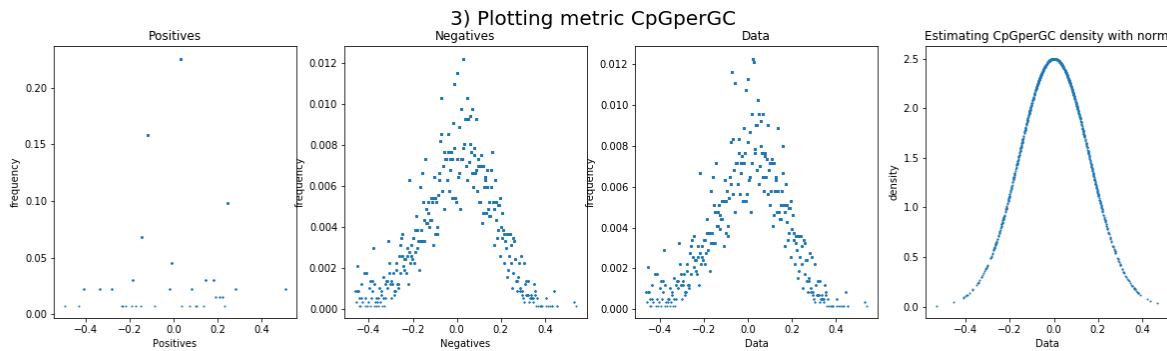


Figura 2.6: Sampling distribution of metric CpGperGC

2.4.2 Metric values

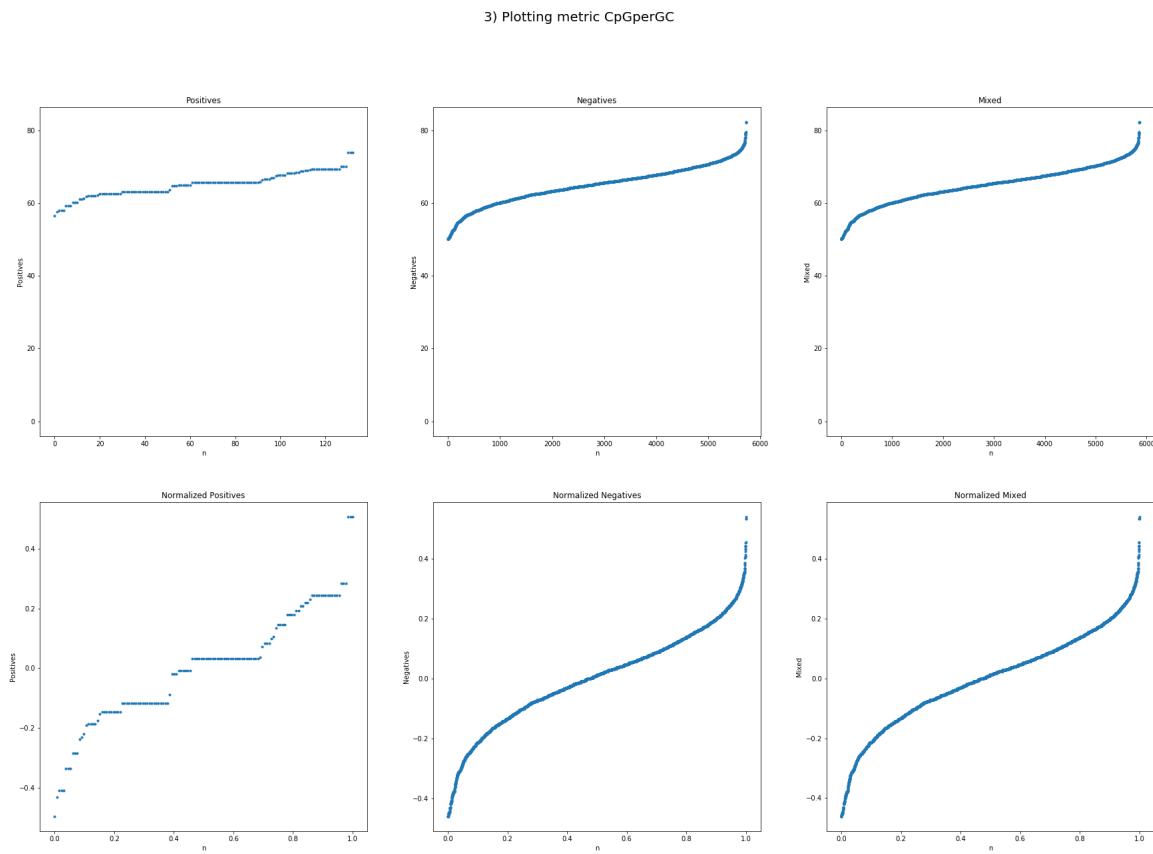


Figura 2.7: Values of metric CpGperGC

2.5 DGVCount

2.5.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

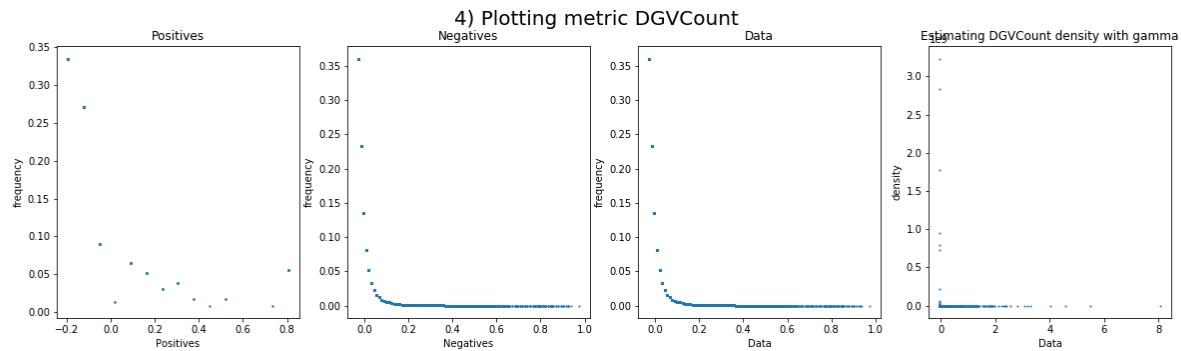


Figura 2.8: Sampling distribution of metric DGVCount

2.5.2 Metric values

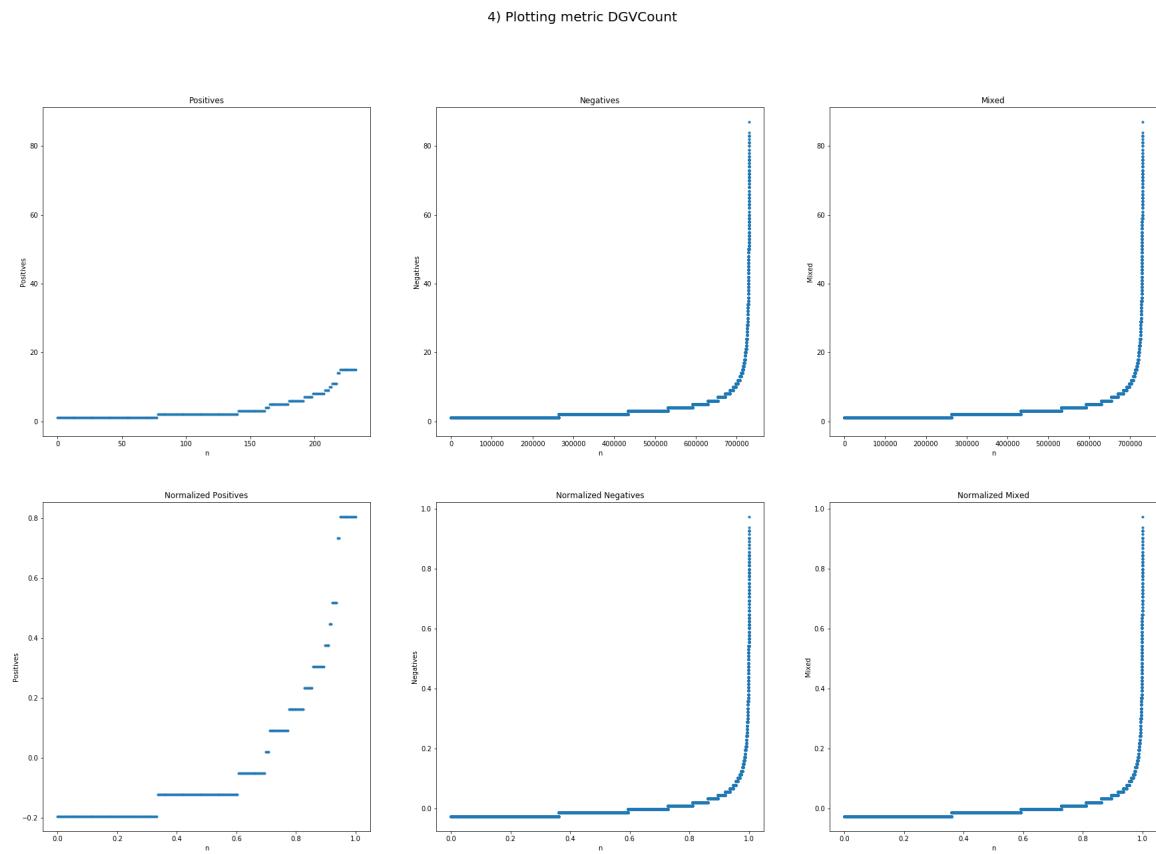


Figura 2.9: Values of metric DGVCount

2.6 DnaseClusteredHyp

2.6.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

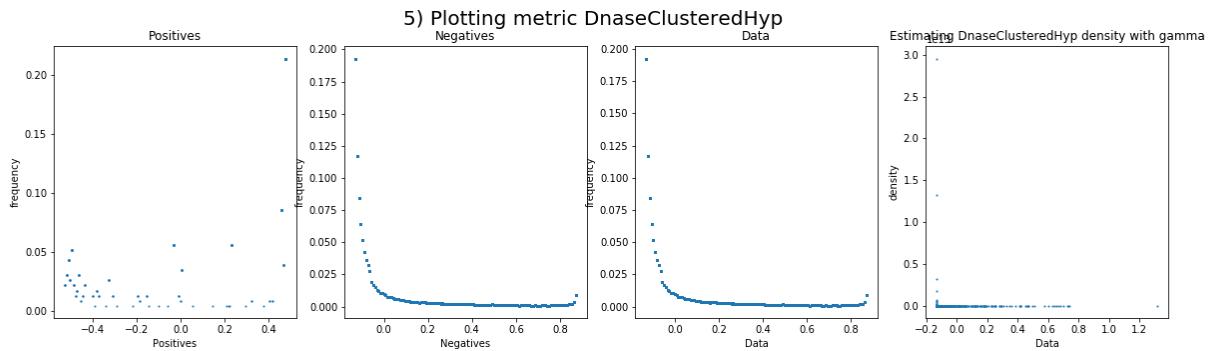


Figura 2.10: Sampling distribution of metric DnaseClusteredHyp

2.6.2 Metric values

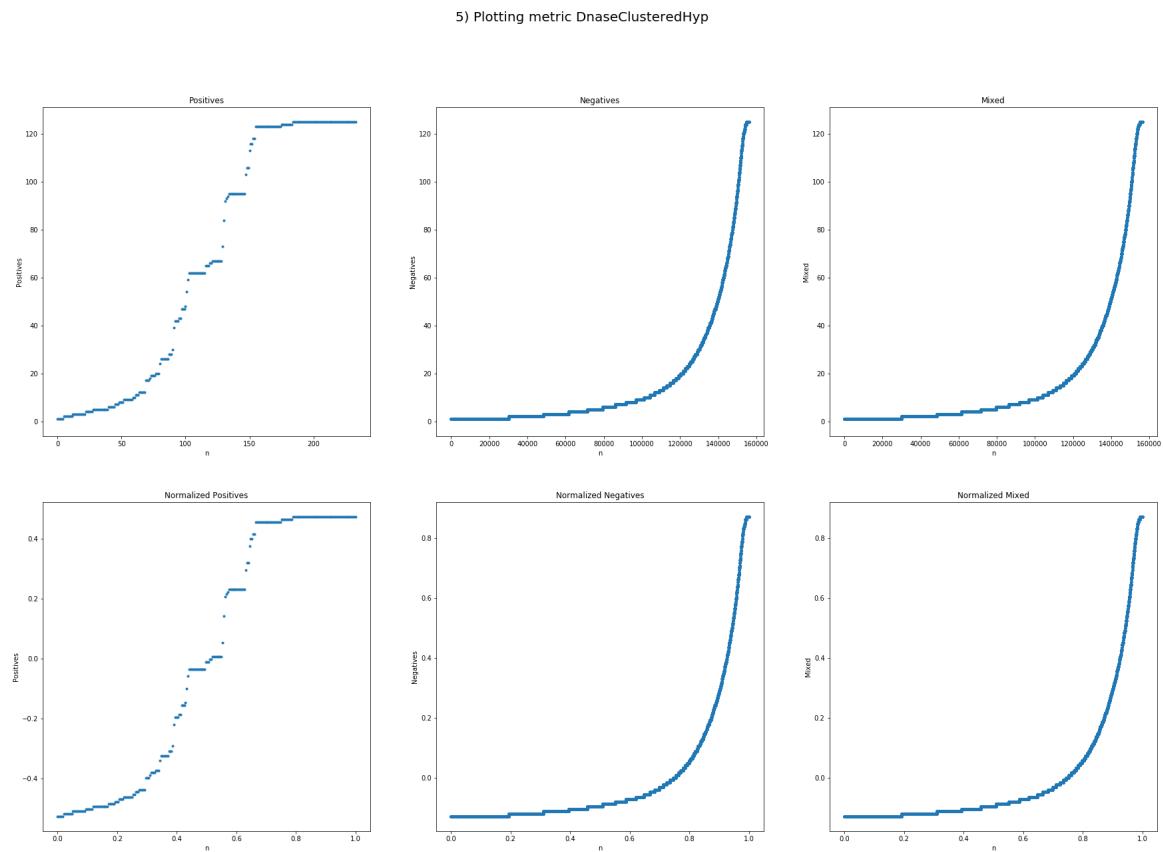


Figura 2.11: Values of metric DnaseClusteredHyp

2.7 DnaseClusteredScore

2.7.1 Metric sample distribution

The data points seem to follow **slightly** a **Beta** distribution.

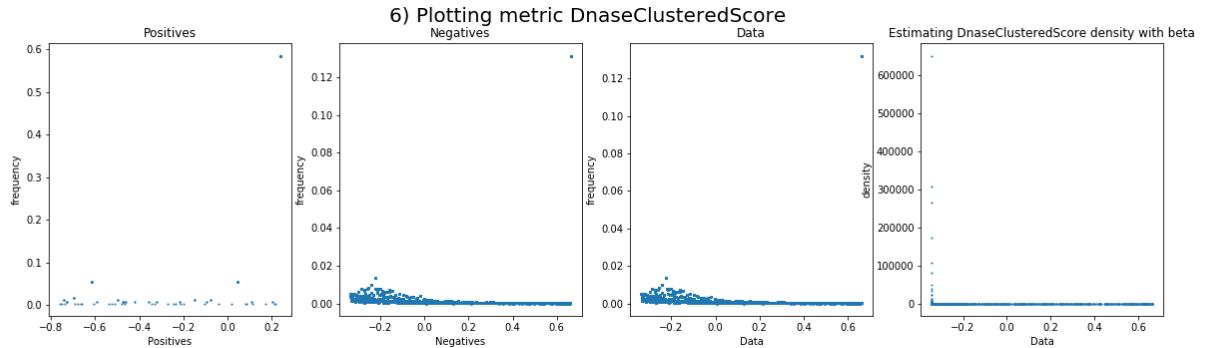


Figura 2.12: Sampling distribution of metric DnaseClusteredScore

2.7.2 Metric values

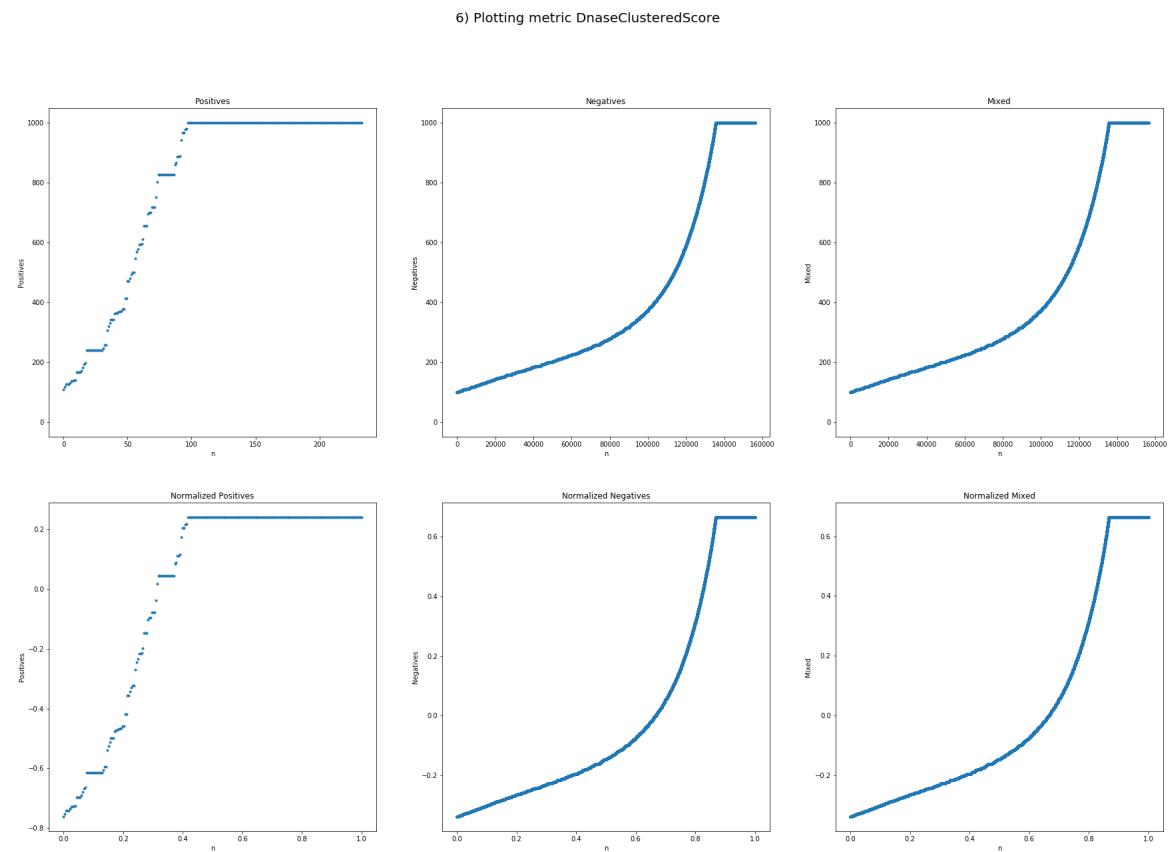


Figura 2.13: Values of metric DnaseClusteredScore

2.8 EncH3K27Ac

2.8.1 Metric sample distribution

The data points seem to follow a family of **Gamma** distributions (a speculation for this distribution could be the different groups from which the data are extracted), we will approximate them to one with a linear combination of the parameters.

