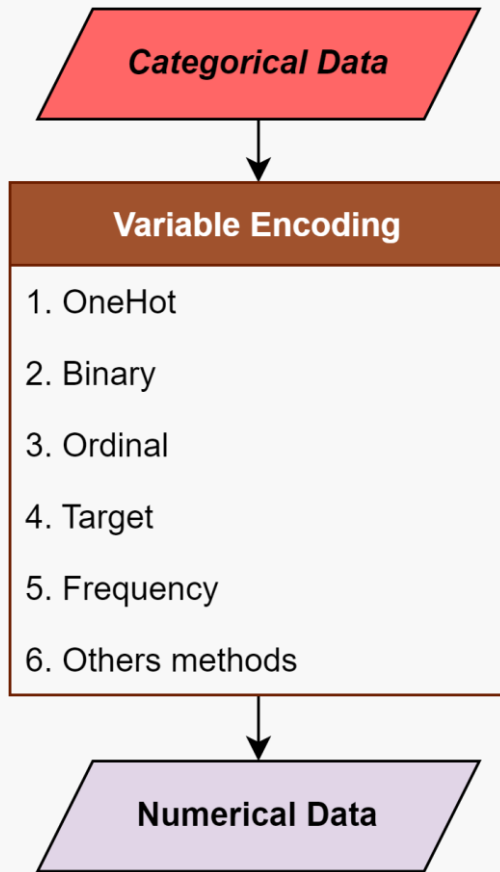


# Supervised & Unsupervised Encoding Schemes of Binary Variables for Prediction Performance Enhancement

Institute of Industrial Engineering, NTU

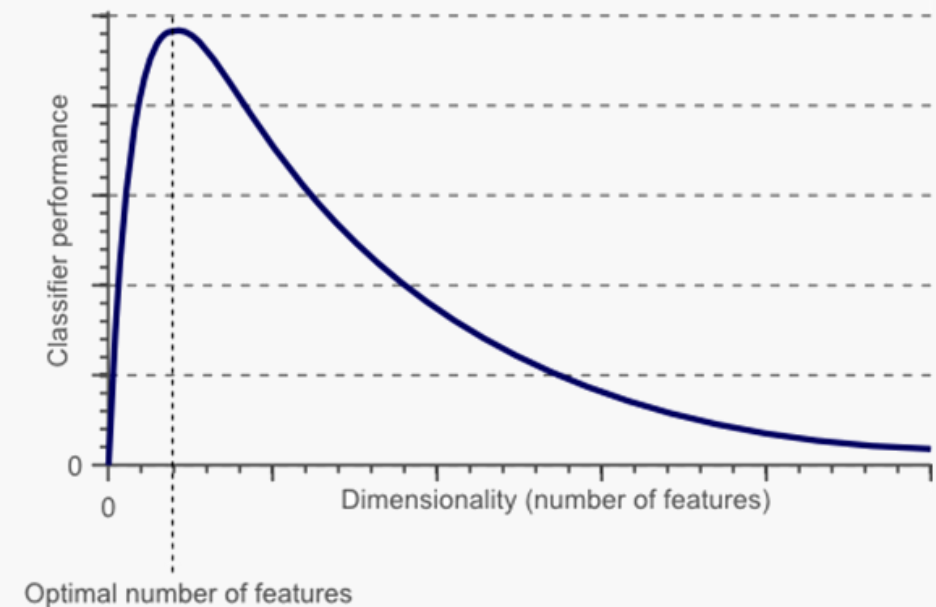
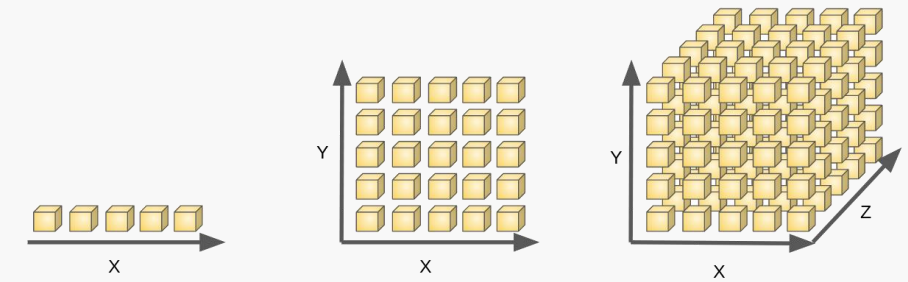
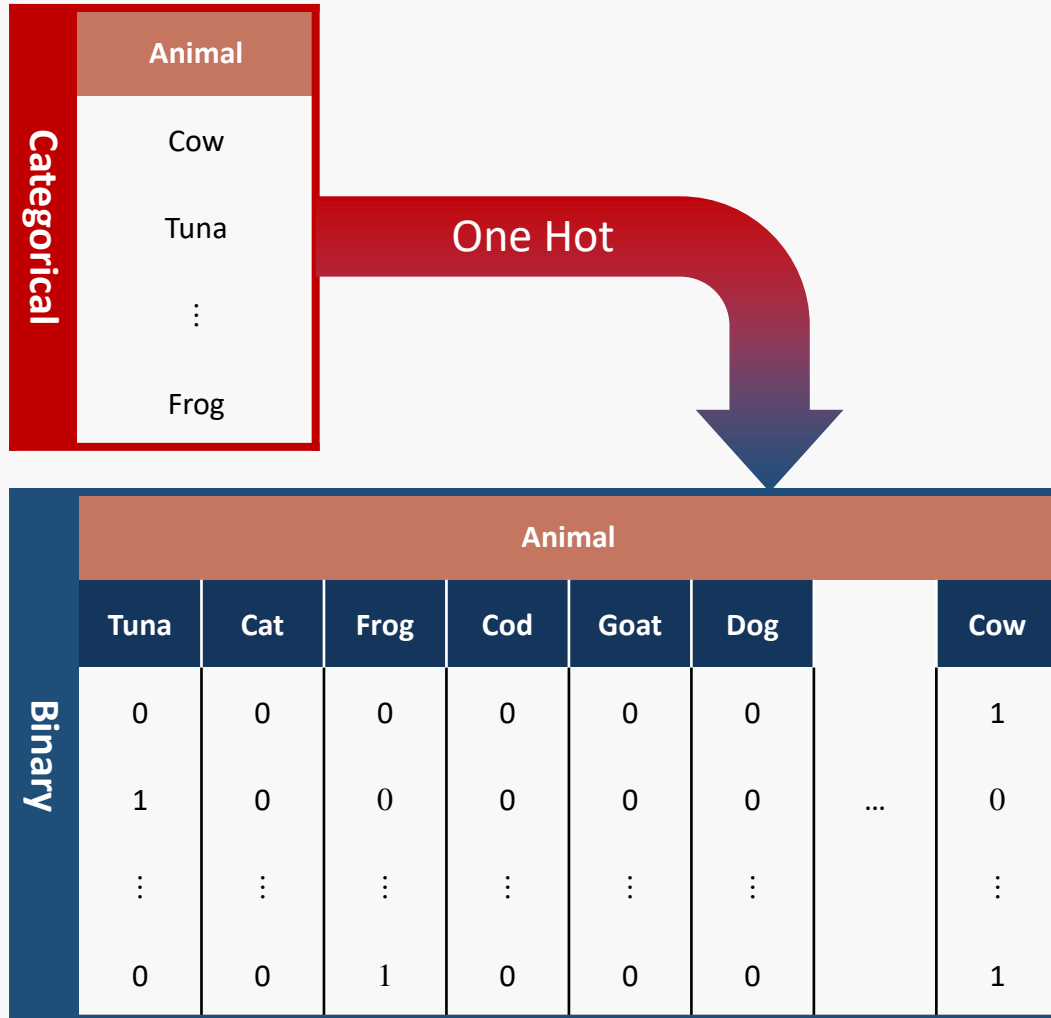
Yun-Hao Yang, Jakey Blue

# Motivation

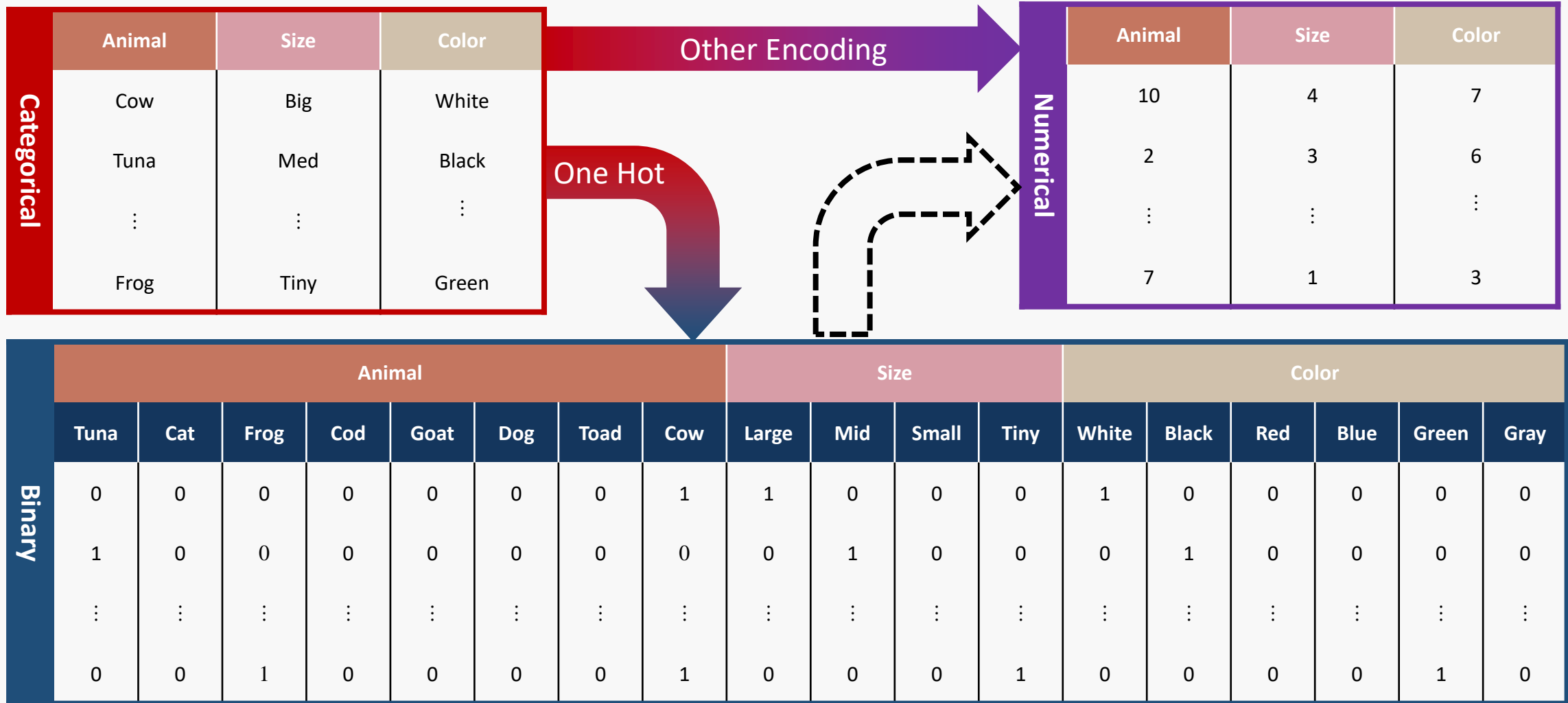


Categorical	Numerical								
City	Ordinal	Binary		One Hot				Frequency	
Taipei	0	0	0	0	0	0	1	0.2	
New York	1	0	1	0	0	1	0	0.2	
London	2	1	0	0	1	0	0	0.2	
Tokyo	3	1	1	1	0	0	0	0.4	
Tokyo	3	1	1	1	0	0	0	0.4	