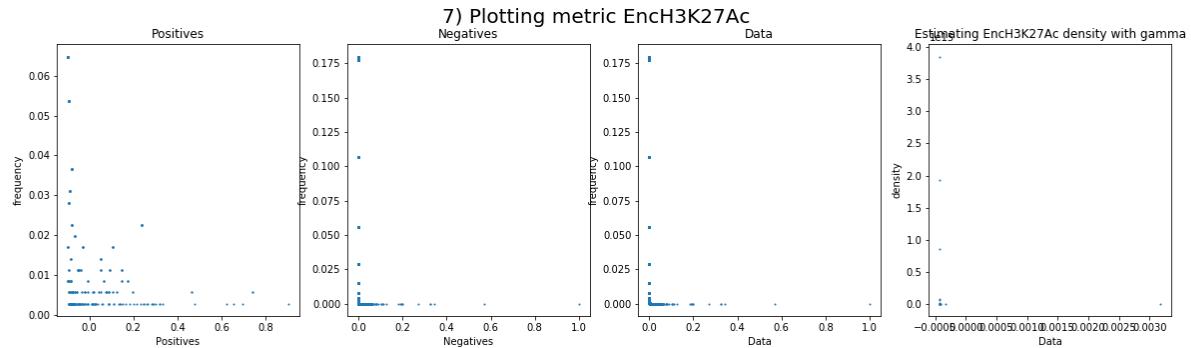


Figura 2.14: Sampling distribution of metric EncH3K27Ac

2.8.2 Metric values

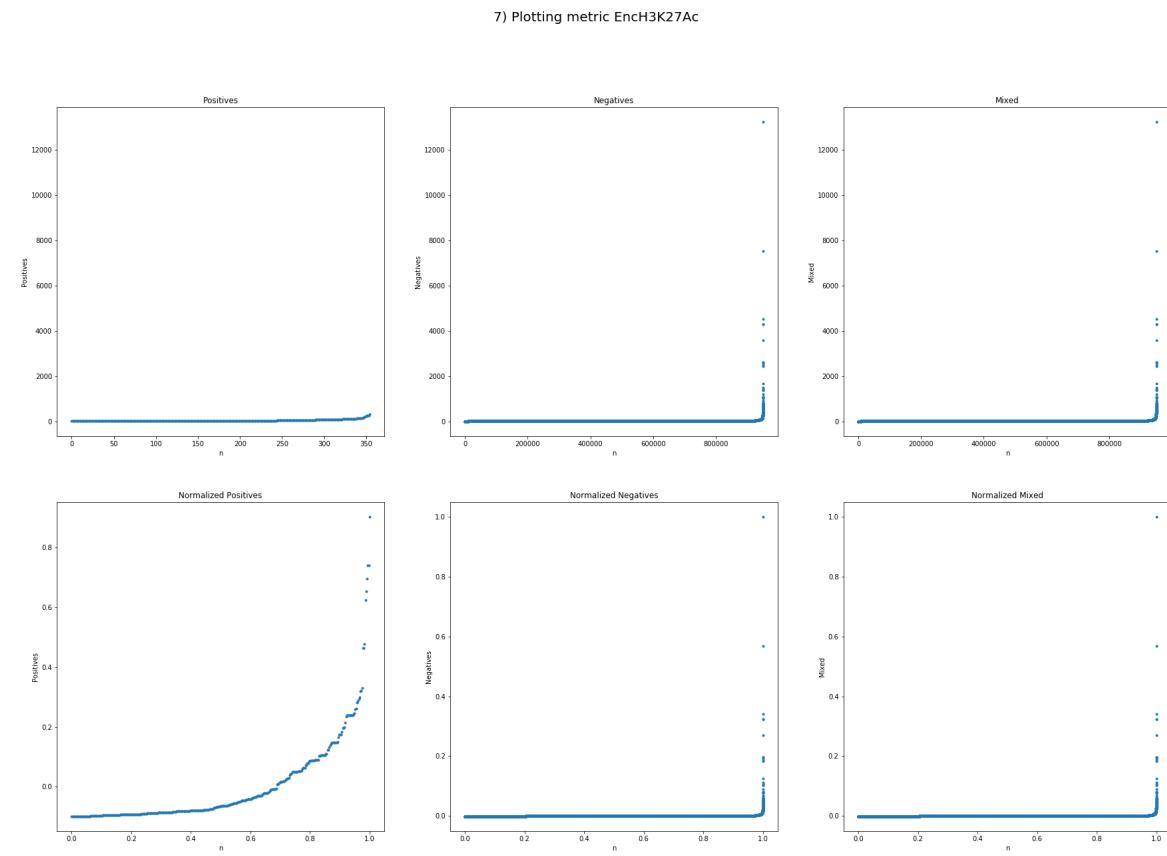


Figura 2.15: Values of metric EncH3K27Ac

2.9 EncH3K4Me1

2.9.1 Metric sample distribution

The data points seem to follow a family of **Gamma** distributions (a speculation for this distribution could be the different groups from which the data are extracted), we will approximate them to one with a linear combination of the parameters.

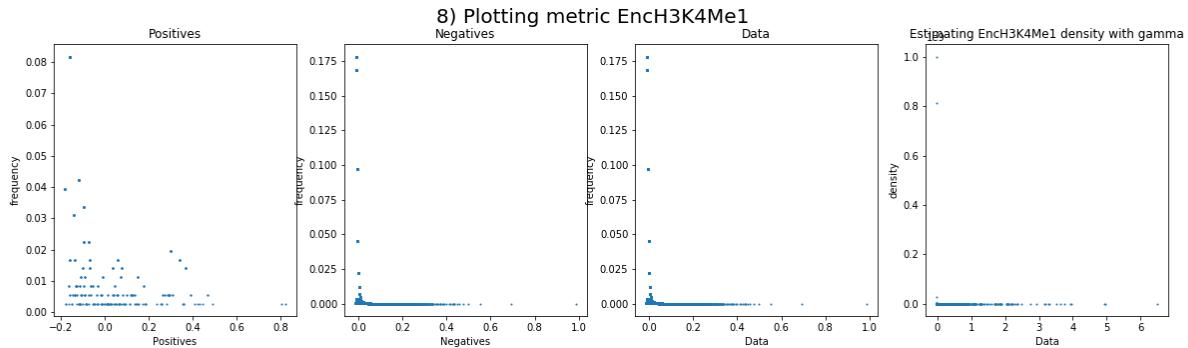


Figura 2.16: Sampling distribution of metric EncH3K4Me1

2.9.2 Metric values

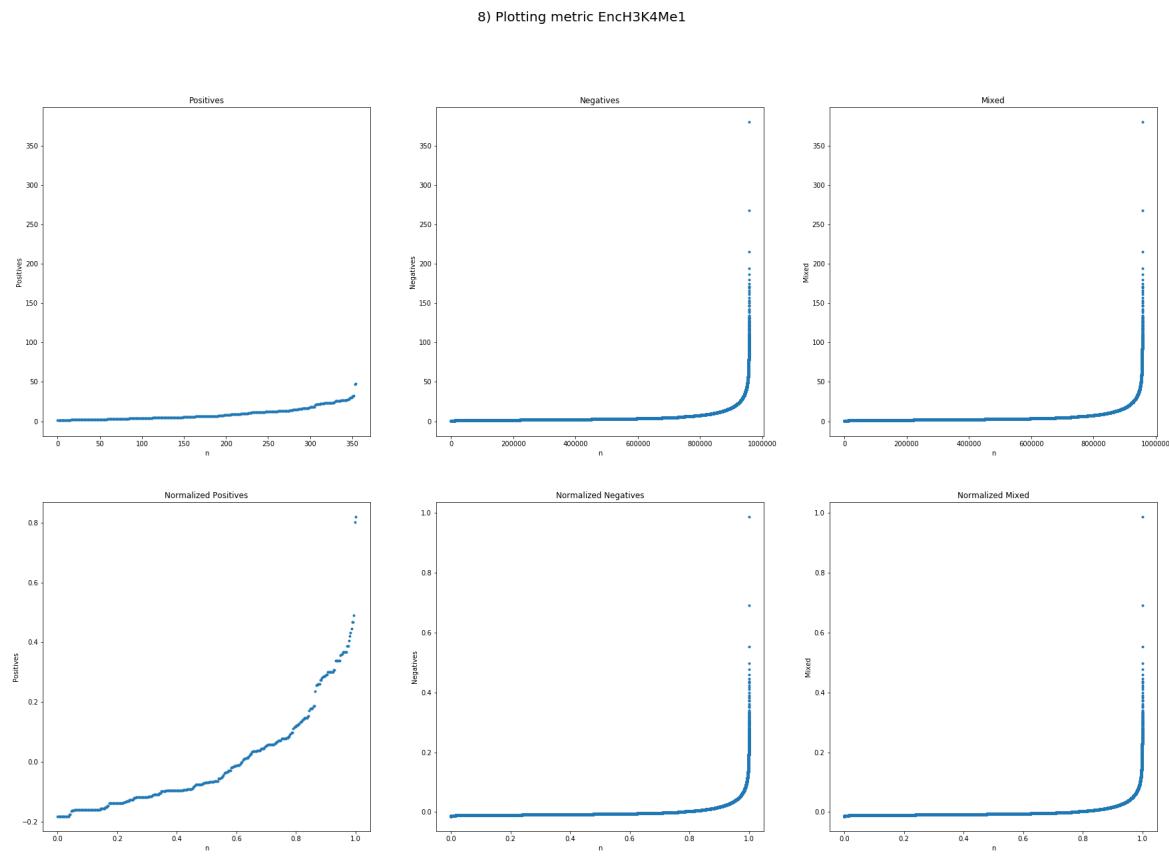


Figura 2.17: Values of metric EncH3K4Me1

2.10 EncH3K4Me3

2.10.1 Metric sample distribution

The data points seem to follow a family of **Gamma** distributions (a speculation for this distribution could be the different groups from which the data are extracted), we will approximate them to one with a linear combination of the parameters.

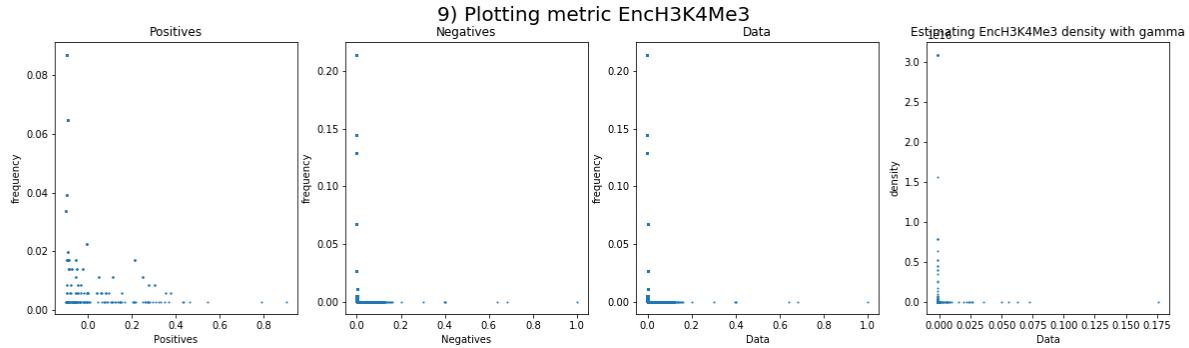


Figura 2.18: Sampling distribution of metric EncH3K4Me3

2.10.2 Metric values

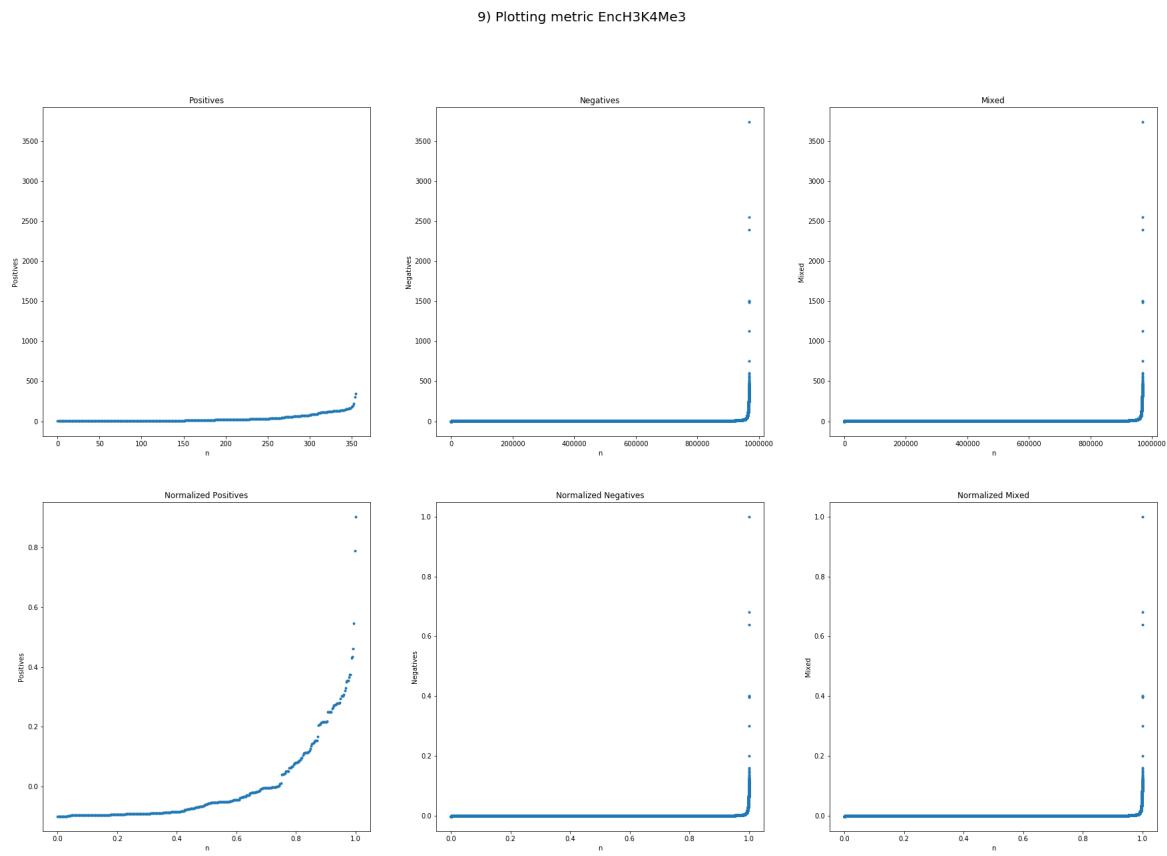


Figura 2.19: Values of metric EncH3K4Me3

2.11 GCContent

2.11.1 Metric sample distribution

The data points seem to be a combination of two **Gaussian** distributions.

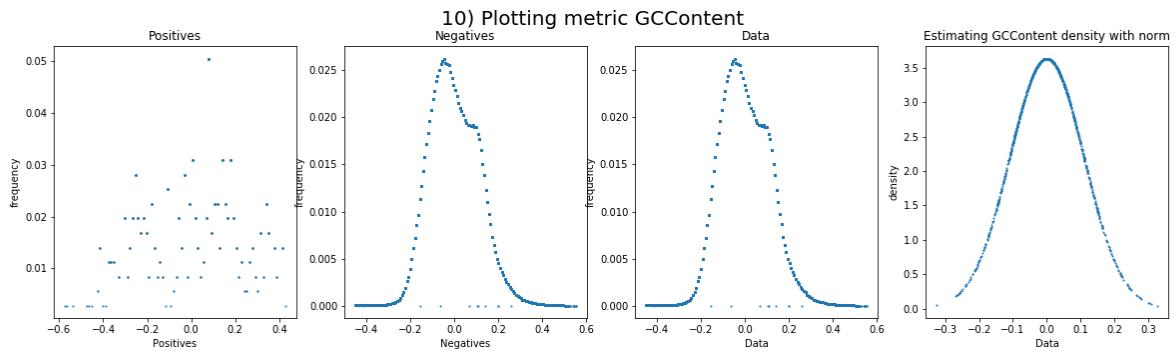


Figura 2.20: Sampling distribution of metric GCContent

2.11.2 Metric values

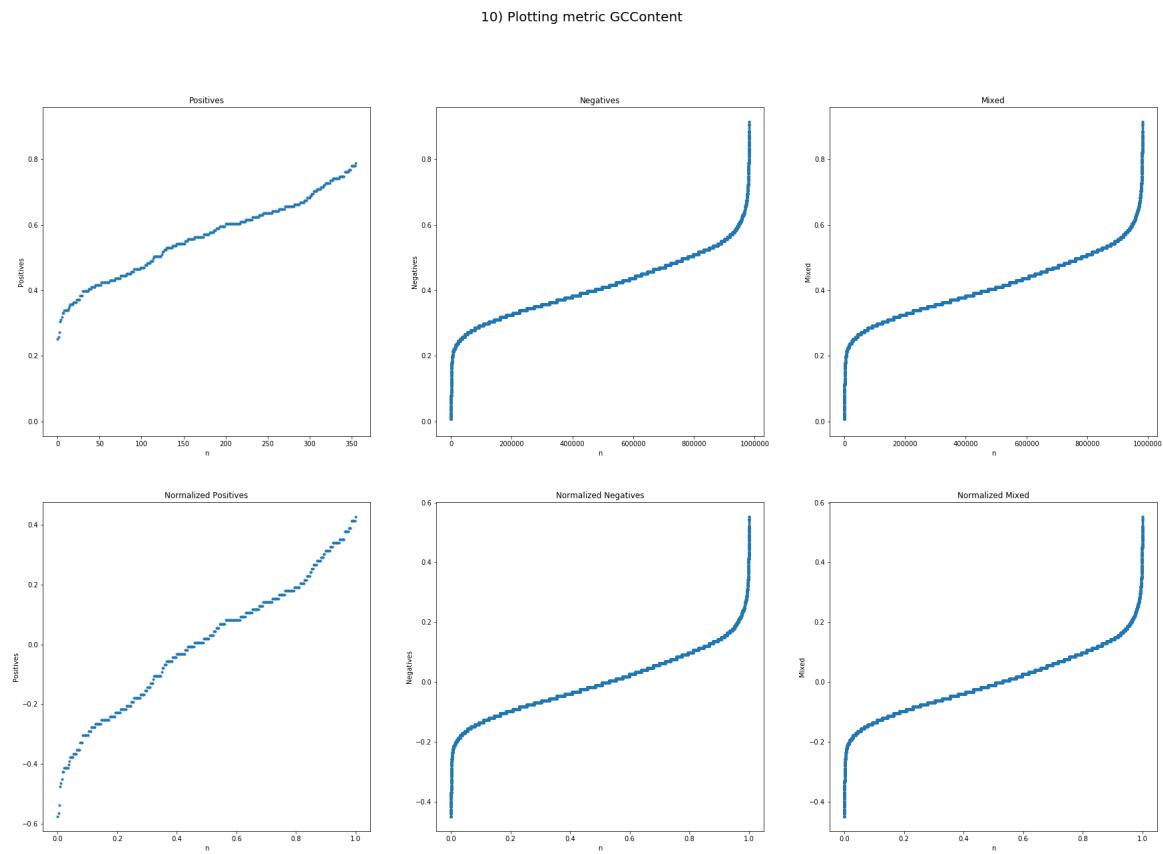


Figura 2.21: Values of metric GCContent

2.12 GerpRS

2.12.1 Metric sample distribution

The data points seem to follow a family of **Gamma** distributions (a speculation for this distribution could be the different groups from which the data are extracted), we will approximate them to one with a linear combination of the parameters.

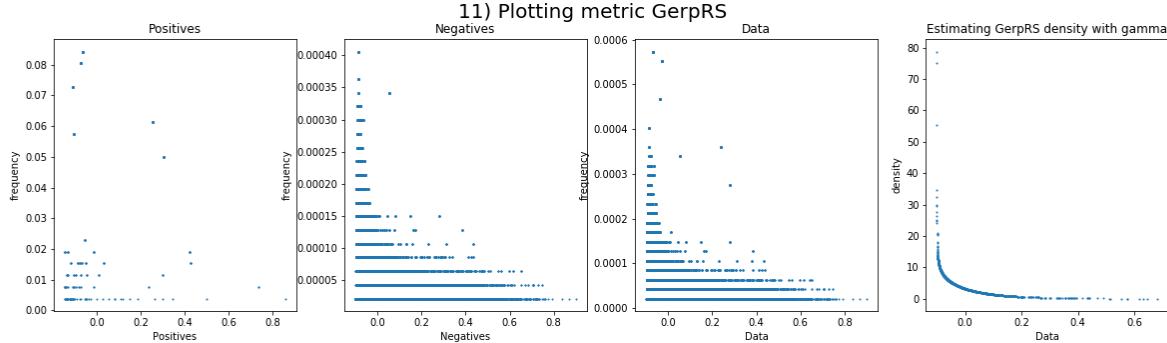


Figura 2.22: Sampling distribution of metric GerpRS

2.12.2 Metric values

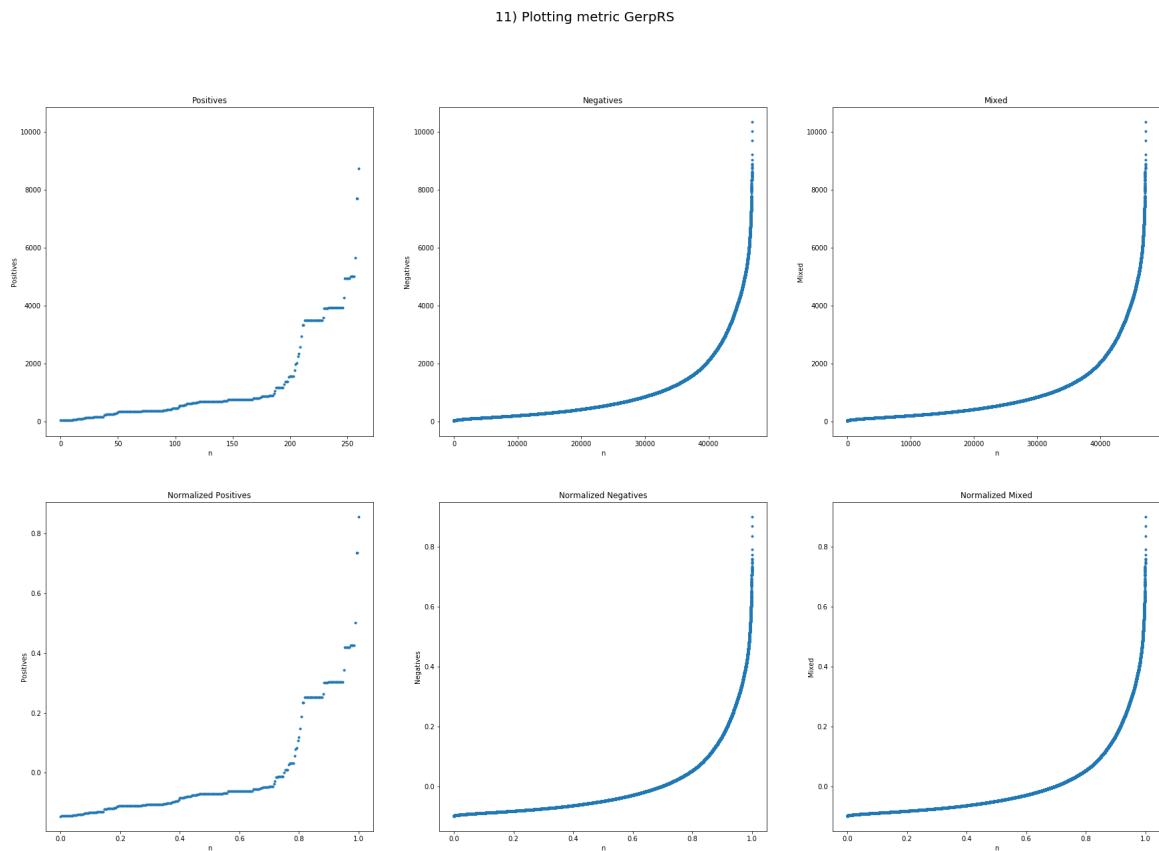


Figura 2.23: Values of metric GerpRS

2.13 GerpRSpv

2.13.1 Metric sample distribution

The data points seem to follow a family of **Gamma** distributions (a speculation for this distribution could be the different groups from which the data are extracted), we will approximate them to one with a linear combination of the parameters.

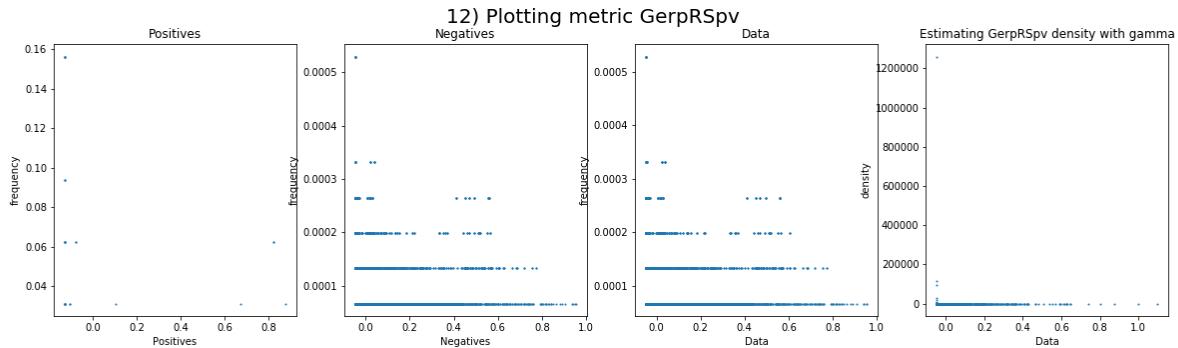


Figura 2.24: Sampling distribution of metric GerpRSpv

2.13.2 Metric values

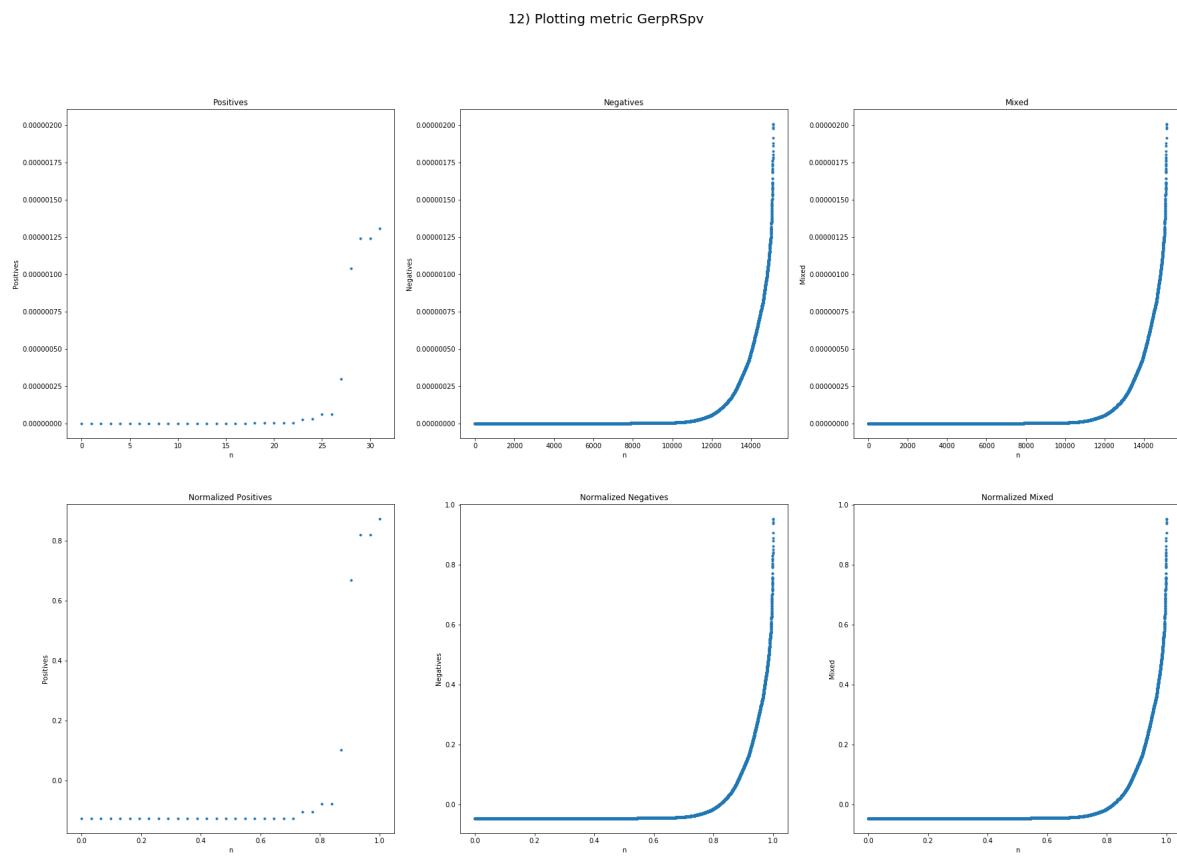


Figura 2.25: Values of metric GerpRSpv

2.14 ISCApath

2.14.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

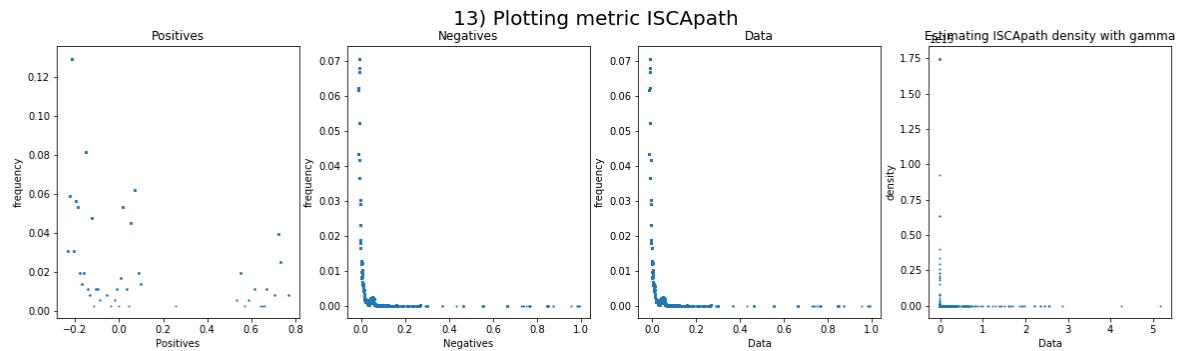


Figura 2.26: Sampling distribution of metric ISCApath

2.14.2 Metric values

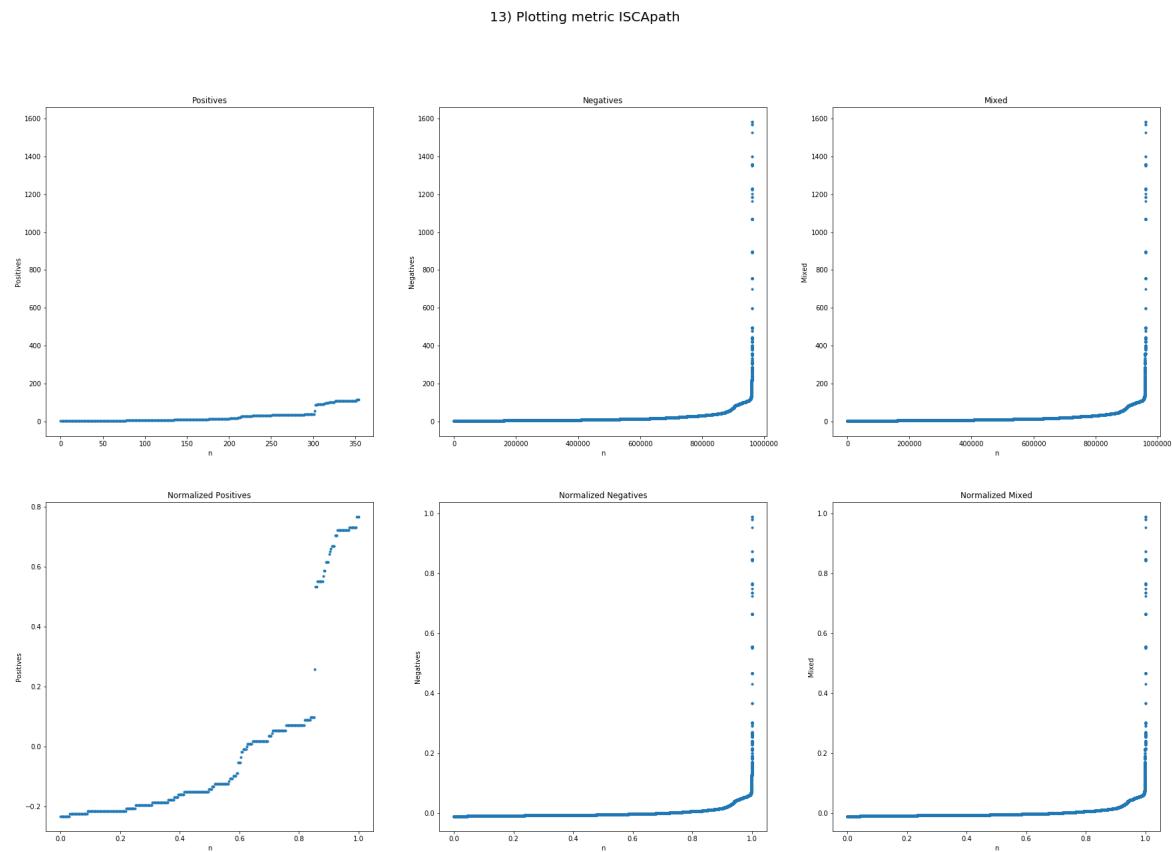


Figura 2.27: Values of metric ISCApath

2.15 commonVar

2.15.1 Metric sample distribution

The data points seem to follow an **Exponential Weibull** distribution.