# Motivation



# Motivation




# Motivation



Numerical	Group 1	Group 2	Group 3
	10	4	7
	2	3	6
	⋮	⋮	⋮
	7	1	3

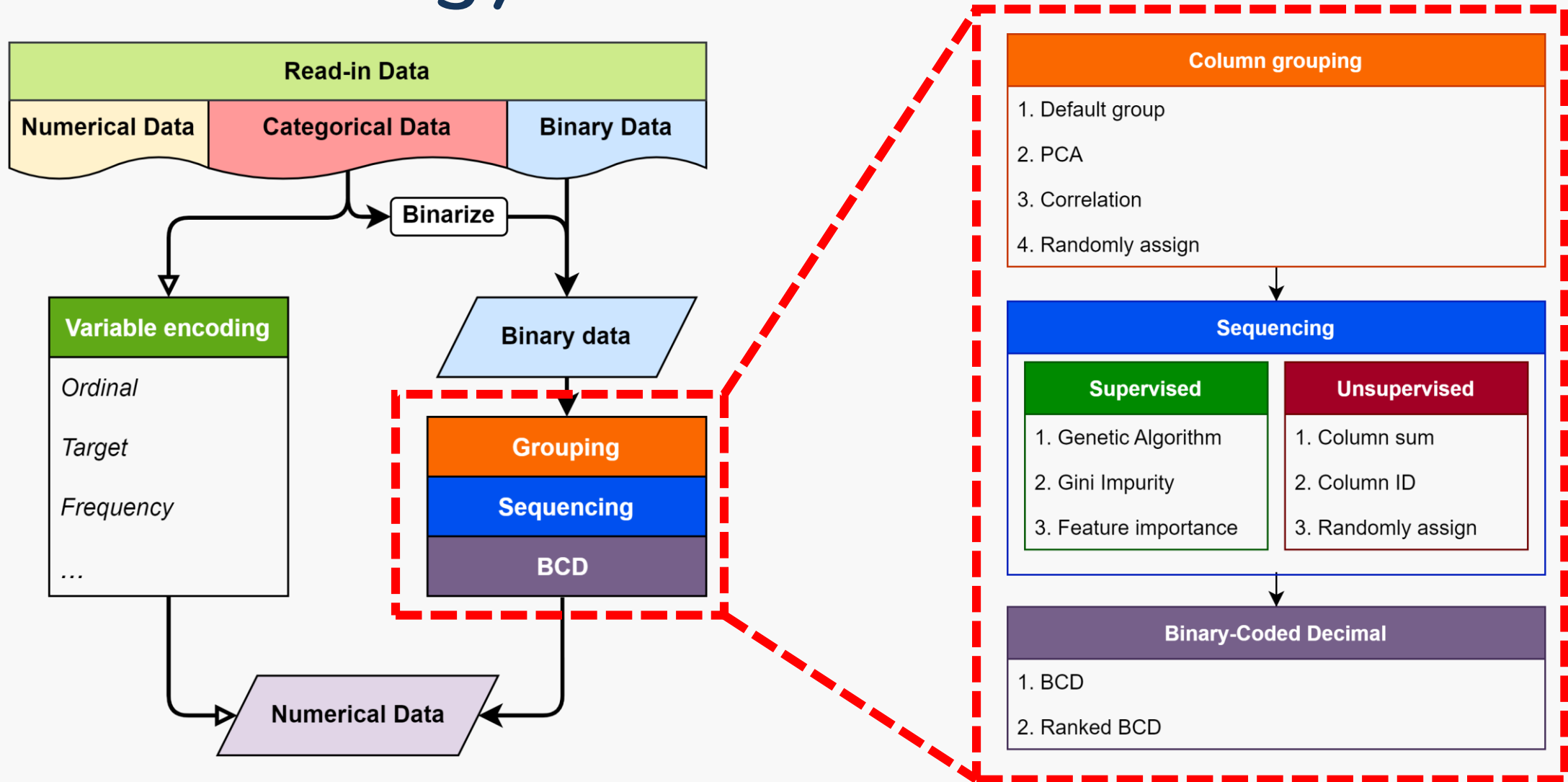


When the data is not follow the One Hot rule,  
can't use other method to transform the binary  
data. If only can we find a way in representing  
the numerous binary features data in a much  
simple way...



Binary	Tuna	Cat	Frog	Cod	Goat	Dog	Toad	Cow	Large	Mid	Small	Tiny	White	Black	Red	Blue	Green	Gray
	0	0	0	0	0	0	0	1	1	0	0	0	1	0	0	0	0	0
	1	0	0	0	0	0	0	0	0	1	0	0	0	1	0	0	0	0
	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮	⋮
	0	0	1	0	0	0	0	1	0	0	0	1	0	0	0	0	1	0

# Methodology



# Methodology

In the research, we propose a method to encode binary data in to numerical data. Via grouping, sequencing and transforming the binary row data with BCD encode.

1. Grouping similar, correlated features
2. Sequencing features in each feature group
3. BCD encode on each feature group

Tiny	Cat	Large	Black	White	Cow
1	1	0	1	0	0
0	0	1	0	1	1
⋮	⋮	⋮	⋮	⋮	⋮
0	0	1	0	1	1

# Methodology

In the research, we propose a method to encode binary data in to numerical data. Via grouping, sequencing and transforming the binary row data with BCD encode.

- 1. Grouping similar, correlated features
- 2. Sequencing features in each feature group
- 3. BCD encode on each feature group

Tiny	Cat	Large	Black	White	Cow
1	1	0	1	0	0
0	0	1	0	1	1
⋮	⋮	⋮	⋮	⋮	⋮
0	0	1	0	1	1



# Methodology

In the research, we propose a method to encode binary data in to numerical data. Via grouping, sequencing and transforming the binary row data with BCD encode.

- 1. Grouping similar, correlated features
- 2. Sequencing features in each feature group
- 3. BCD encode on each feature group

Cat	Cow	Black	White	Tiny	Large
1	0	1	0	1	0
0	1	0	1	0	1
⋮	⋮	⋮	⋮	⋮	⋮
0	1	0	1	0	1

# Methodology

In the research, we propose a method to encode binary data in to numerical data. Via grouping, sequencing and transforming the binary row data with BCD encode.

1. Grouping similar, correlated features
2. Sequencing features in each feature group
3. BCD encode on each feature group

Cow	Cat	White	Black	Tiny	Large
0	1	0	1	1	0
1	0	1	0	0	1
⋮	⋮	⋮	⋮	⋮	⋮
1	0	1	0	0	1

# Methodology

In the research, we propose a method to encode binary data in to numerical data. Via grouping, sequencing and transforming the binary row data with BCD encode.

- 1. Grouping similar, correlated features
- 2. Sequencing features in each feature group
- 3. BCD encode on each feature group

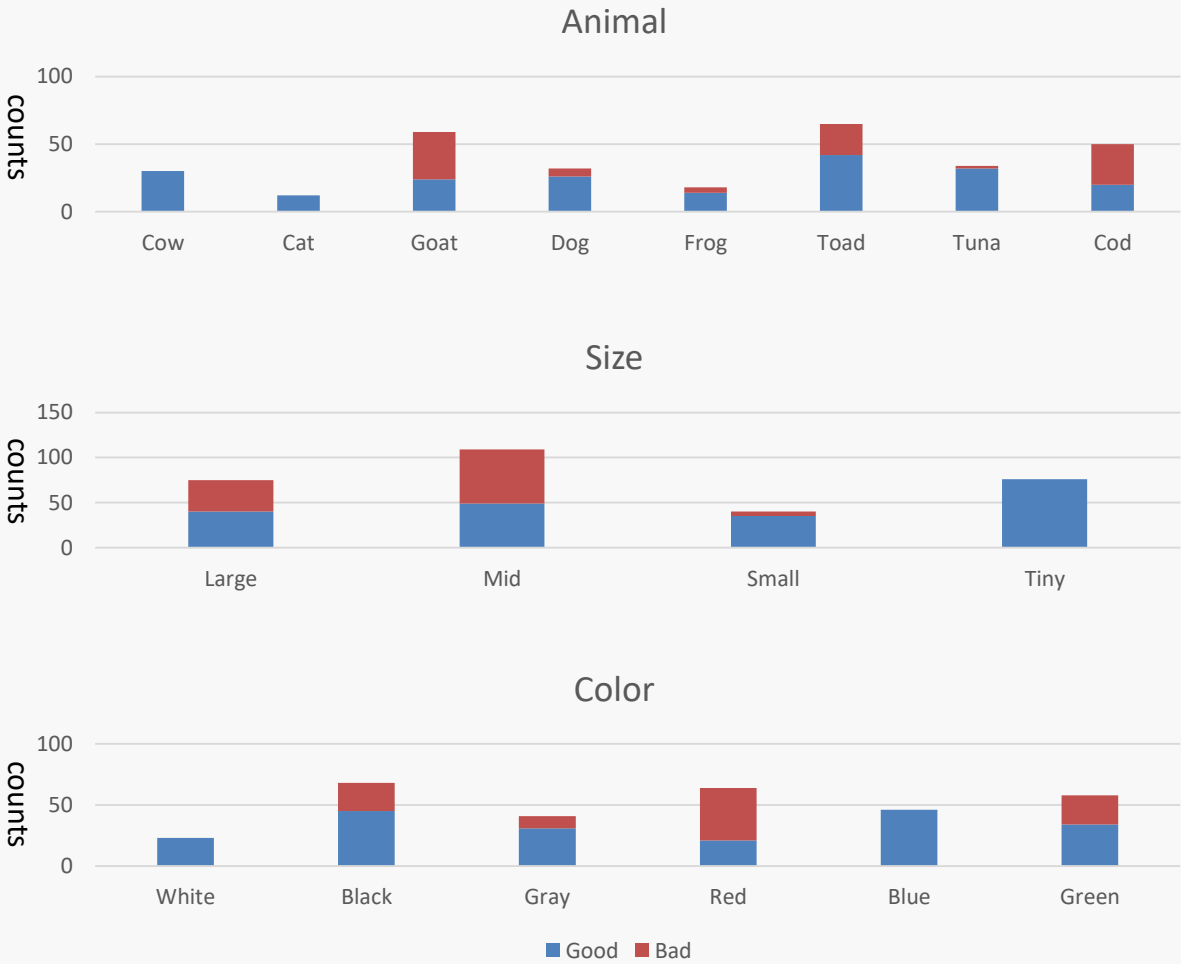
Animal	Color	Size
1	1	2
2	2	1
⋮	⋮	⋮
2	2	1