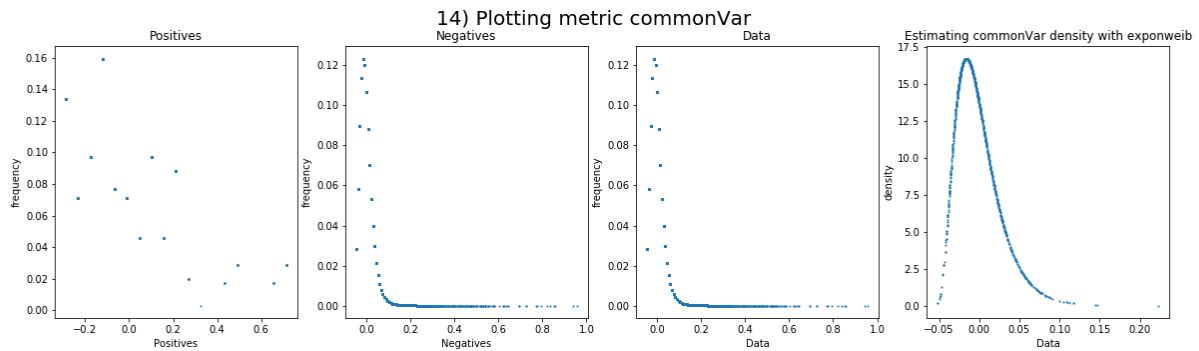


Figura 2.28: Sampling distribution of metric commonVar

2.15.2 Metric values

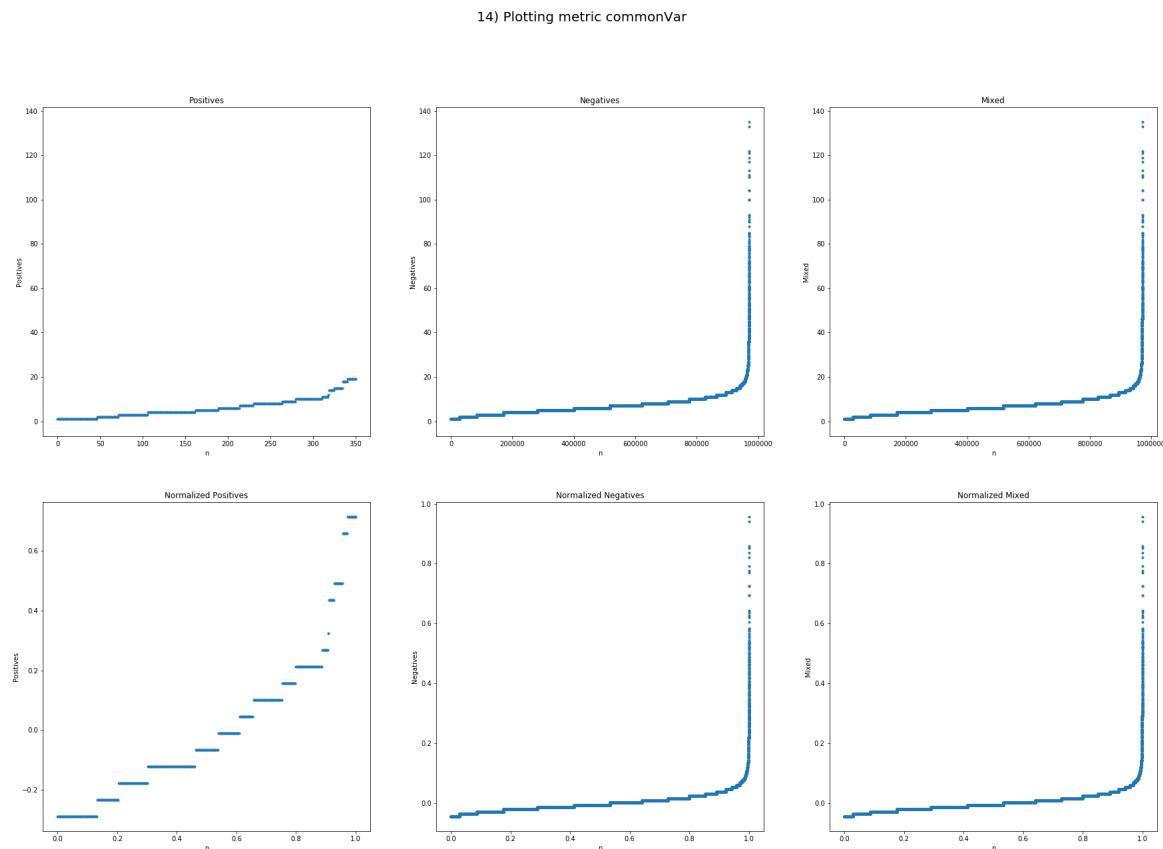


Figura 2.29: Values of metric commonVar

2.16 dbVARCount

2.16.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

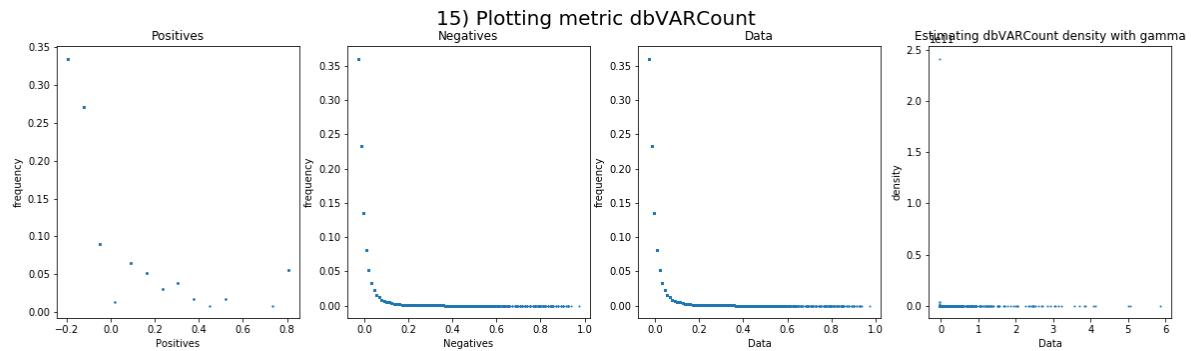


Figura 2.30: Sampling distribution of metric dbVARCount

2.16.2 Metric values

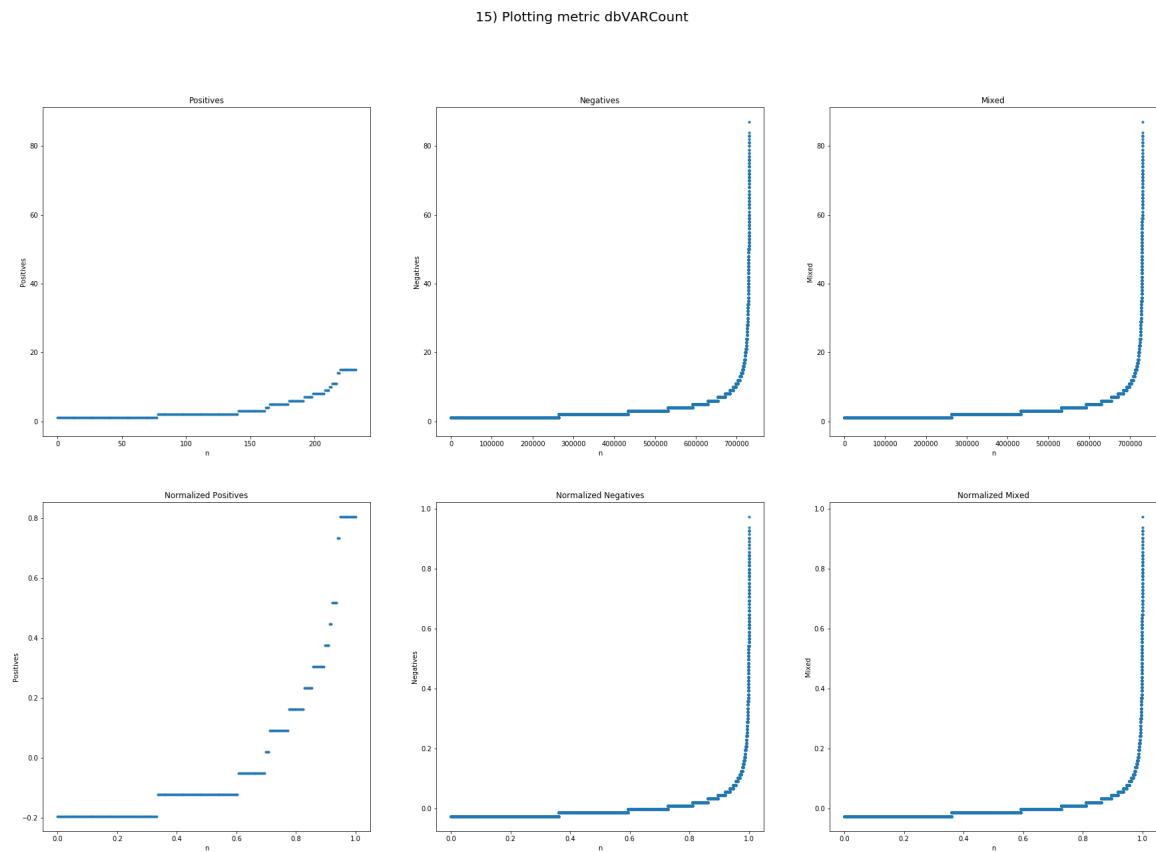


Figura 2.31: Values of metric dbVARCount

2.17 fantom5Perm

2.17.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

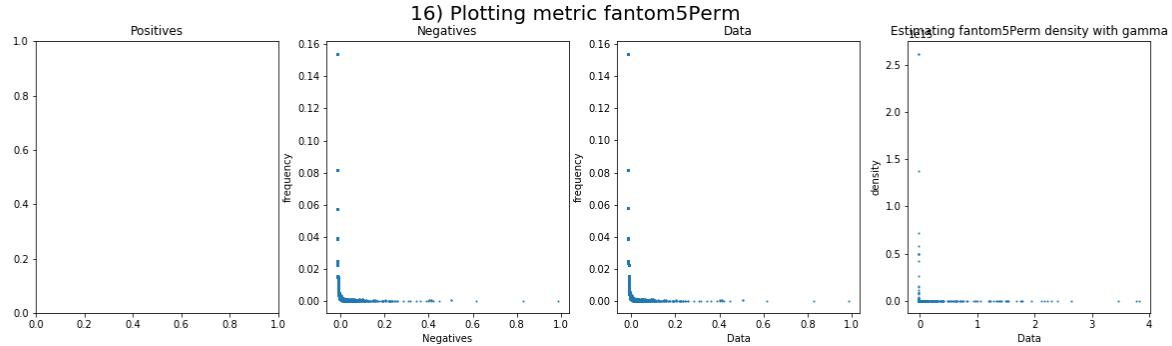


Figura 2.32: Sampling distribution of metric fantom5Perm

2.17.2 Metric values

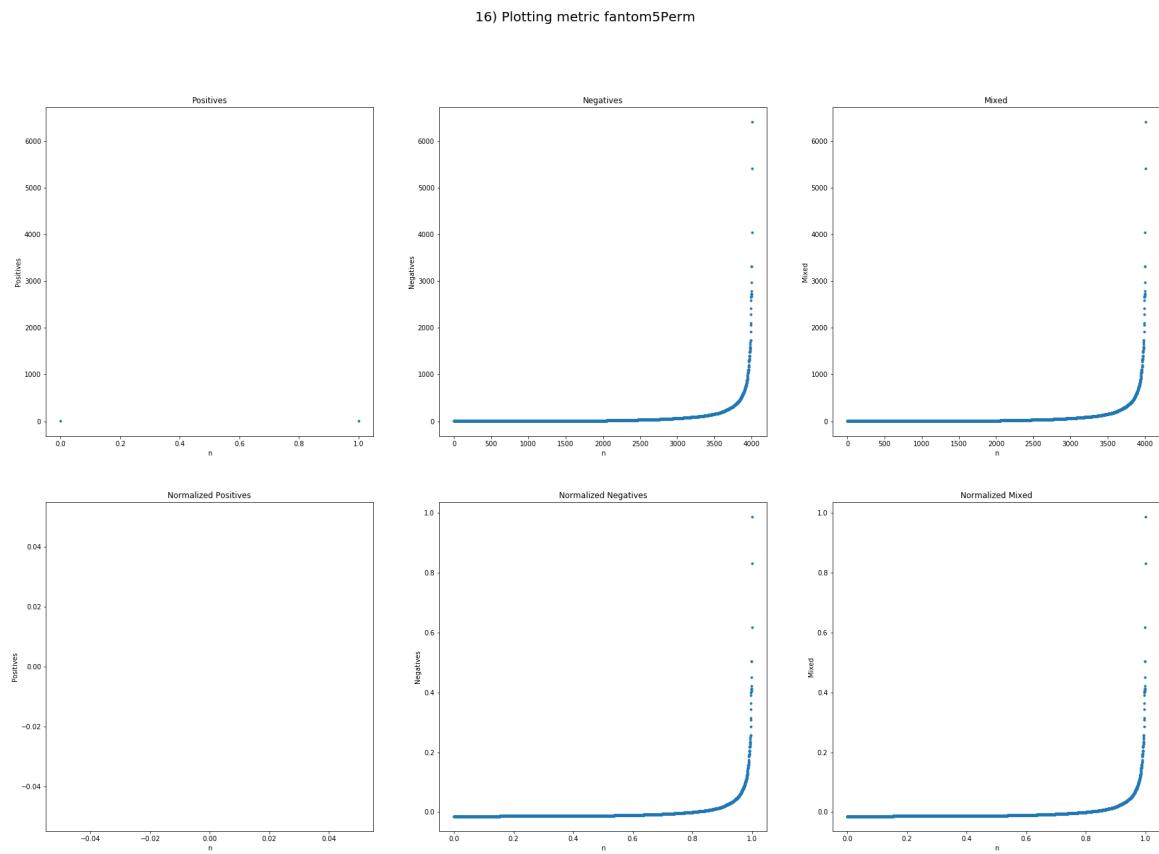


Figura 2.33: Values of metric fantom5Perm

2.18 fantom5Robust

2.18.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

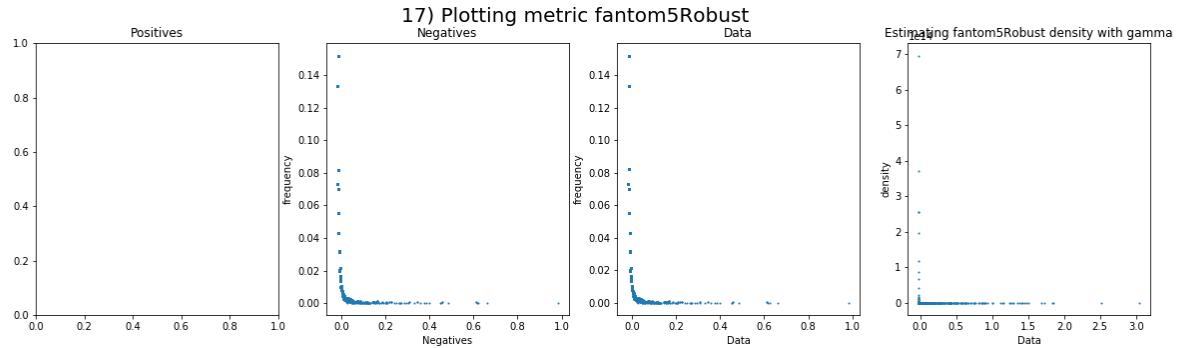


Figura 2.34: Sampling distribution of metric fantom5Robust

2.18.2 Metric values

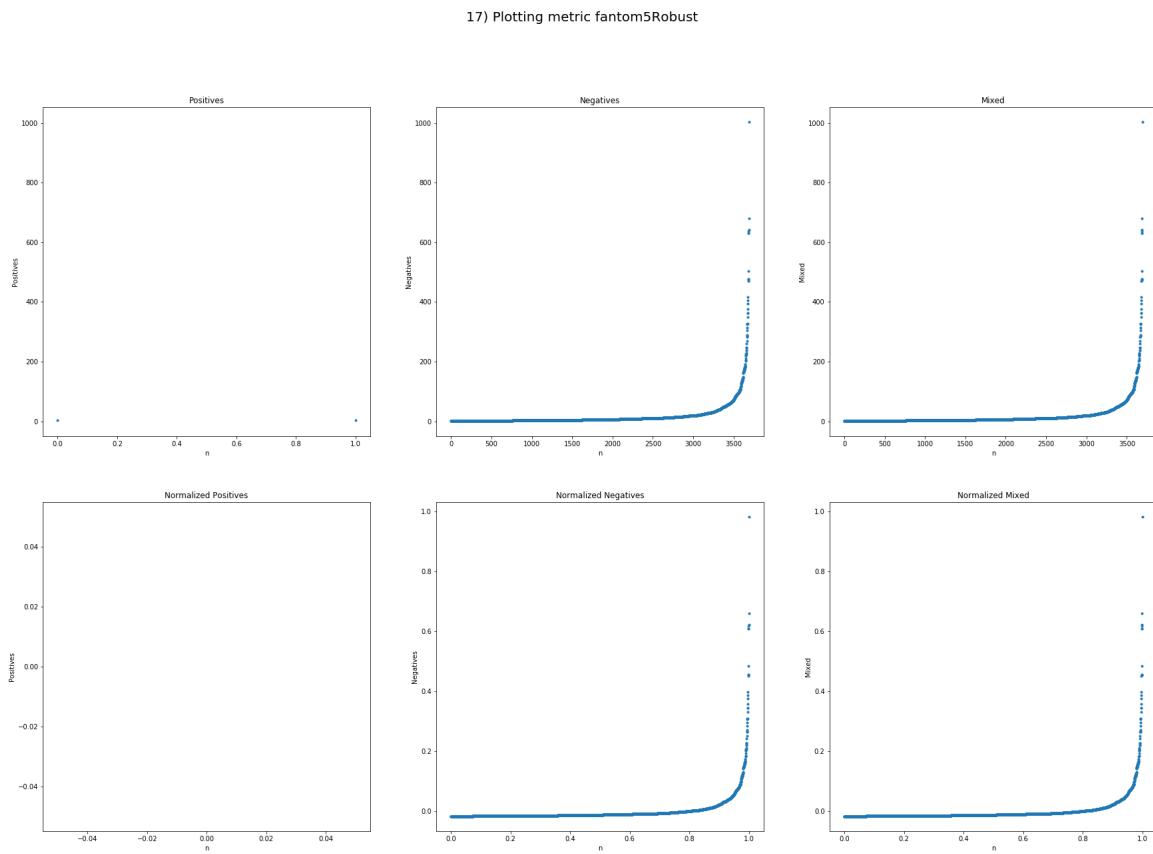


Figura 2.35: Values of metric fantom5Robust

2.19 fracRareCommon

2.19.1 Metric sample distribution

The data points seem to follow an **Beta** distribution.

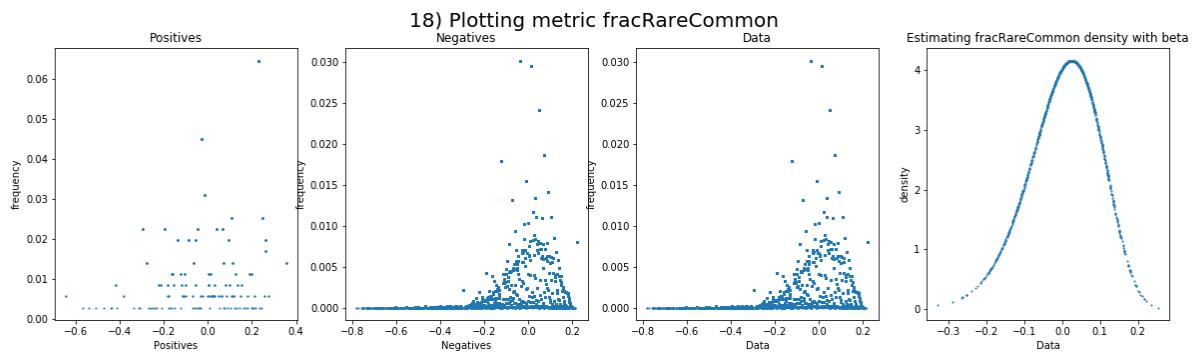


Figura 2.36: Sampling distribution of metric fracRareCommon

2.19.2 Metric values

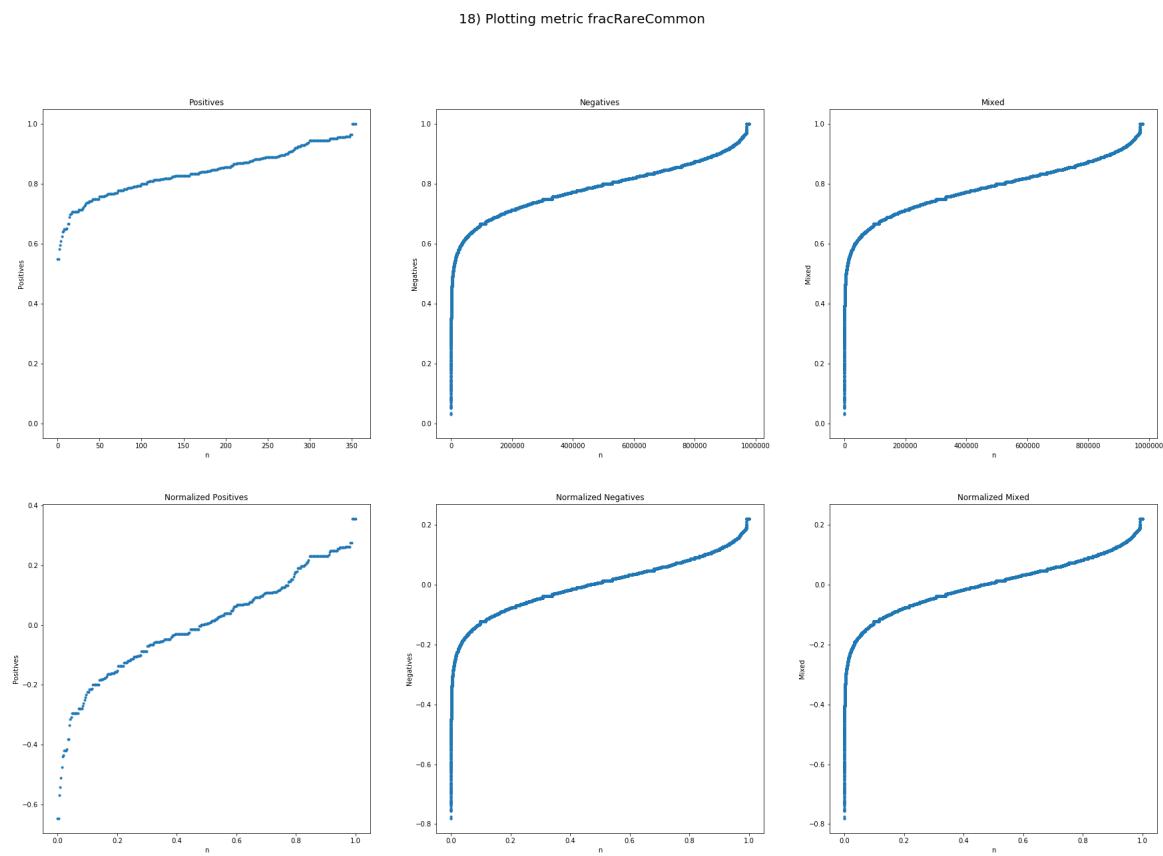


Figura 2.37: Values of metric fracRareCommon

2.20 mamPhastCons46way

2.20.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

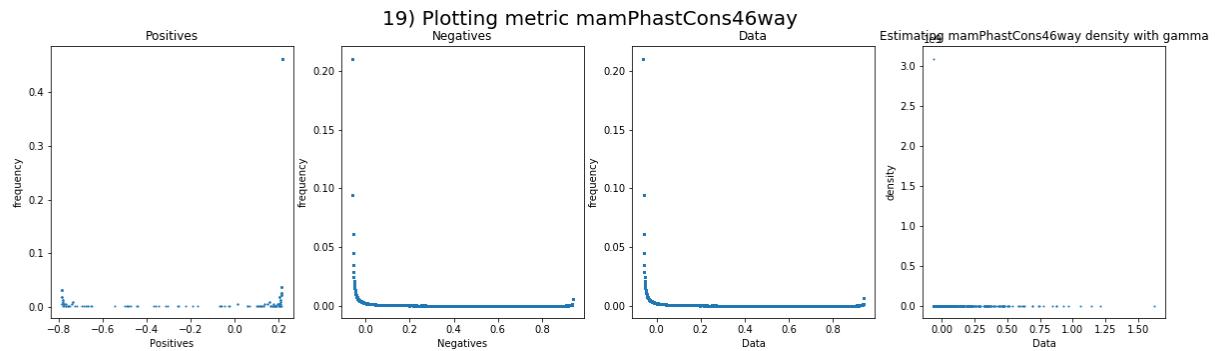


Figura 2.38: Sampling distribution of metric mamPhastCons46way

2.20.2 Metric values

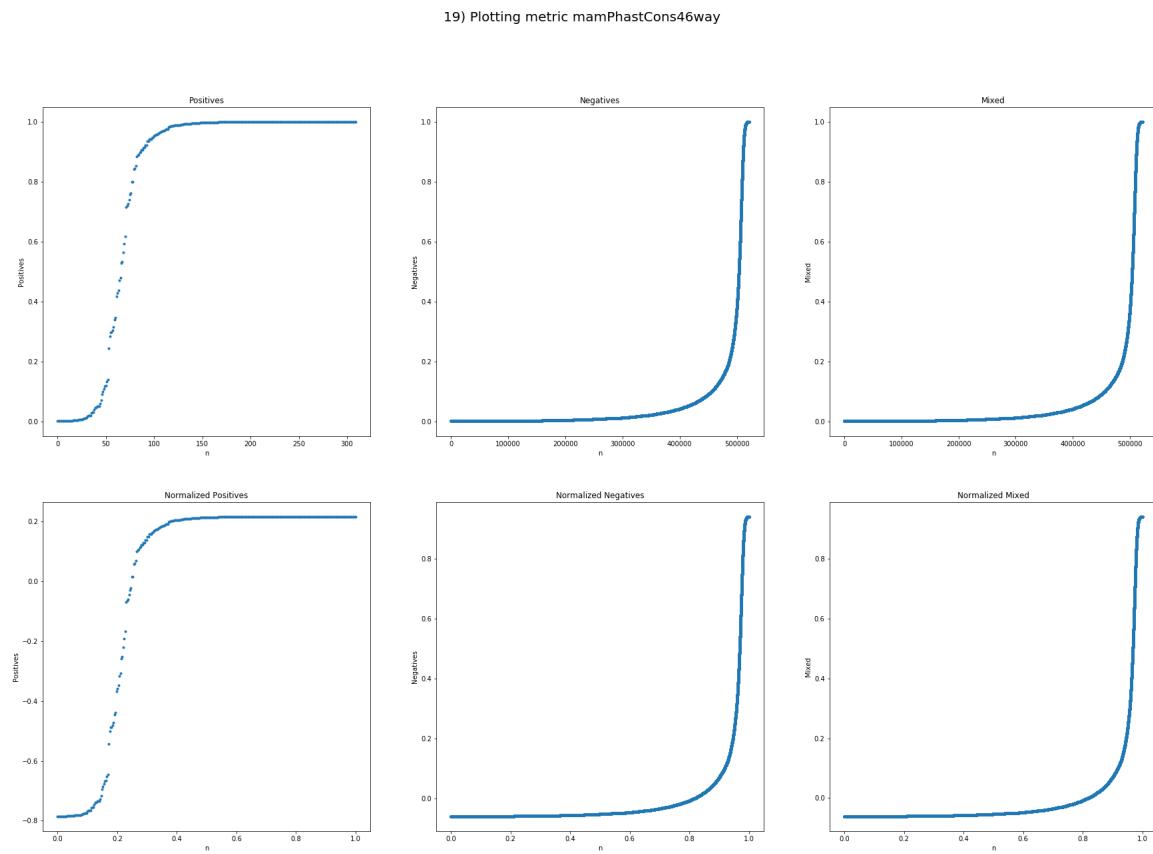


Figura 2.39: Values of metric mamPhastCons46way

2.21 mamPhyloP46way

2.21.1 Metric sample distribution

The data points seem to follow a **Gaussian** distribution.

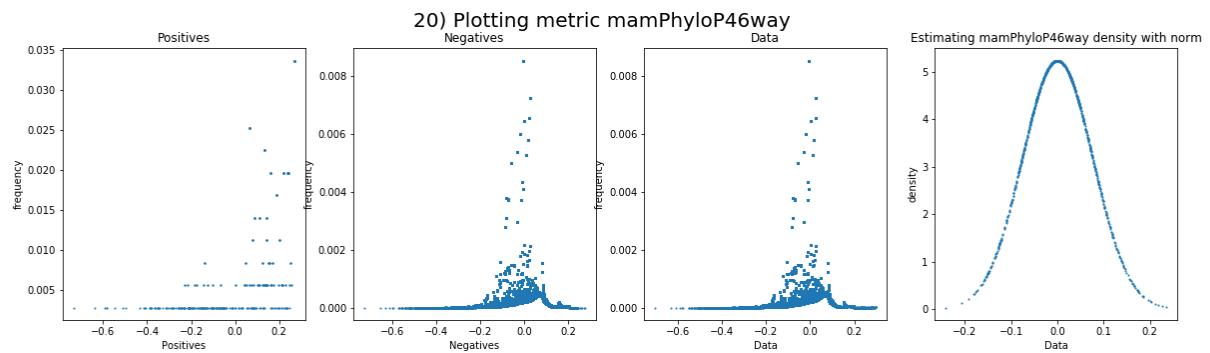


Figura 2.40: Sampling distribution of metric mamPhyloP46way

2.21.2 Metric values

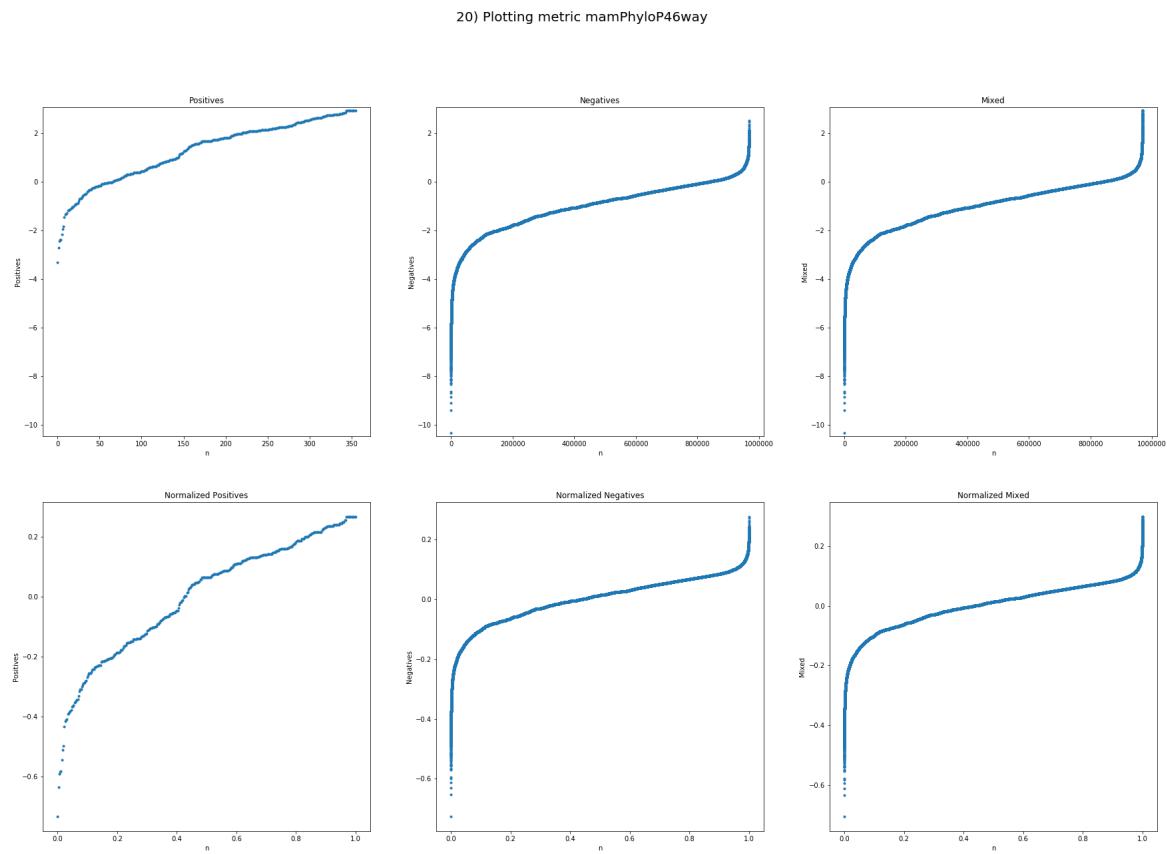


Figura 2.41: Values of metric mamPhyloP46way

2.22 numTFBSConserved

2.22.1 Metric sample distribution

The data points seem to follow a **exponential** distribution.

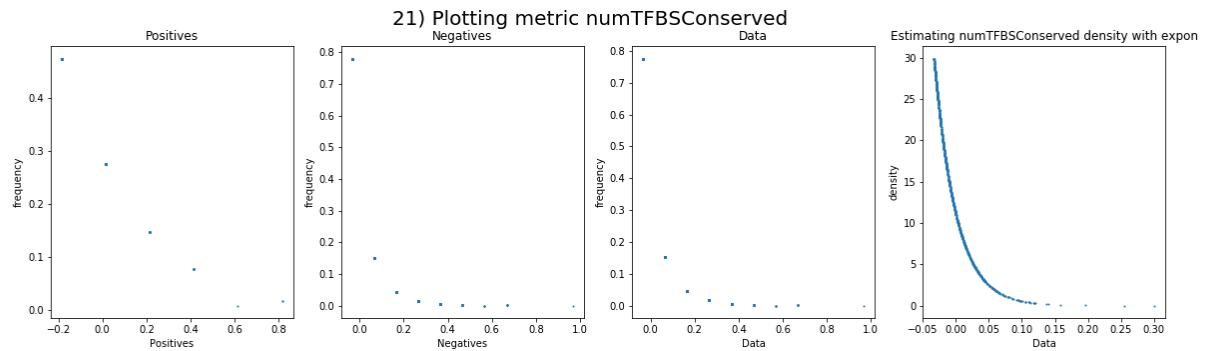


Figura 2.42: Sampling distribution of metric numTFBSConserved

2.22.2 Metric values

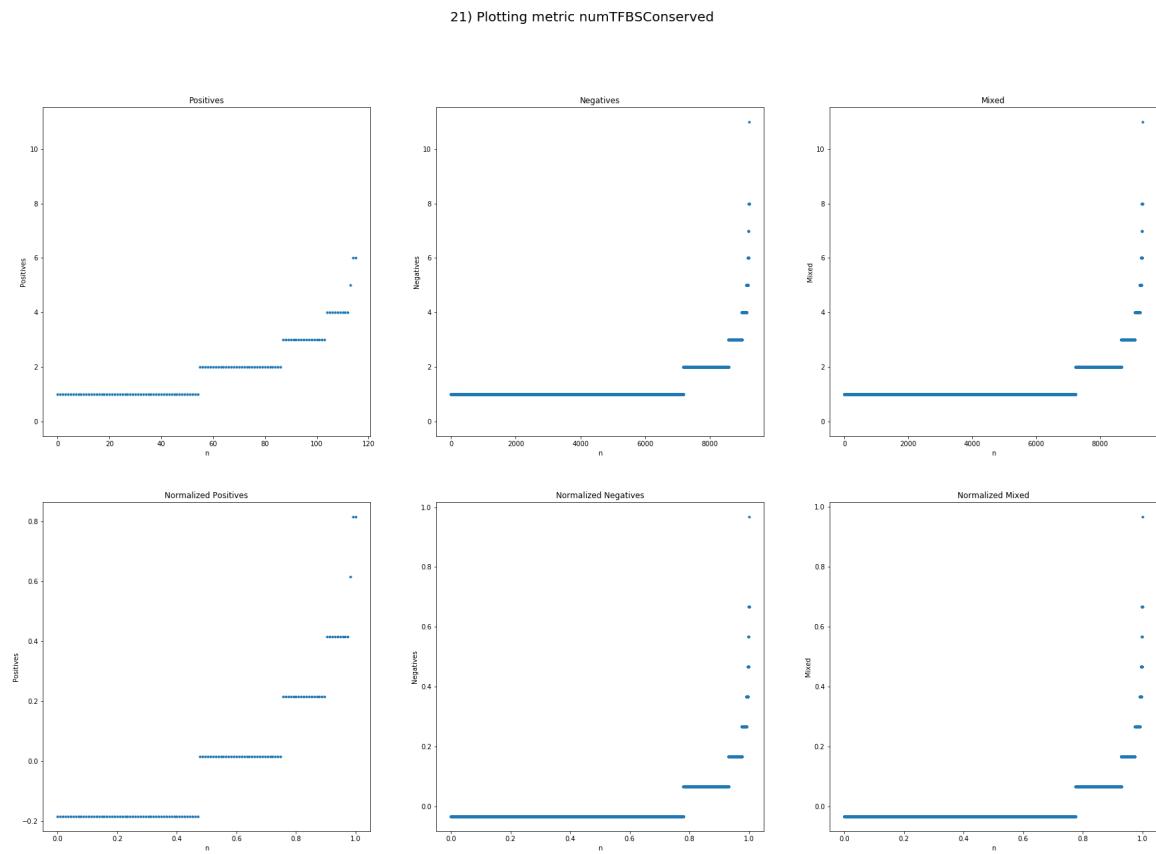


Figura 2.43: Values of metric numTFBSConserved

2.23 priPhastCons46way

2.23.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

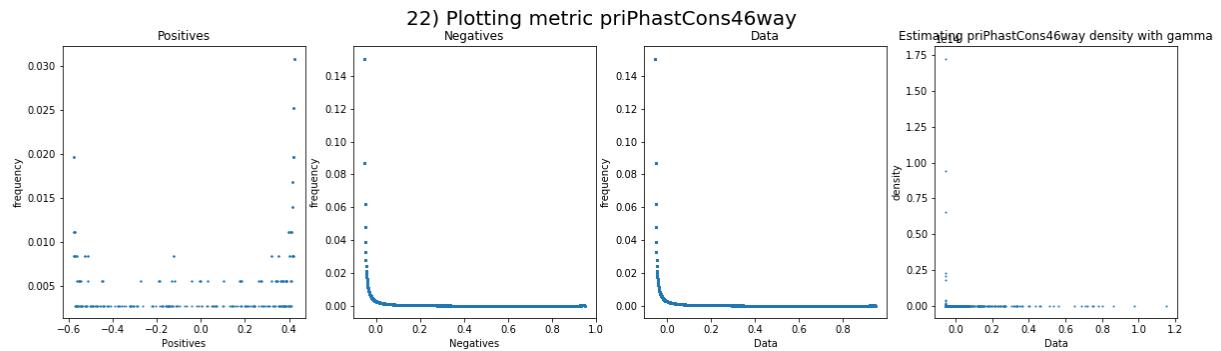


Figura 2.44: Sampling distribution of metric priPhastCons46way

2.23.2 Metric values

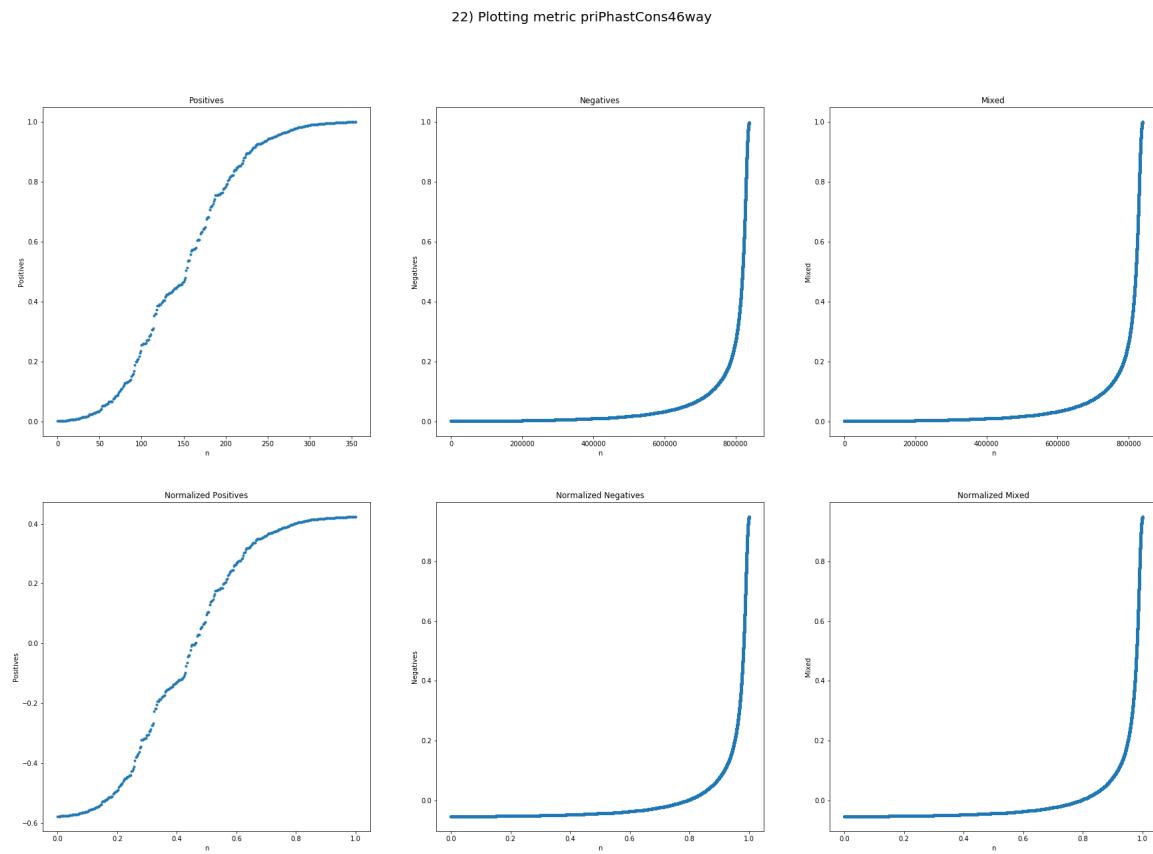


Figura 2.45: Values of metric priPhastCons46way

2.24 priPhyloP46way

2.24.1 Metric sample distribution

The data points seem to follow an **Beta** distribution.