# Methodology - Grouping



Animal							
Cow	Cat	Goat	Dog	Frog	Toad	Tuna	Cod
		V					
Size							
Large		Mid		Small		Tiny	
		V					
Color							
White	Black	Gray	Red	Blue	Green		
	V						
Health Status							
Good				Bad			
				V			

Binary Data (Total:300 sample, 200 good, 100 bad)



# Methodology - Grouping

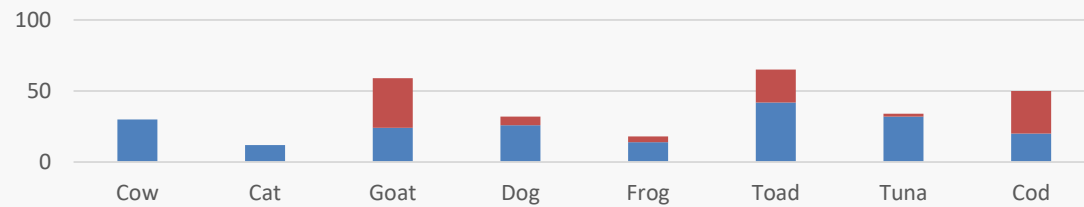
Group

Sequence

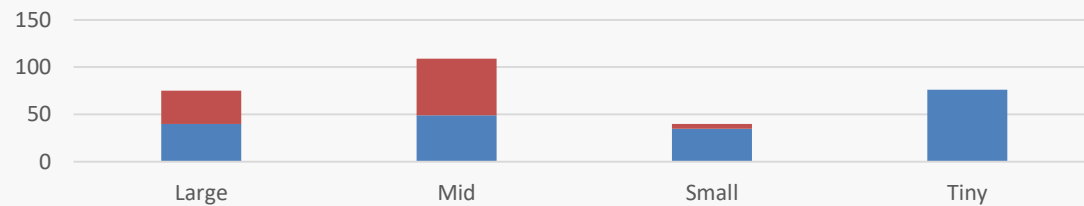
BCD

Binary Data

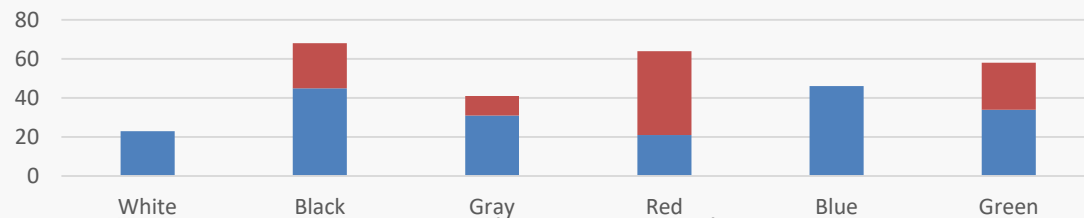
Animal



Size

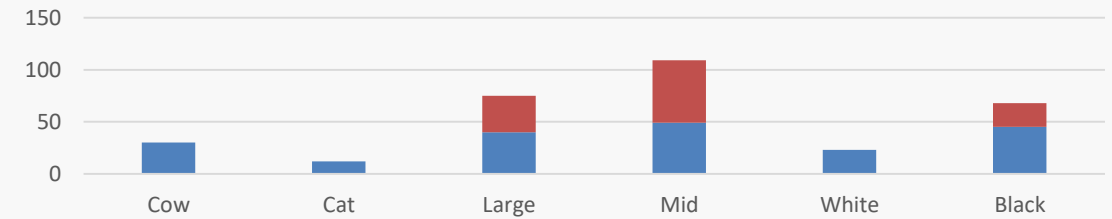


Color

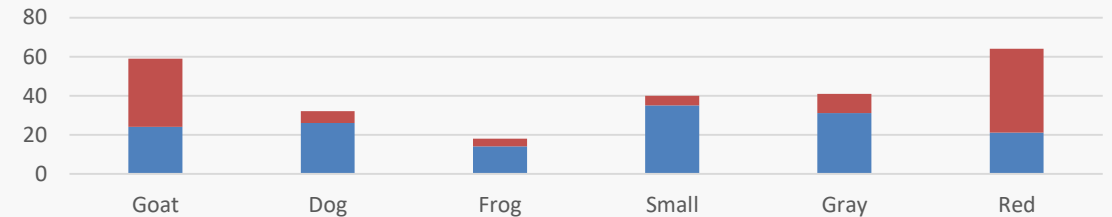


Grouped by PCA

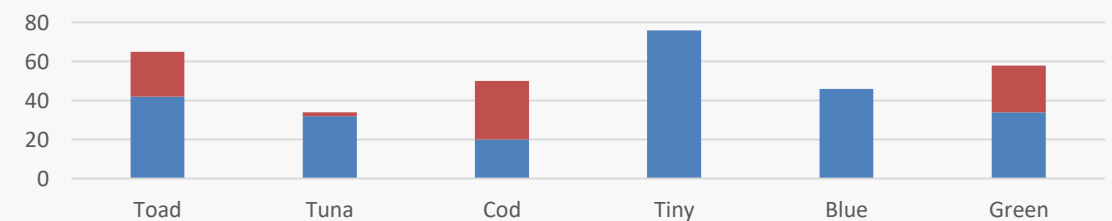
PC 1



PC 2



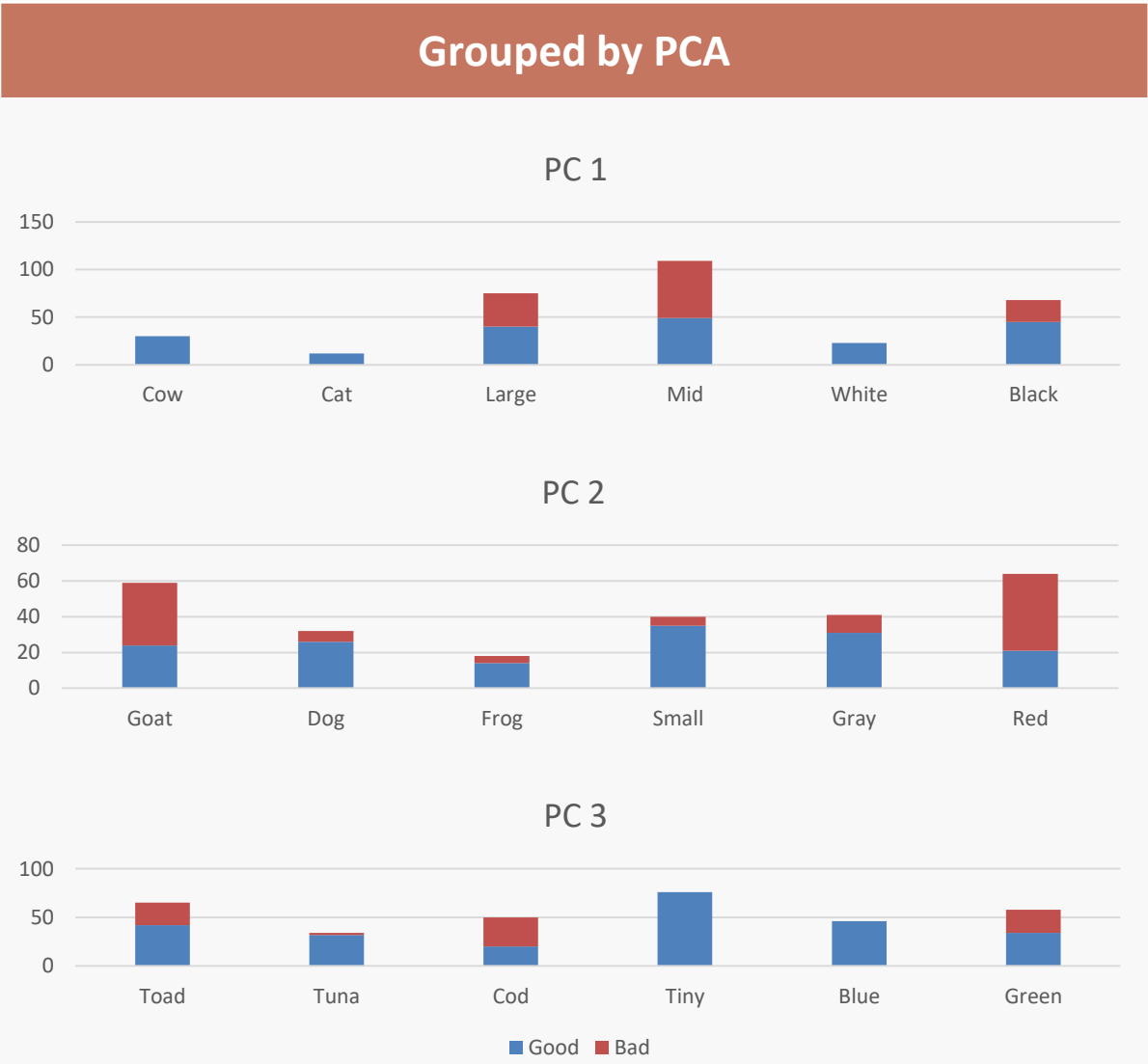
PC 3



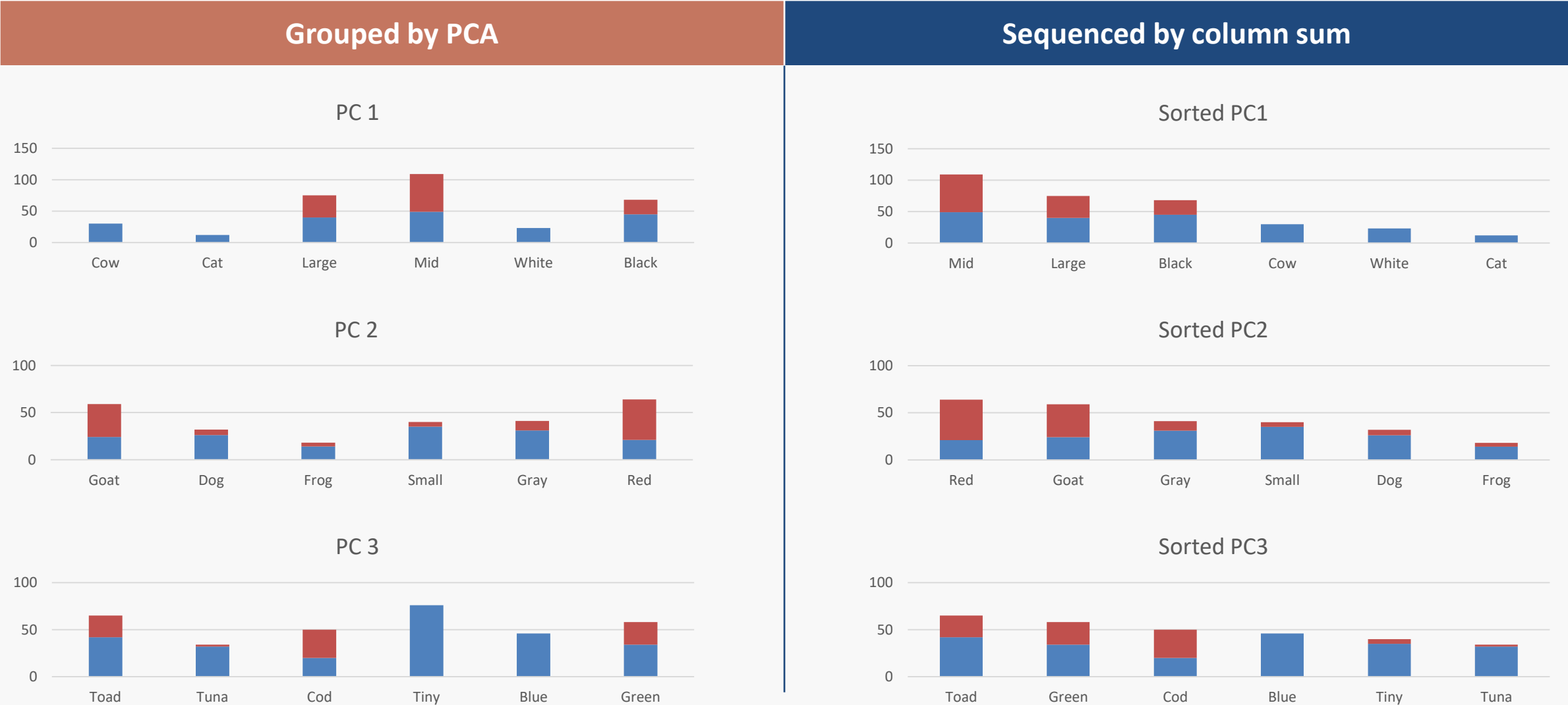
# Methodology - Sequencing

Secondly, sequencing features in each group by columns' attributes, for outputting better BCD values after encode.

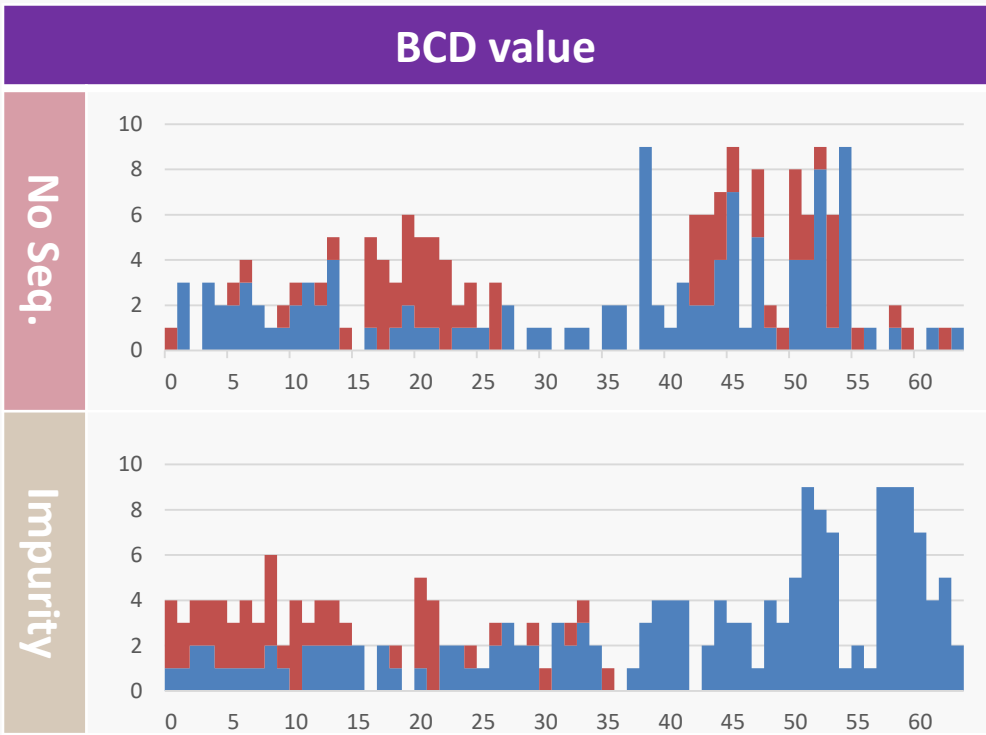
- 1. Column sum
- 2. Type impurity
- 3. Feature importance