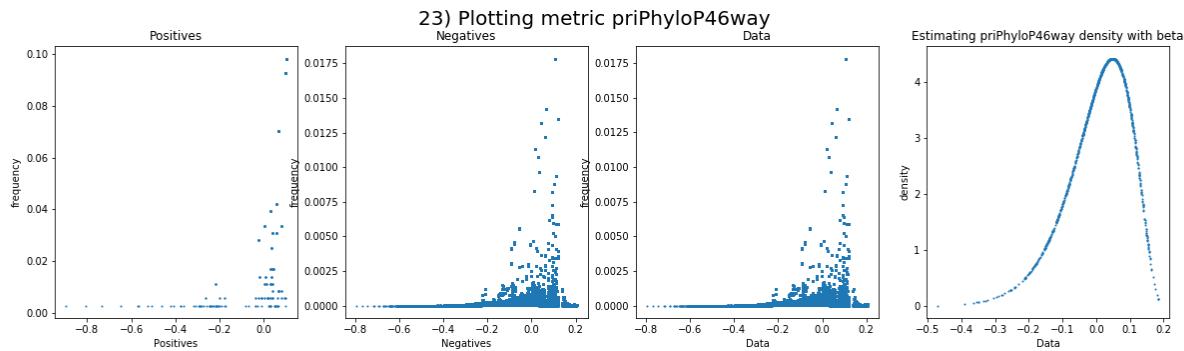


Figura 2.46: Sampling distribution of metric priPhyloP46way

2.24.2 Metric values

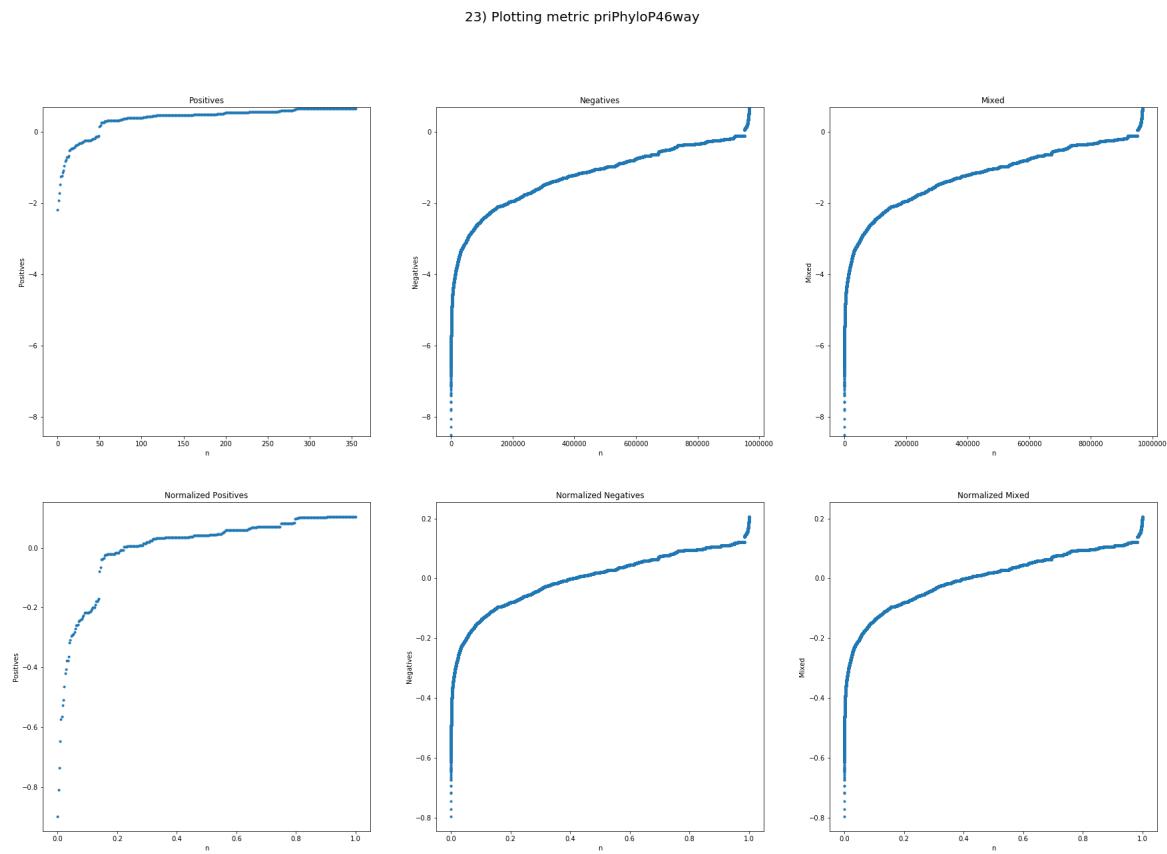


Figura 2.47: Values of metric priPhyloP46way

2.25 rareVar

2.25.1 Metric sample distribution

The data points seem to follow an **Beta** distribution.

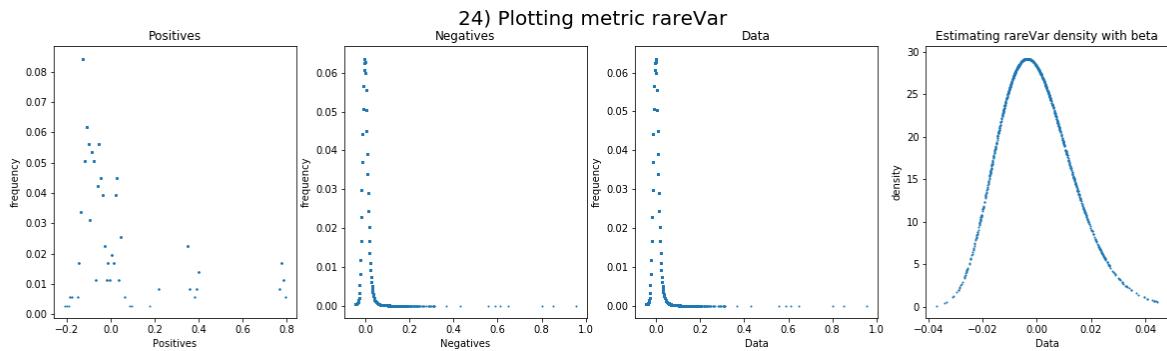


Figura 2.48: Sampling distribution of metric rareVar

2.25.2 Metric values

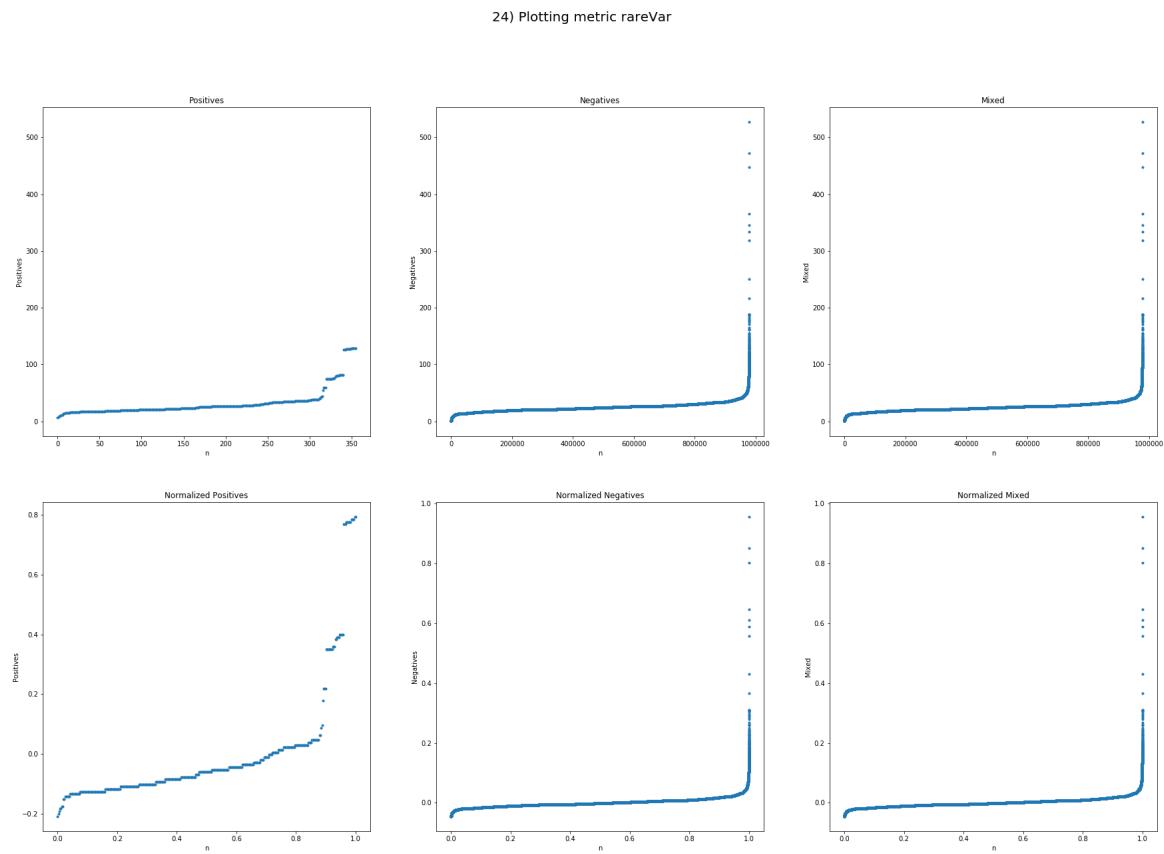


Figura 2.49: Values of metric rareVar

2.26 verPhastCons46way

2.26.1 Metric sample distribution

The data points seem to follow a **Gamma** distribution.

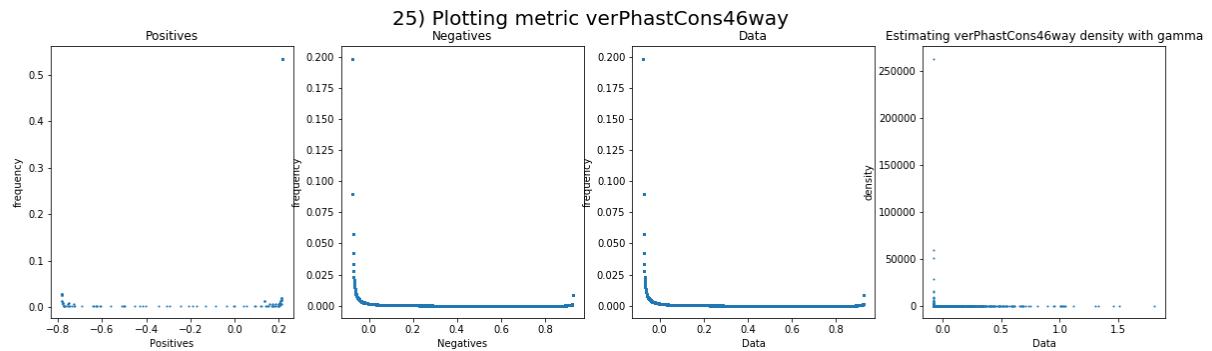


Figura 2.50: Sampling distribution of metric verPhastCons46way

2.26.2 Metric values

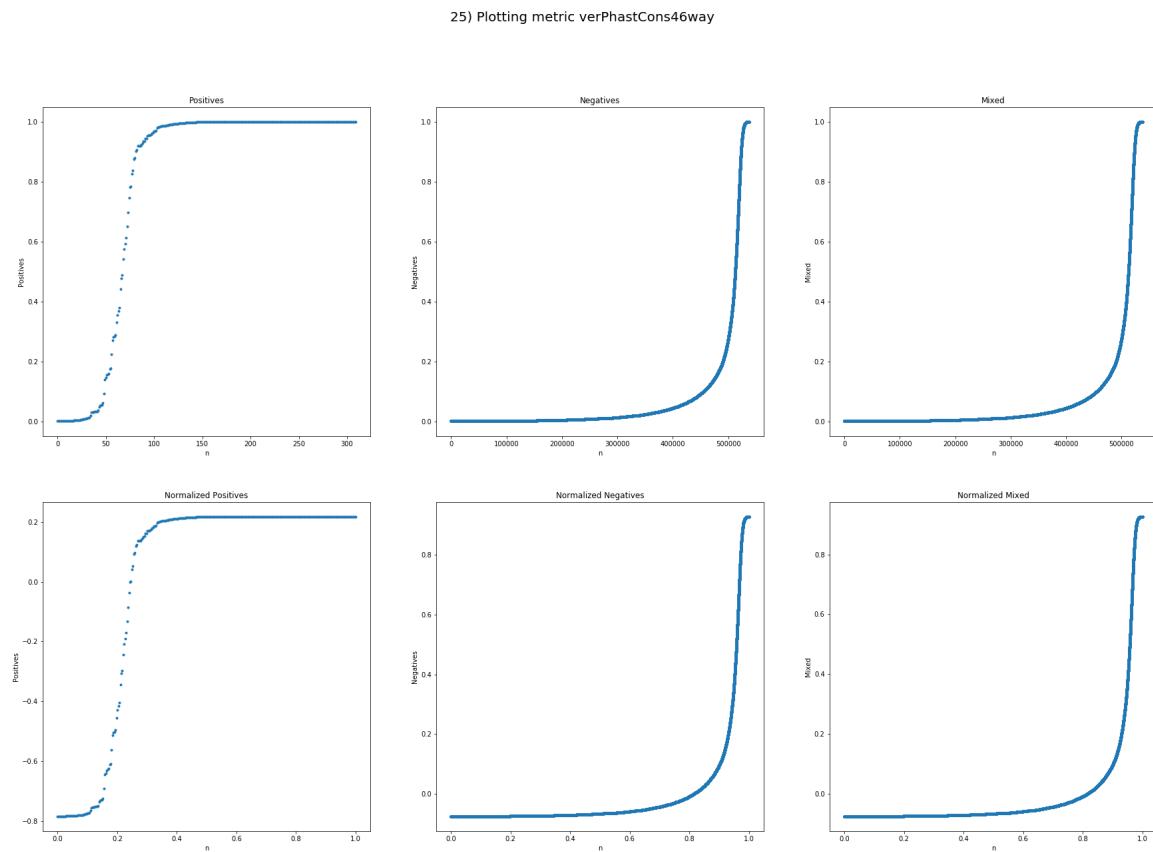


Figura 2.51: Values of metric verPhastCons46way

2.27 verPhyloP46way

2.27.1 Metric sample distribution

The data points seem to follow a **Gaussian** distribution.