# Methodology - Sequencing

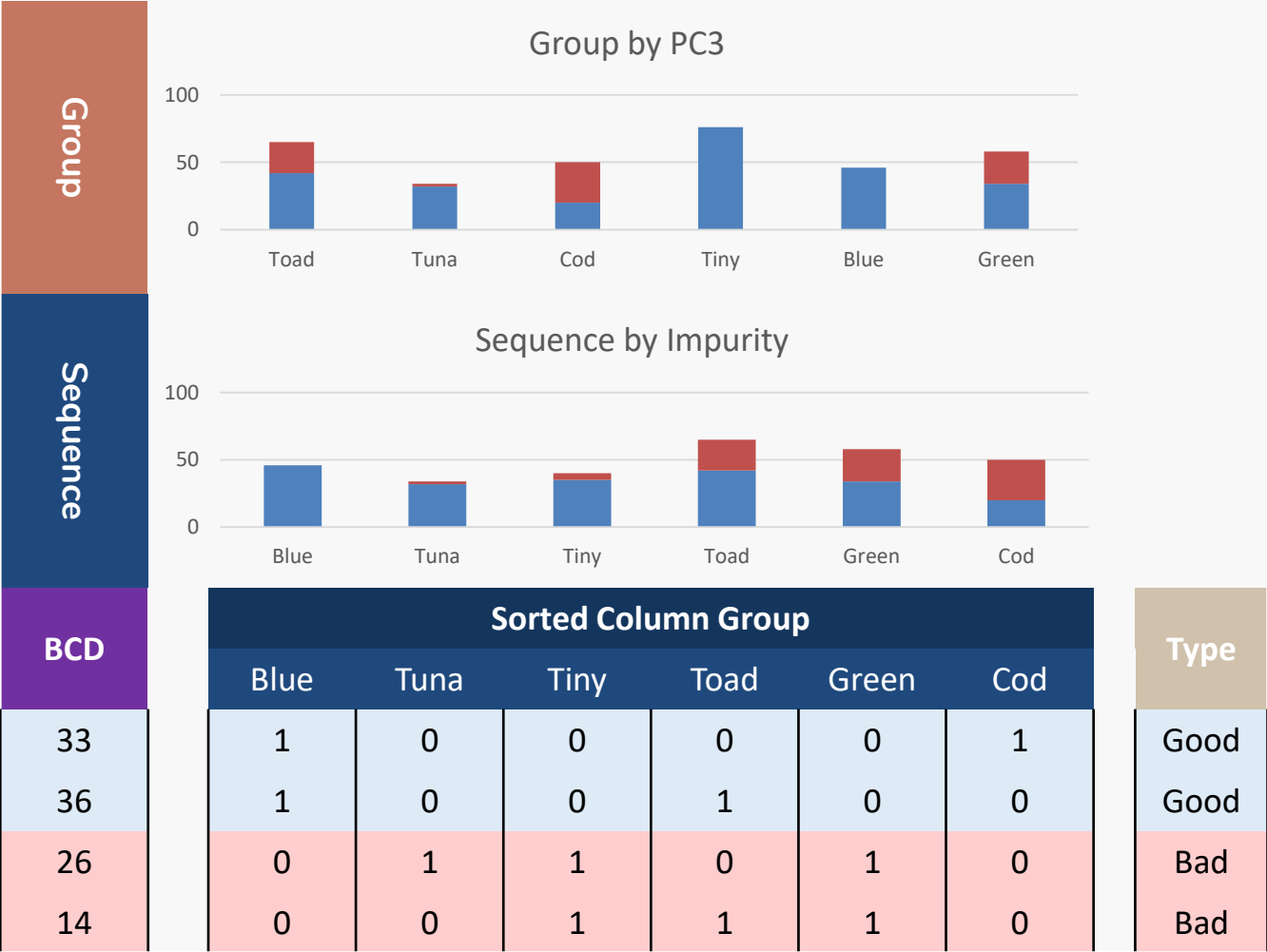


# Methodology - Sequencing



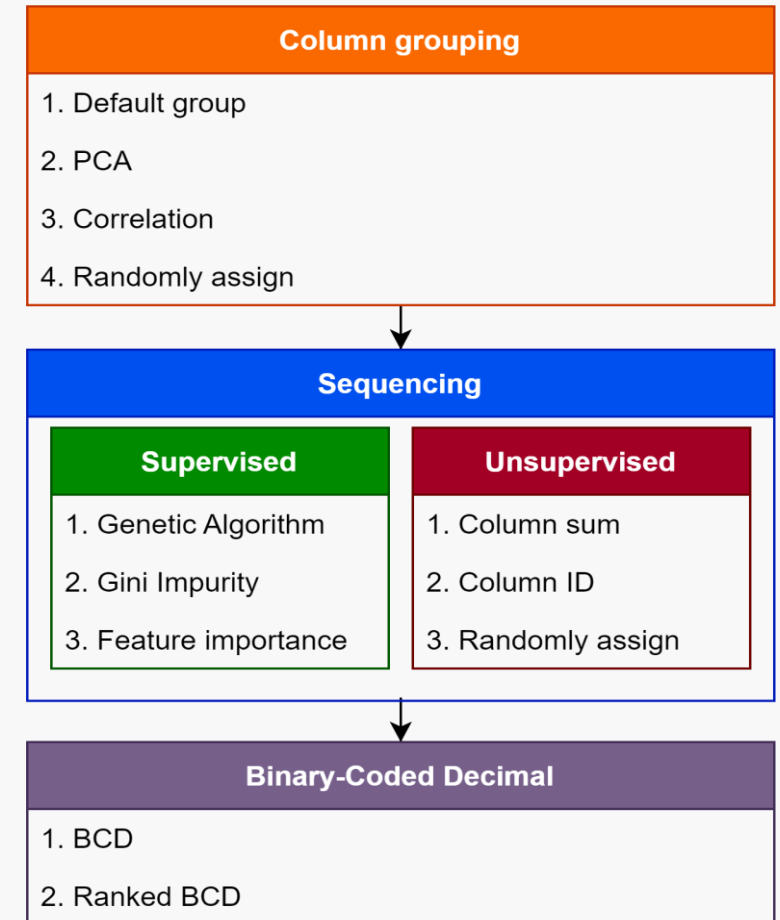
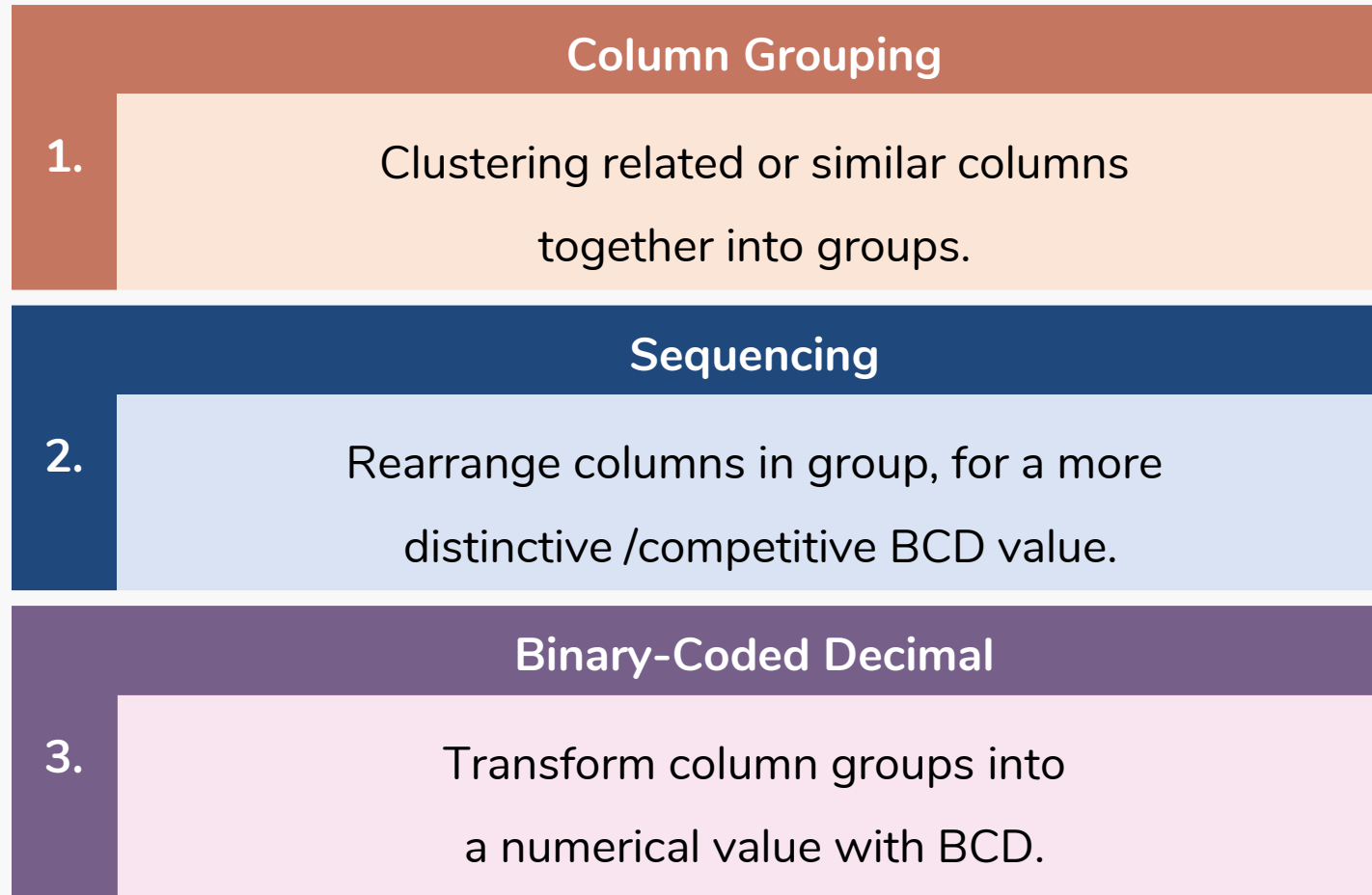


# Methodology - BCD code



Decimal digit	BCD			
	8	4	2	1
0	0	0	0	0
1	0	0	0	1
2	0	0	1	0
3	0	0	1	1
4	0	1	0	0
5	0	1	0	1
6	0	1	1	0
7	0	1	1	1
8	1	0	0	0
9	1	0	0	1

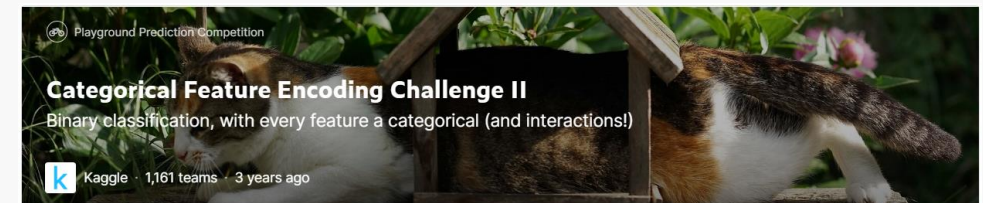
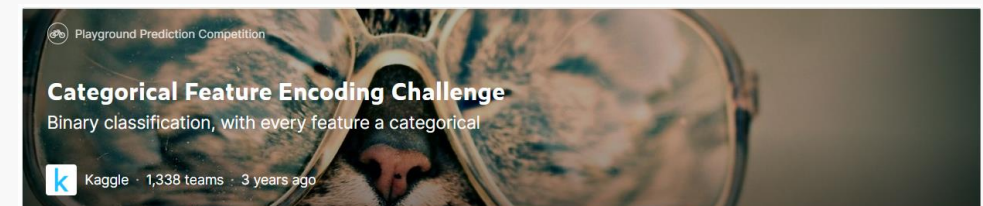
# Methodology - sum up



# Case study

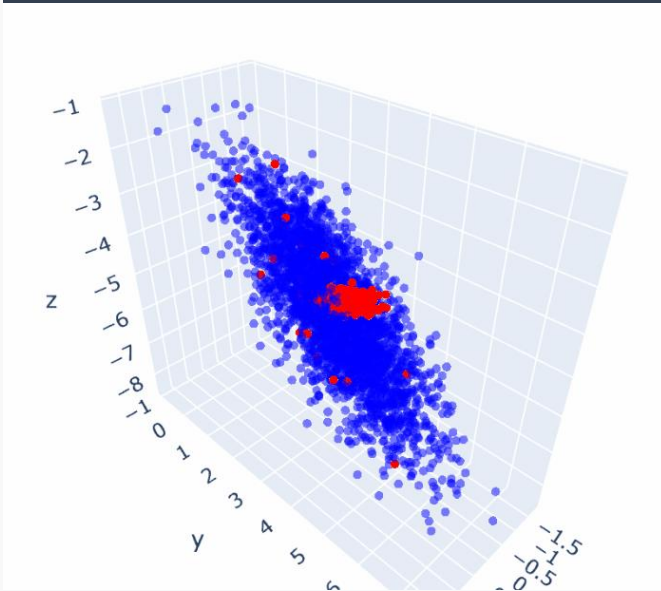
In Case study, we compare classification results with the commonly used variable encoding method under different datasets.

1. Simulated datasets
2. Kaggle dataset

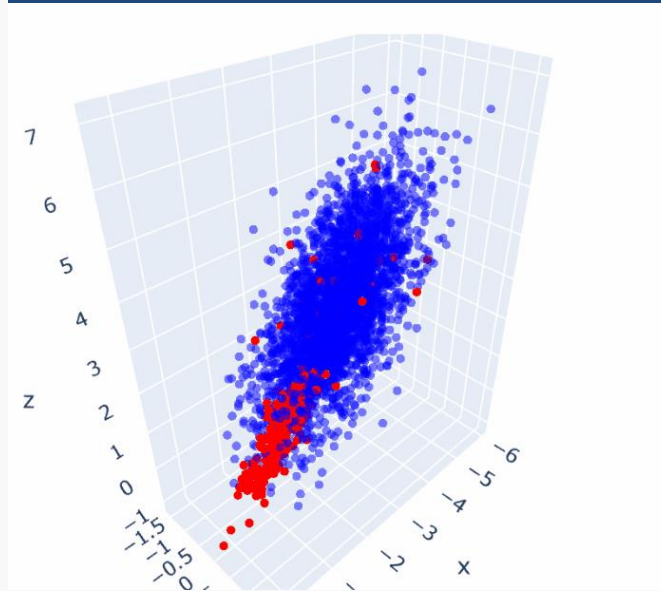


# Simulated data (3300 samples, 3000 Good, 300 Bad)

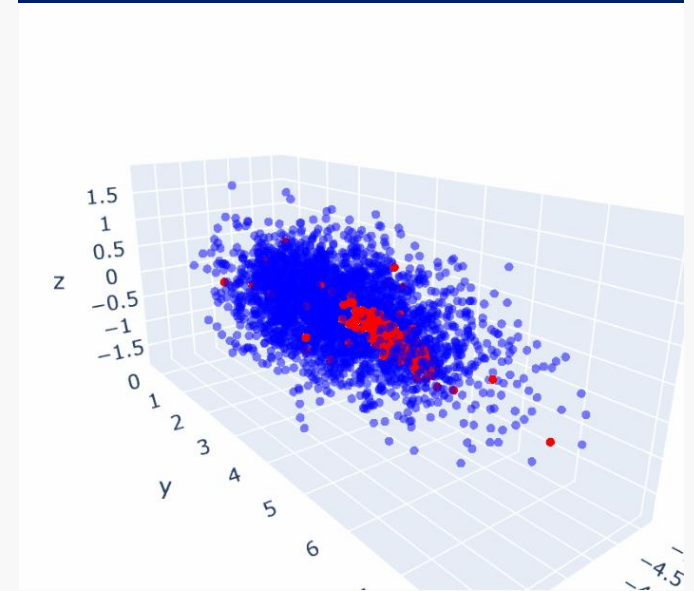
Dataset 1



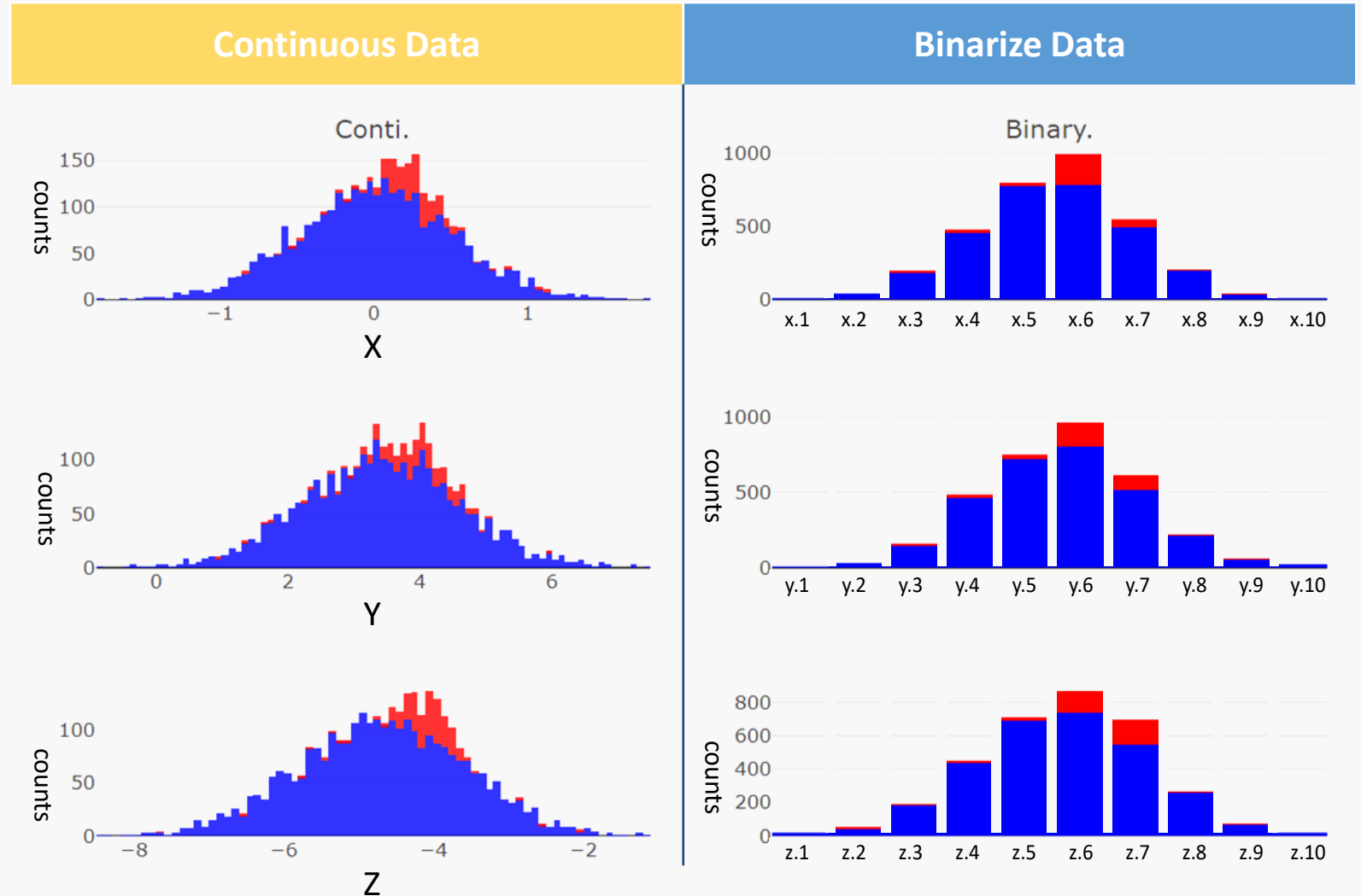
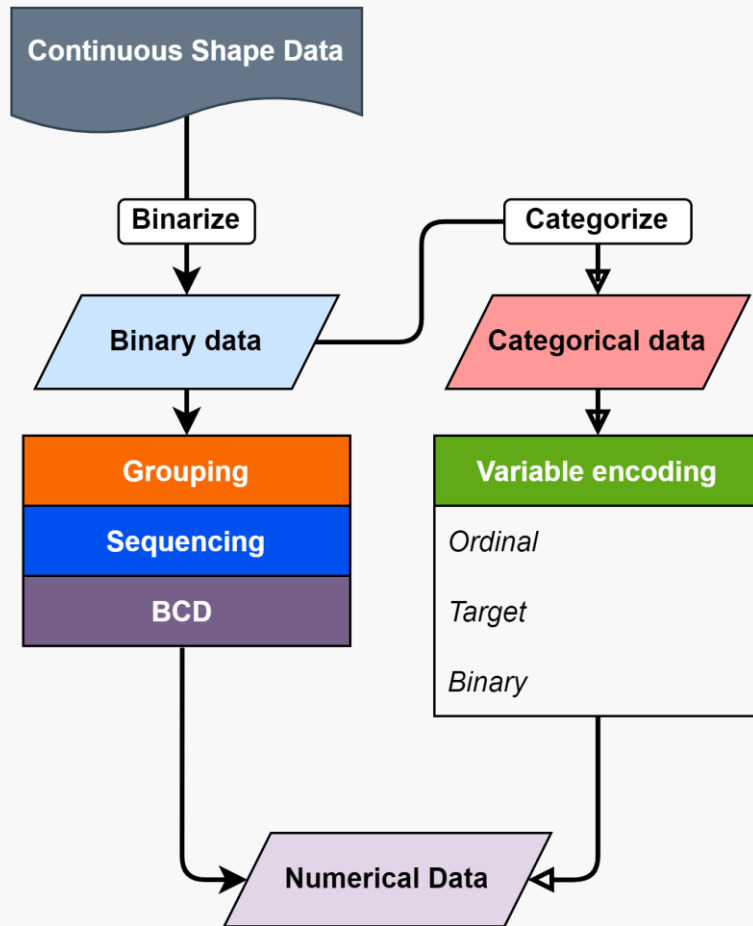
Dataset 2