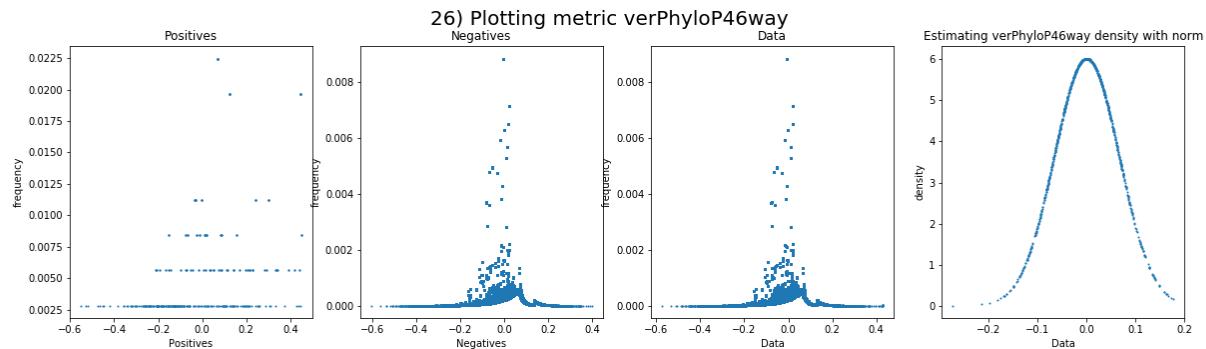


Figura 2.52: Sampling distribution of metric verPhyloP46way

2.27.2 Metric values

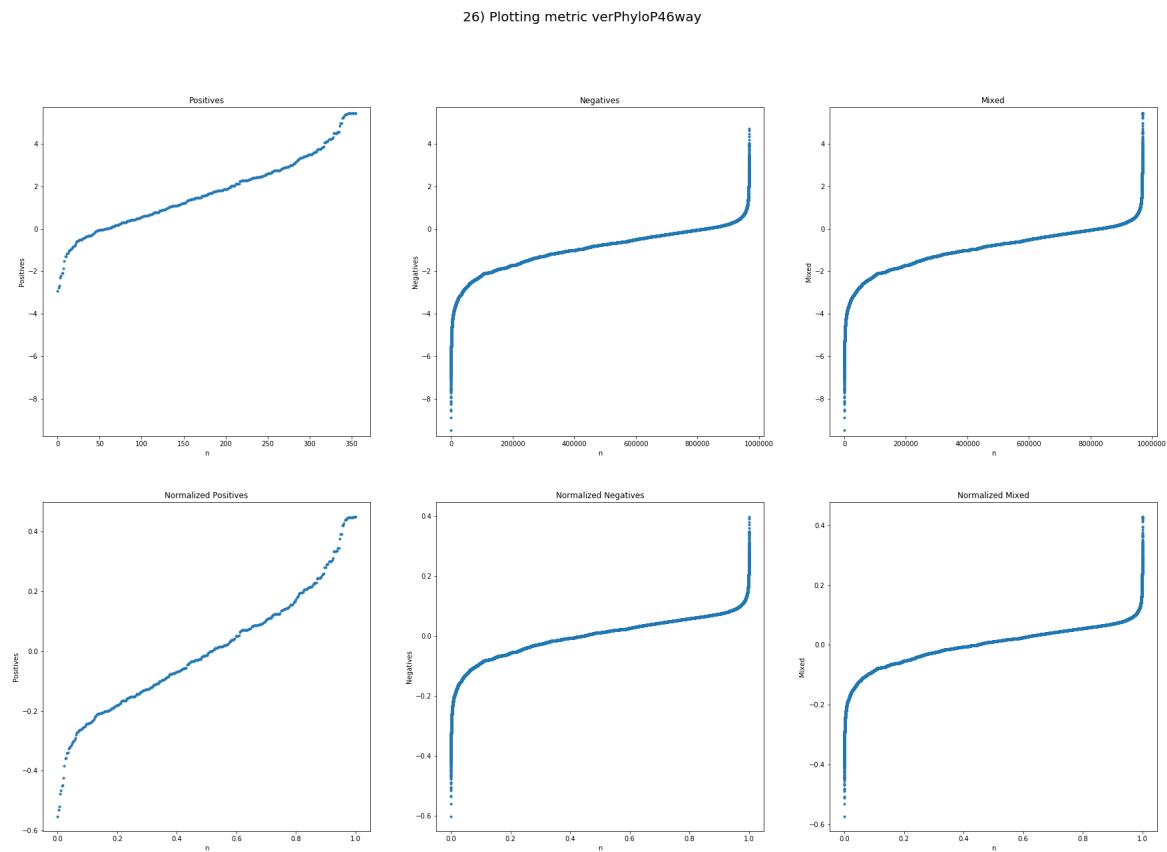


Figura 2.53: Values of metric verPhyloP46way

3

Metric distribution summary

The metrics seem to follow these sample distributions:

Metric	Distribution
CpGobsExp	Beta
CpGperCpG	Beta
CpGperGC	Gaussian
DGVCount	Gamma
DnaseClusteredHyp	Gamma
EncH3K27Ac	Gamma
GCContent	Gaussian
EncH3K4Me3	Gamma
ISCApath	Gamma
DnaseClusteredScore	Beta
EncH3K4Me1	Gamma
GerpRS	Gamma
GerpRSpv	Gamma
commonVar	Exponential Weibull
dbVARCount	Gamma
fantom5Perm	Gamma
fantom5Robust	Gamma
mamPhastCons46way	Gamma
priPhastCons46way	Gamma
rareVar	Beta
verPhastCons46way	Gamma
numTFBSConserved	Exponential
fracRareCommon	Beta
priPhyloP46way	Beta
verPhyloP46way	Gaussian
mamPhyloP46way	Gaussian

Tabella 3.1: Metrics and their distribution

4

Data correlation

We now proceed to try and identify eventual data correlations.

4.1 Scatter plot

A scatter plot with higher resolution is available in the project repository.

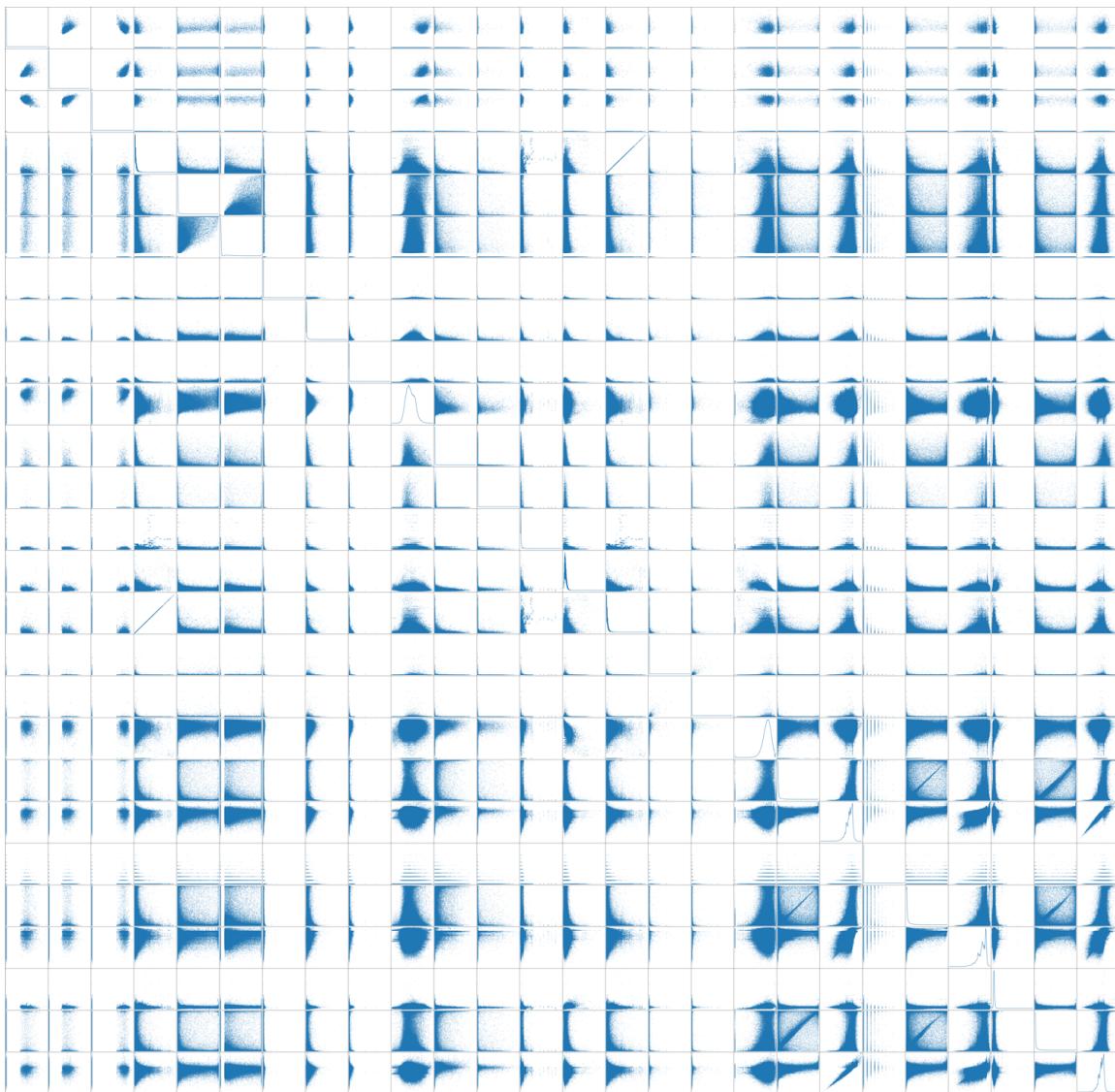


Figura 4.1: Scatter plot

4.2 Correlation coefficient matrix

A correlation matrix with higher resolution is available in the project repository.

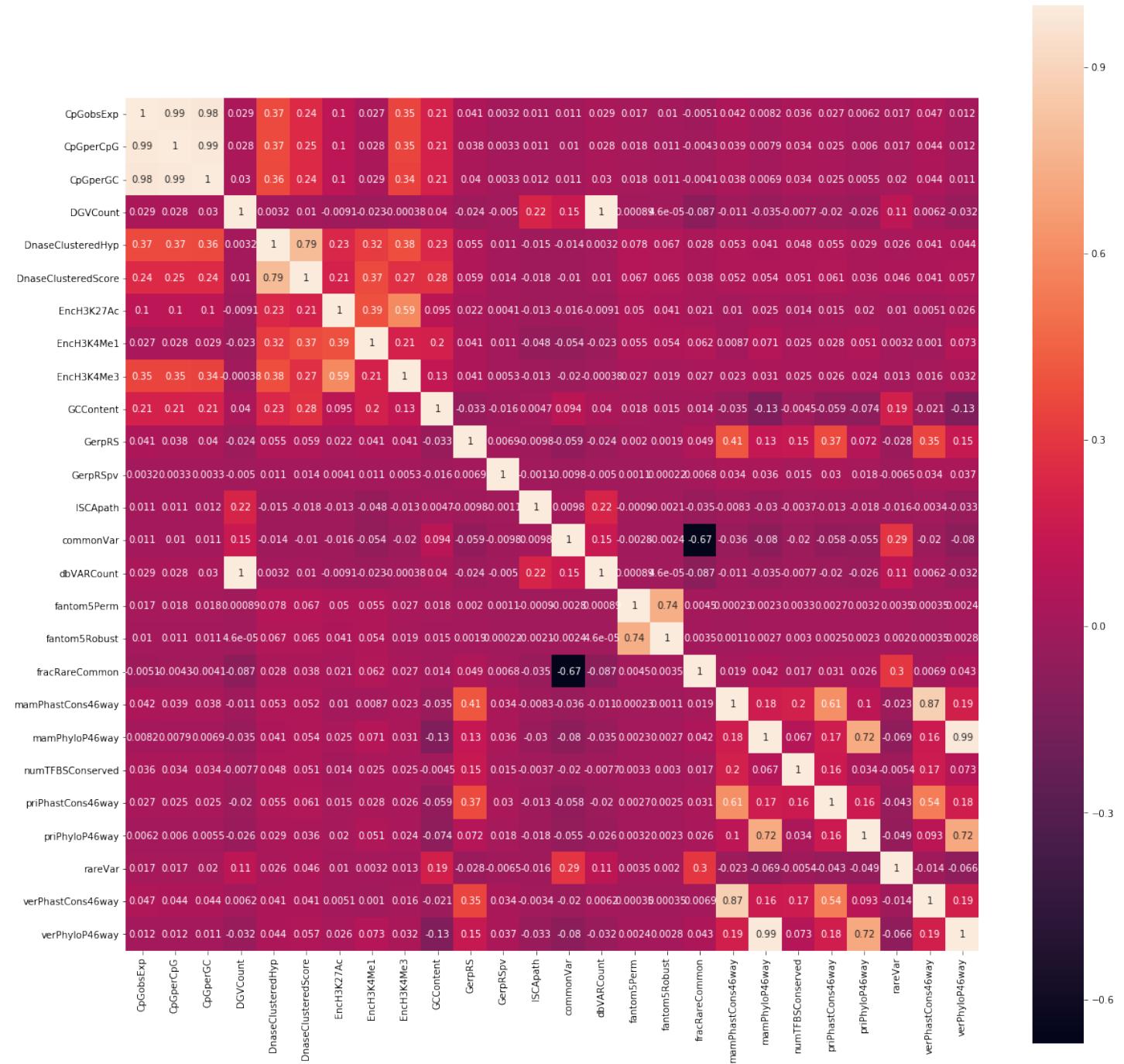


Figura 4.2: Correlation matrix

4.2.1 CpGobsExp and CpGperCpG

The two metric CpGobsExp and CpGperCpG have correlation index 0.9856203442596099.

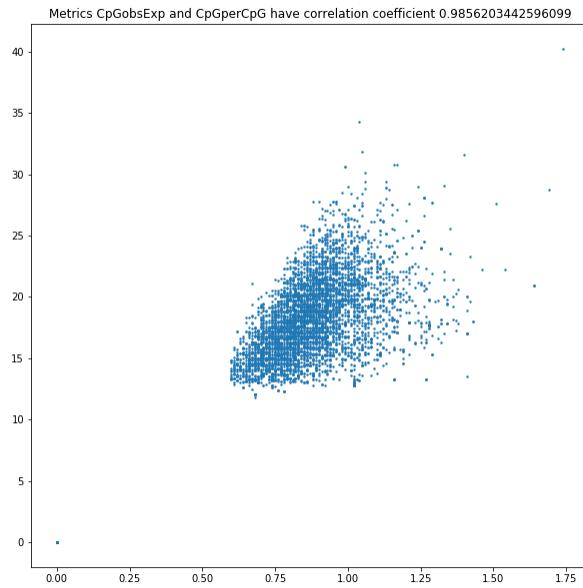


Figura 4.3: CpGobsExp and CpGperCpG

These two metrics have an extremely high correlation, so one of the two will be removed from the dataset, arbitrarily **CpGobsExp**.

4.2.2 CpGobsExp and CpGperGC

The two metric CpGobsExp and CpGperGC have correlation index 0.9785748818167331.

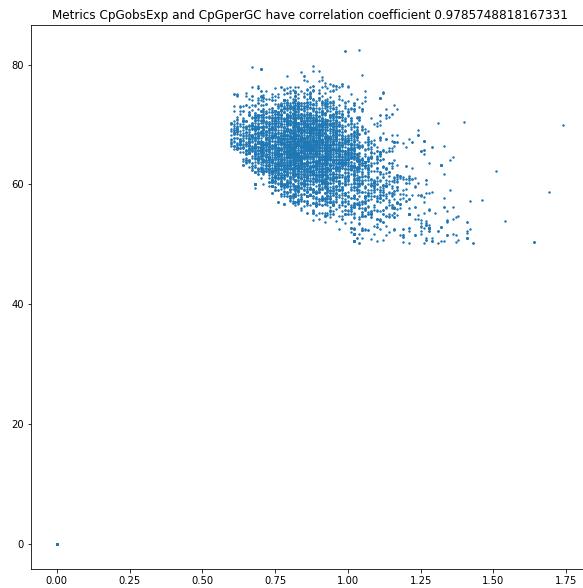


Figura 4.4: CpGobsExp and CpGperGC

These two metrics have an extremely high correlation, so one of the two will be removed from the dataset, arbitrarily **CpGobsExp**.

4.2.3 CpGperCpG and CpGperGC

The two metric CpGperCpG and CpGperGC have correlation index 0.9897887253514923.

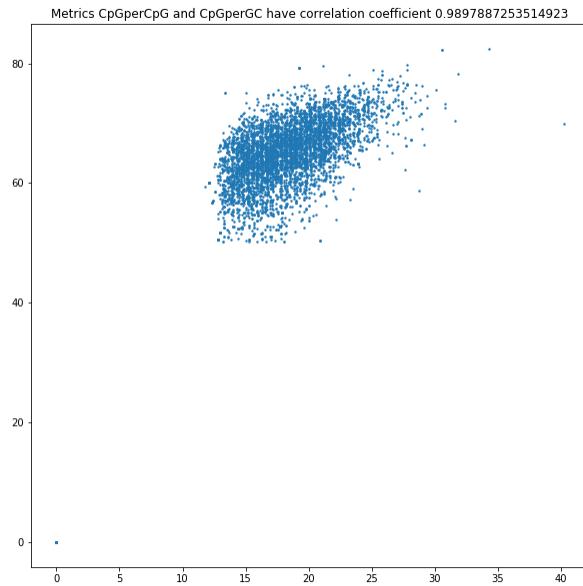


Figura 4.5: CpGperCpG and CpGperGC

These two metrics have an extremely high correlation, so one of the two will be removed from the dataset, arbitrarily **CpGperCpG**.

4.2.4 dbVARCount and DGVCount

The two metric dbVARCount and DGVCount have correlation index 1.0.

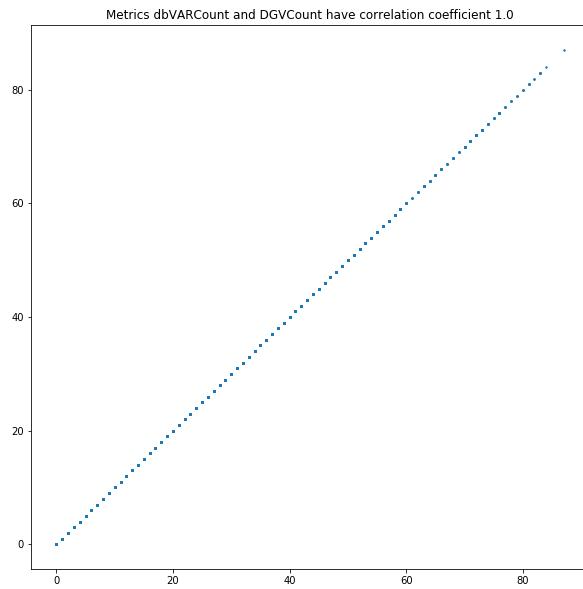


Figura 4.6: dbVARCount and DGVCount

These two metrics have an extremely high correlation, so one of the two will be removed from the dataset, arbitrarily **dbVARCount**.

4.2.5 mamPhyloP46way and verPhyloP46way

The two metric mamPhyloP46way and verPhyloP46way have correlation index 0.9902257463490804.

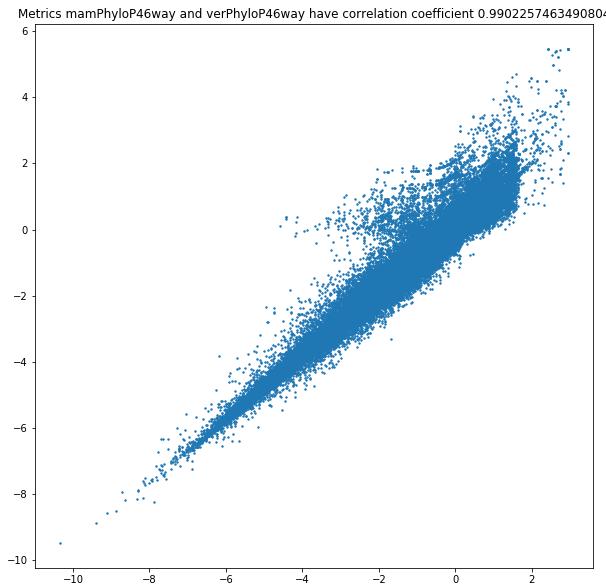


Figura 4.7: mamPhyloP46way and verPhyloP46way

These two metrics have an extremely high correlation, so one of the two will be removed from the dataset, arbitrarily **mamPhyloP46way**.