Dataset 3

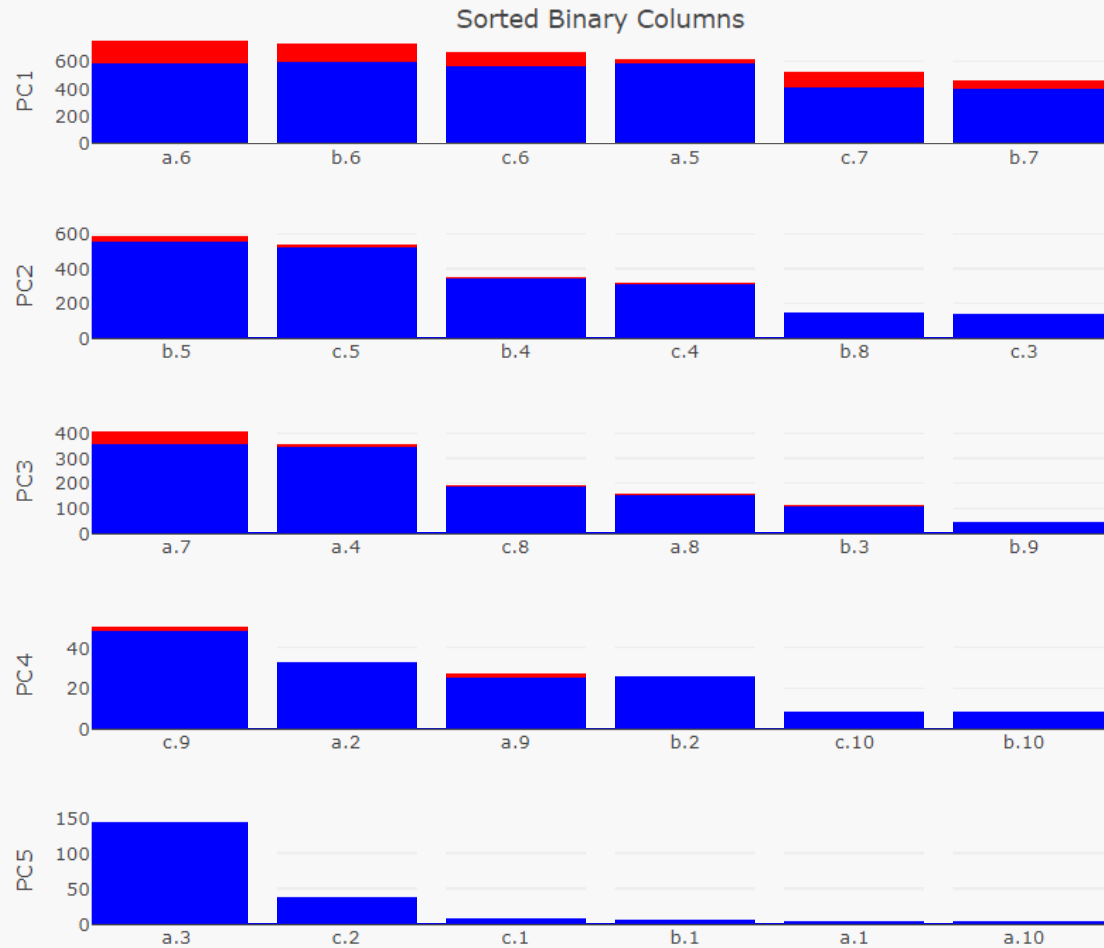


# Simulated data - binarize

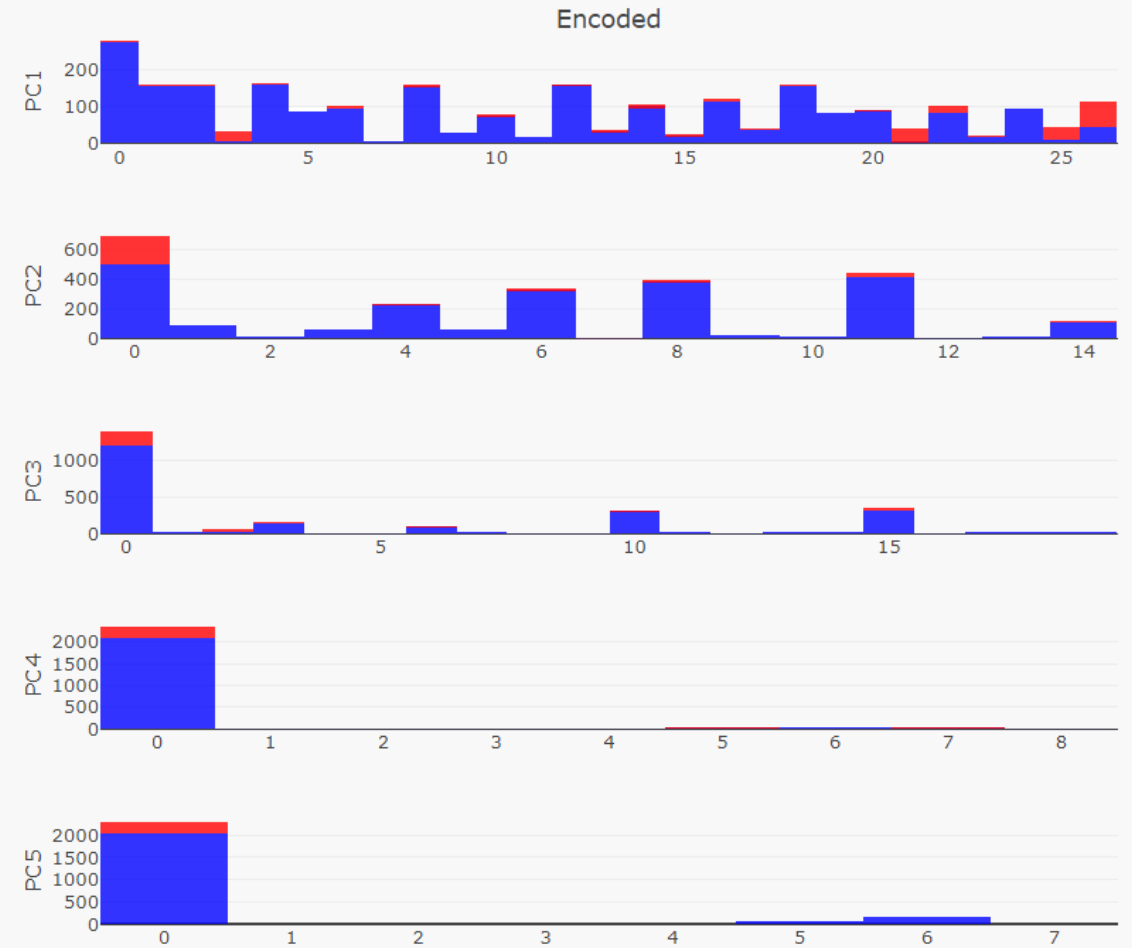


# Simulated data - seq. by sum

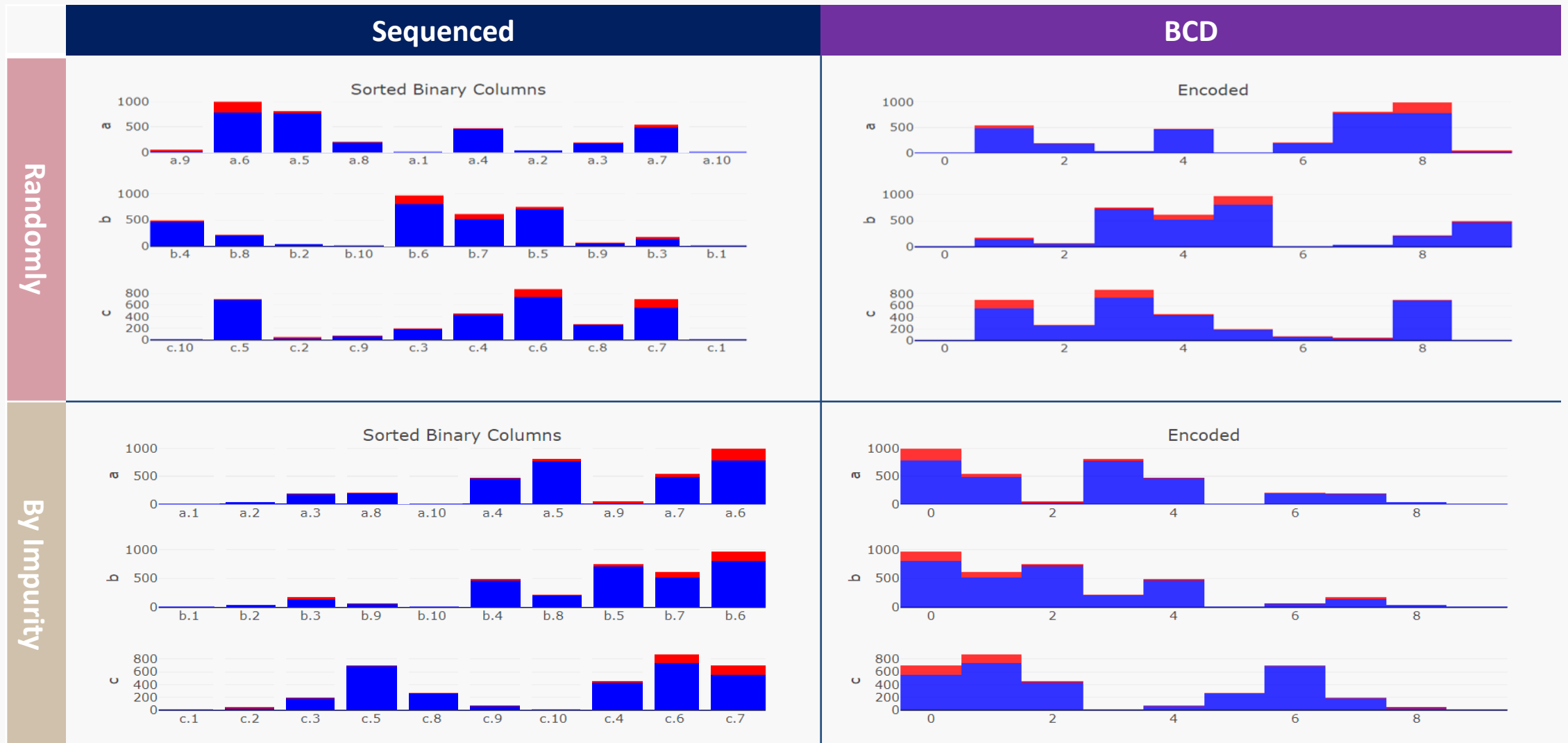
Sequence by column sum



BCD



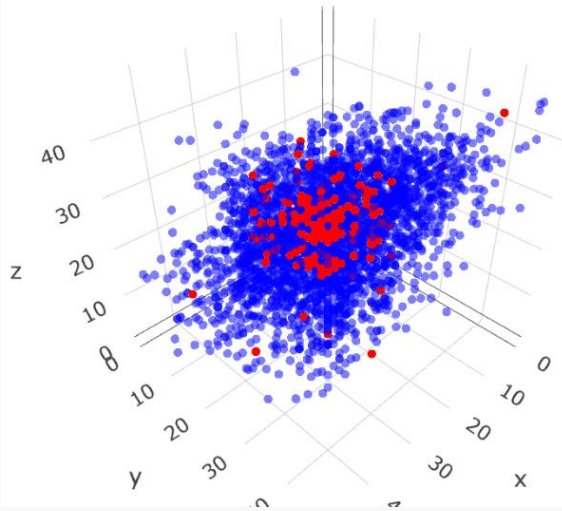
# Simulated data - seq. comparison



# Simulated data - Dimension reduction

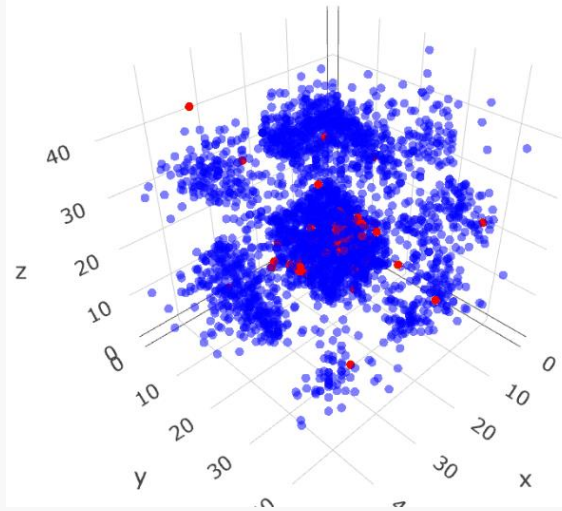
Sum

Dimension Reduction Data  
shape 1 / Default / Sum



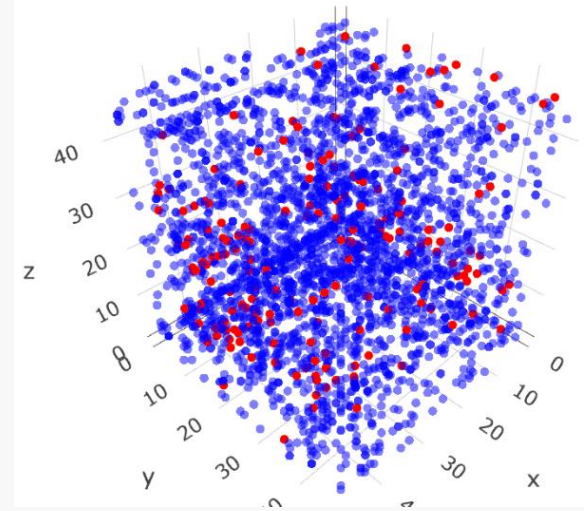
Impurity

Dimension Reduction Data  
shape 1 / Default / gini



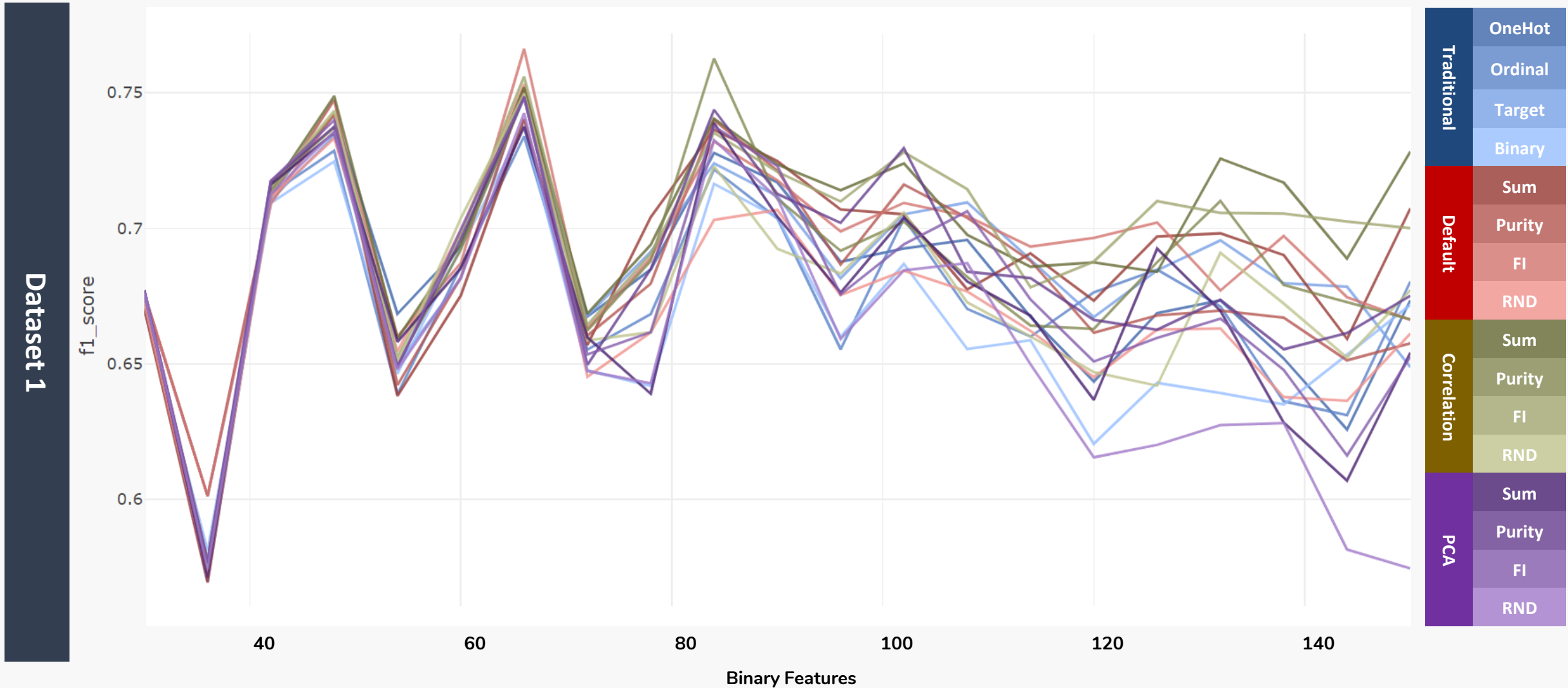