4.2.6 DnaseClusteredHyp and DnaseClusteredScore

The two metric DnaseClusteredHyp and DnaseClusteredScore have correlation index 0.7863337778663062.

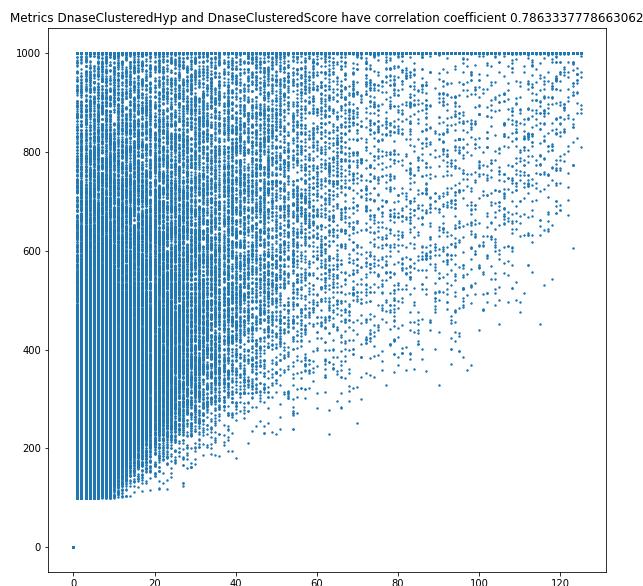


Figura 4.8: DnaseClusteredHyp and DnaseClusteredScore

The correlation value is high, but not enough to motivate actions such as the removal from the dataset.

4.2.7 mamPhastCons46way and verPhastCons46way

The two metric mamPhastCons46way and verPhastCons46way have correlation index 0.8700263713585867.

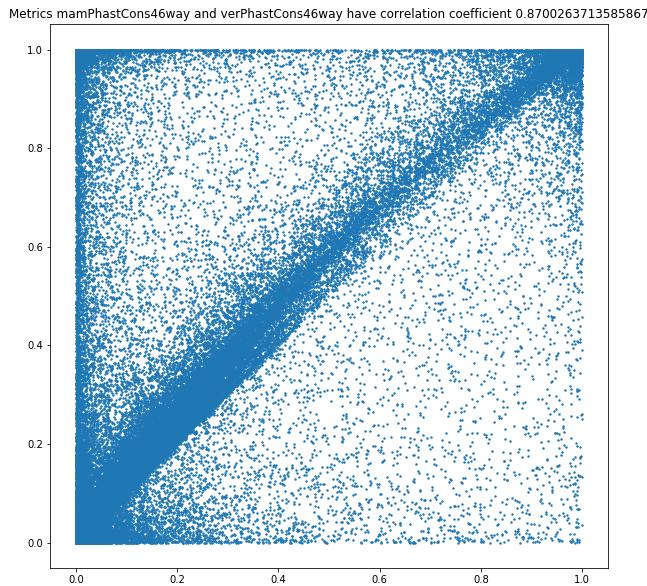


Figura 4.9: mamPhastCons46way and verPhastCons46way

The correlation value is high, but not enough to motivate actions such as the removal from the dataset.

4.3 Identified data correlations

Data correlations exist between:

First metric	Second metric	Correlation coefficient
CpGobsExp	CpGperCpG	0.9856203442596099
CpGobsExp	CpGperGC	0.9785748818167331
CpGperCpG	CpGperGC	0.9897887253514923
dbVARCount	DGVCount	1.0
mamPhyloP46way	verPhyloP46way	0.9902257463490804
DnaseClusteredHyp	DnaseClusteredScore	0.7863337778663062
mamPhastCons46way	verPhastCons46way	0.8700263713585867

4.4 Correlation table after removing highly correlated data

After having removed from the dataset metrics **CpGobsExp**, **CpGperCpG**, **dbVARCount** and **mamPhyloP46way** looks like:

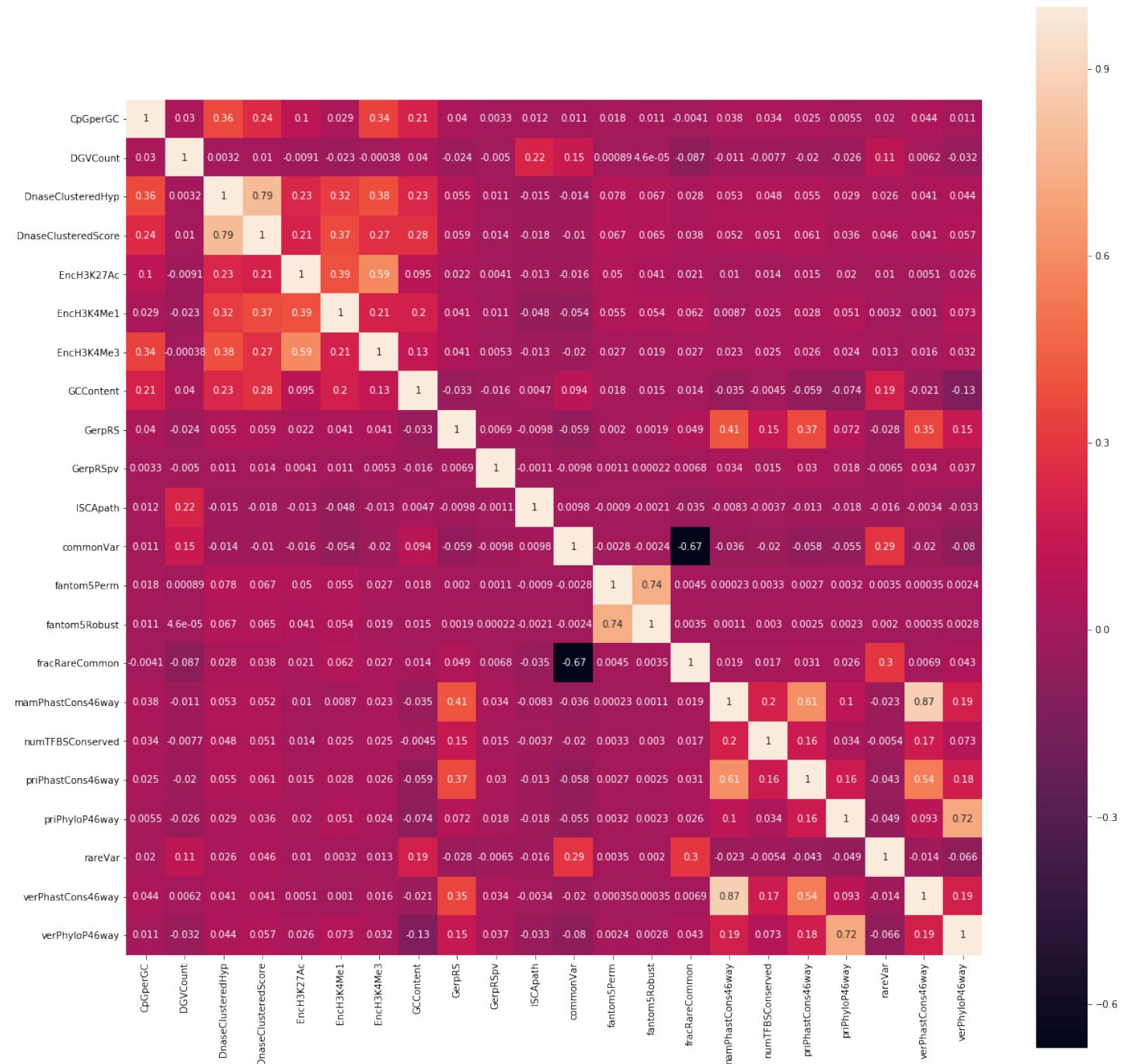


Figura 4.10: Correlation matrix

This correlation matrix with higher resolution is available in the project repository.

5

Dataset visualization

After having removed the correlated metrics we can proceed to use techniques of dimensionality reduction for visualization to see if the dataset is valid for a clustering or machine learning approach:

5.1 PCA

5.1.1 Training dataset visualization

The positive data are clustered inside the negative data: a simple clustering approach would, most probably, not be enough for separating the two classes, but easily separable for any multi-parameters ML approach like networks.

The points isolated on the right are most probably errors in the realization of the dataset.

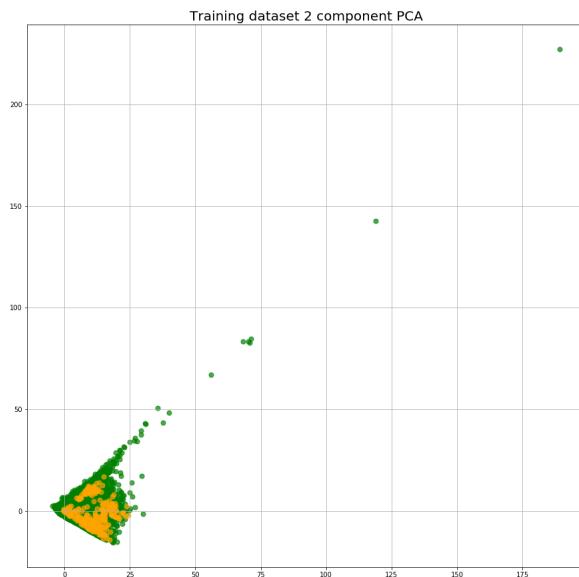


Figura 5.1: Training dataset visualization

5.1.2 Testing dataset visualization

The testing dataset is probably malformed: all the positives point are visibly clustered. A simple clustering approach such as K-Means could probably separate most of them successfully.

An ML approach with a network will probably reach sooner an high accuracy on the testing dataset than on the training dataset for this reason, being able to classify most of the positive points immediately.



Figura 5.2: Testing dataset visualization

5.1.3 Mixed dataset visualization

The two datasets overlap correctly.

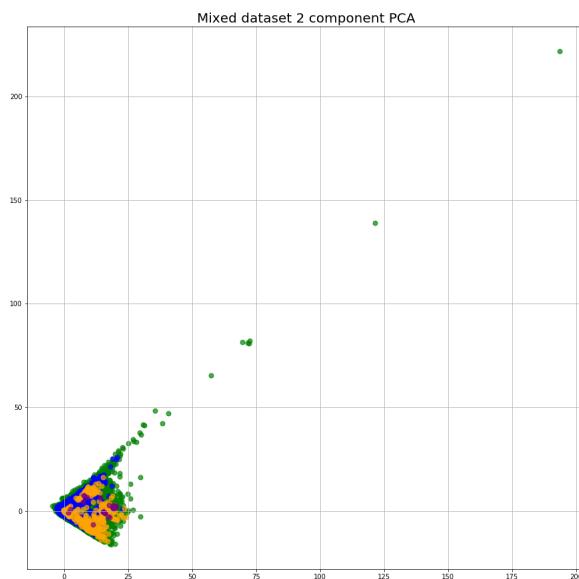


Figura 5.3: Mixed dataset visualization

6

Dataset issues

6.1 Possible dataset errors

The PCA visualization shows points extremely out of the dataset cluster. It is possible that these points are errors.

6.2 Biased testing dataset

The testing dataset does not seem to reflect the training dataset distribution, but agglomerates the two classes in two separated clusters. It is therefore unhelpful when calculating the classification success in networks, as the training dataset result harder to classify than the testing one.

Parte II

Network implementation

7

Model architecture

7.1 Input

$$m' = \frac{\text{metric} - \mathbb{E}(\text{metric values})}{\max\{\text{metric values}\} - \min\{\text{metric values}\}}$$

Figura 7.1: Input normalization

7.2 Output

The output layer of the neural network is modelled by one neuron with a **sigmoid** as activation function. When active models the positive class, and when inactive models the negative class.

$$\text{sigmoid}(x) = \frac{e^x}{e^x + 1}$$

Figura 7.2: Sigmoid

7.3 Weight distribution based on input distribution

Since input values are not from any particular distribution or hold properties such as $\mathbb{E}(X) = 0$ or $\text{Var}(X) = 1$ (in some metrics mean and variance are far from these values) they do not suggest to use any specific distribution.

7.4 Weight distribution based on activation functions and regularization layers

The codomain values from the activation functions, being SELU for most of the network, tend to hold the properties of $\mathbb{E}(X) = 0$ and $\text{Var}(X) = 1$ (<http://arxiv.org/abs/1706.02515>). These values are then regularized to penalize extreme weights that may appear when variance starts with a value significantly away from 1.

For these properties weight will be initialized by extracting them from a Gaussian with $\mathbb{E}(X) = 0$ and $\text{Var}(X) = 1$.

7.5 Batch Size

7.6 Hidden layers

7.7 Possibility: Locally connected dense layers

The first two layers could be locally connected dense layers, to exploit the positional information of the input values.

Other than the group of triples, input will be sorted by distribution kind so that the initial interpolations happen mostly with data from the same distribution family.

Sadly, it does not seem that with the current implementation of Keras this is possible for vector inputs, so Dense layers will be used.

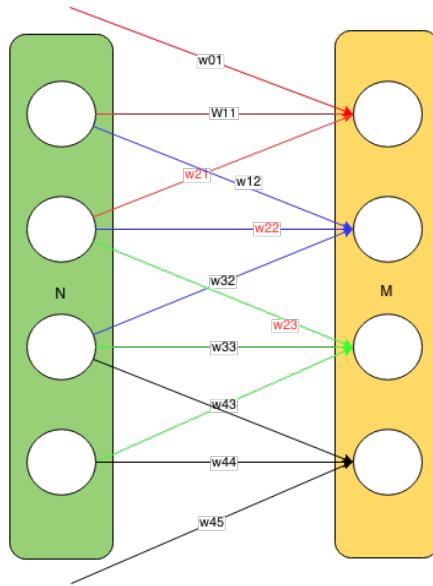


Figura 7.3: Locally connected layer

For the following hidden layers we will be using dense connected layers, with a pyramidal structure (reducing the number of the neurons from 75 to 1).

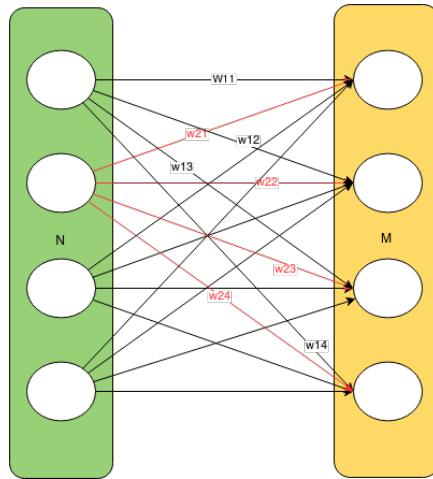


Figura 7.4: Dense connected layer

7.8 Activation function

We'll be using **SELU** for the hidden layers:

$$\text{selu}(x) = \lambda \begin{cases} x & x > 0 \\ \alpha e^x - \alpha & x \leq 0 \end{cases}$$

Figura 7.5: SELU

7.9 Regularization

Regularization layers will be alternated to the dense layers to penalize weight extreme growth.

7.10 Drop out

In addition to regularization, also **drop out** of 10% of neurons per hidden layer will be applied.

7.11 Loss function

Since the task assigned to the network is a binary classification the loss function will be the **cross entropy**.

7.12 Update policy

As update policy we are going to use a form of gradient **back propagation** with **adam**.

7.13 Network model representation

Graphical model of the network. A version in higher resolution is available in the repository.

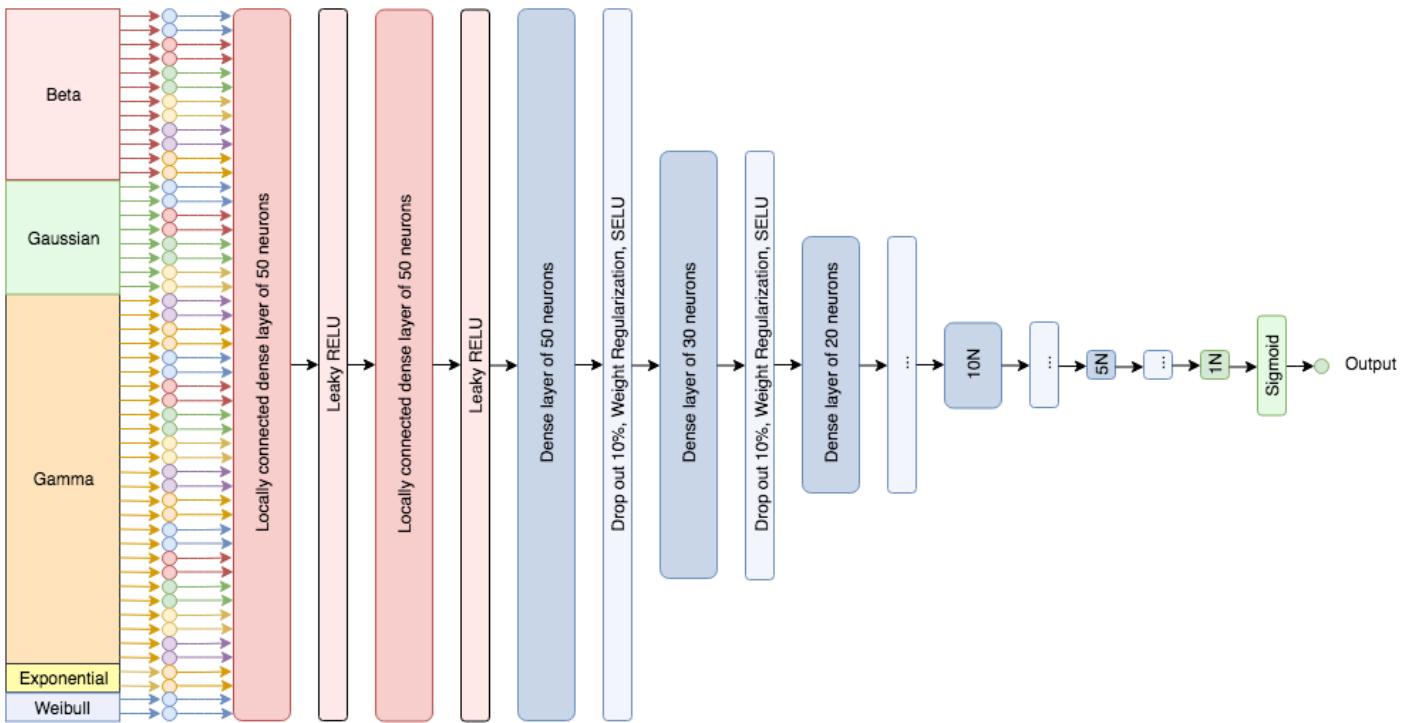


Figura 7.6: Model of the network

8

References

LatexTools does not compile references at this time.