RND

Dimension Reduction Data  
shape 1 / Default / Rnd

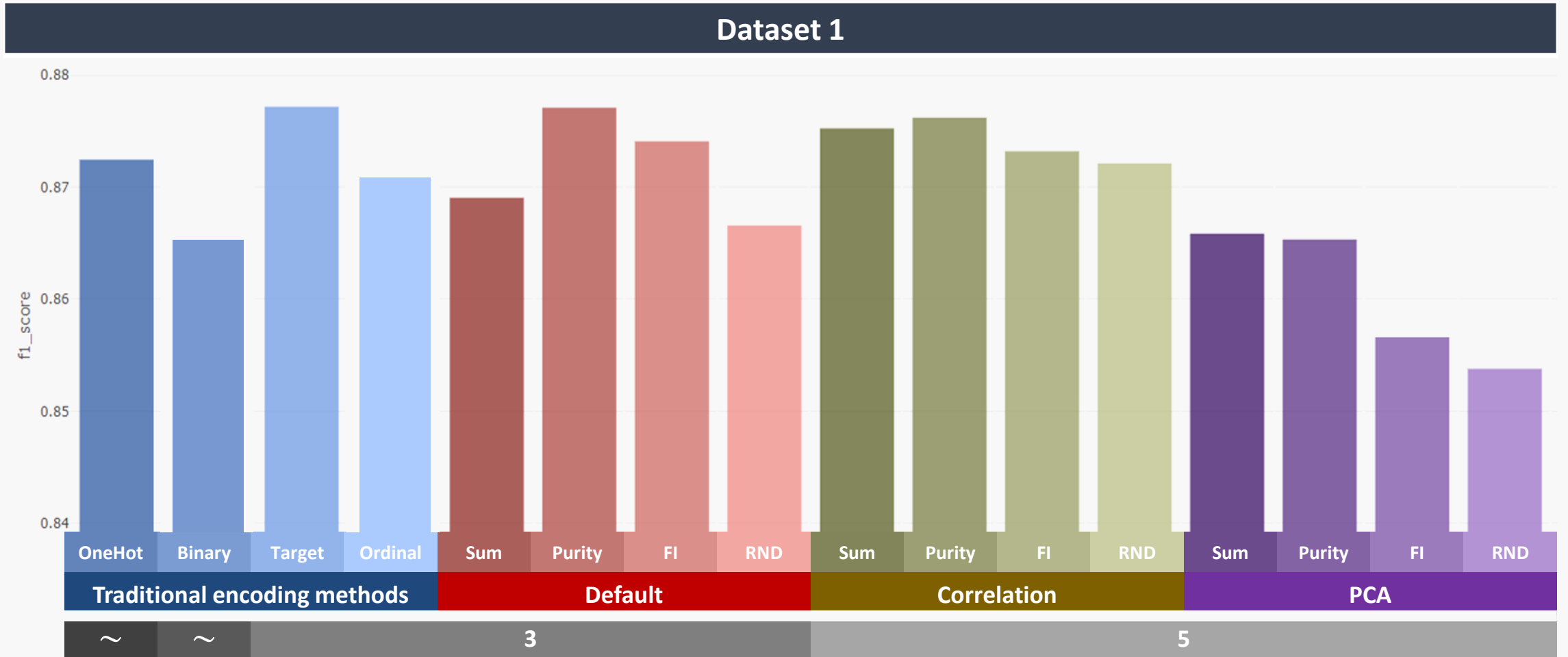




# Continuous data - Classification

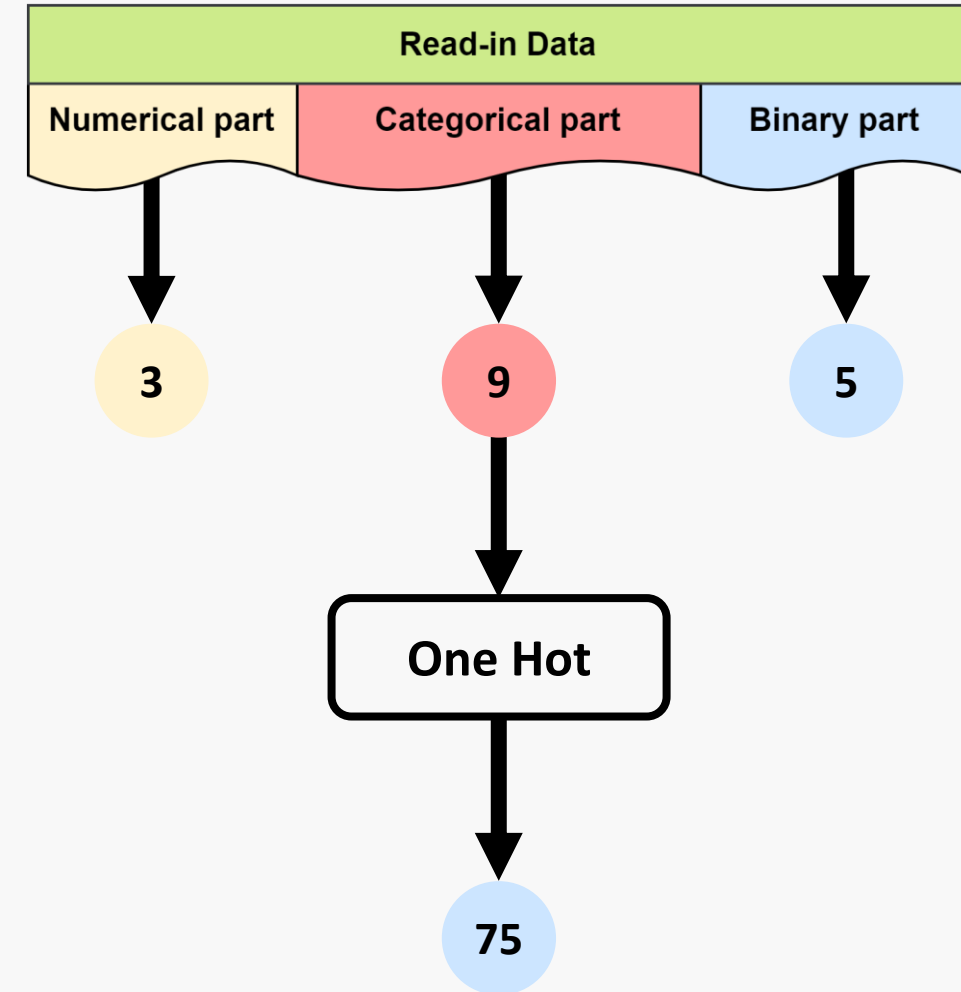


# Simulated data - Classification

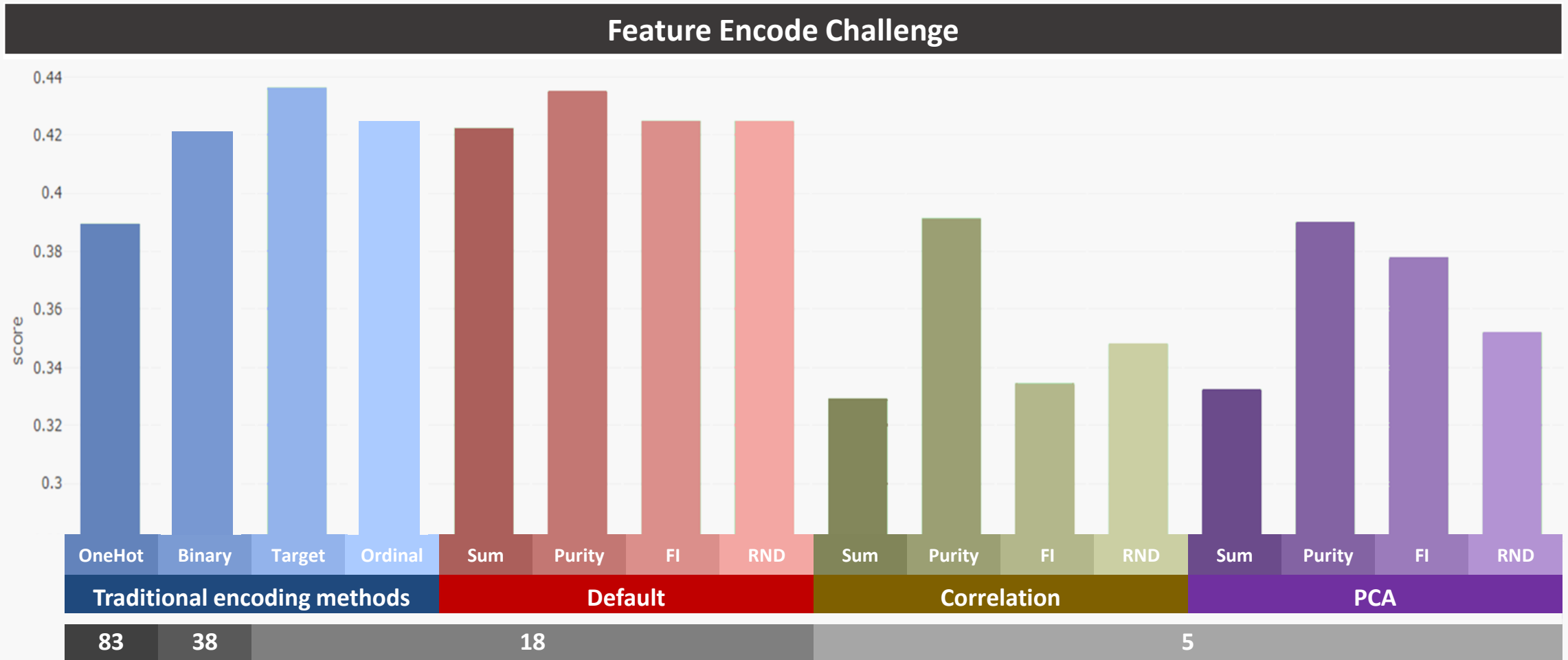


# Kaggle dataset

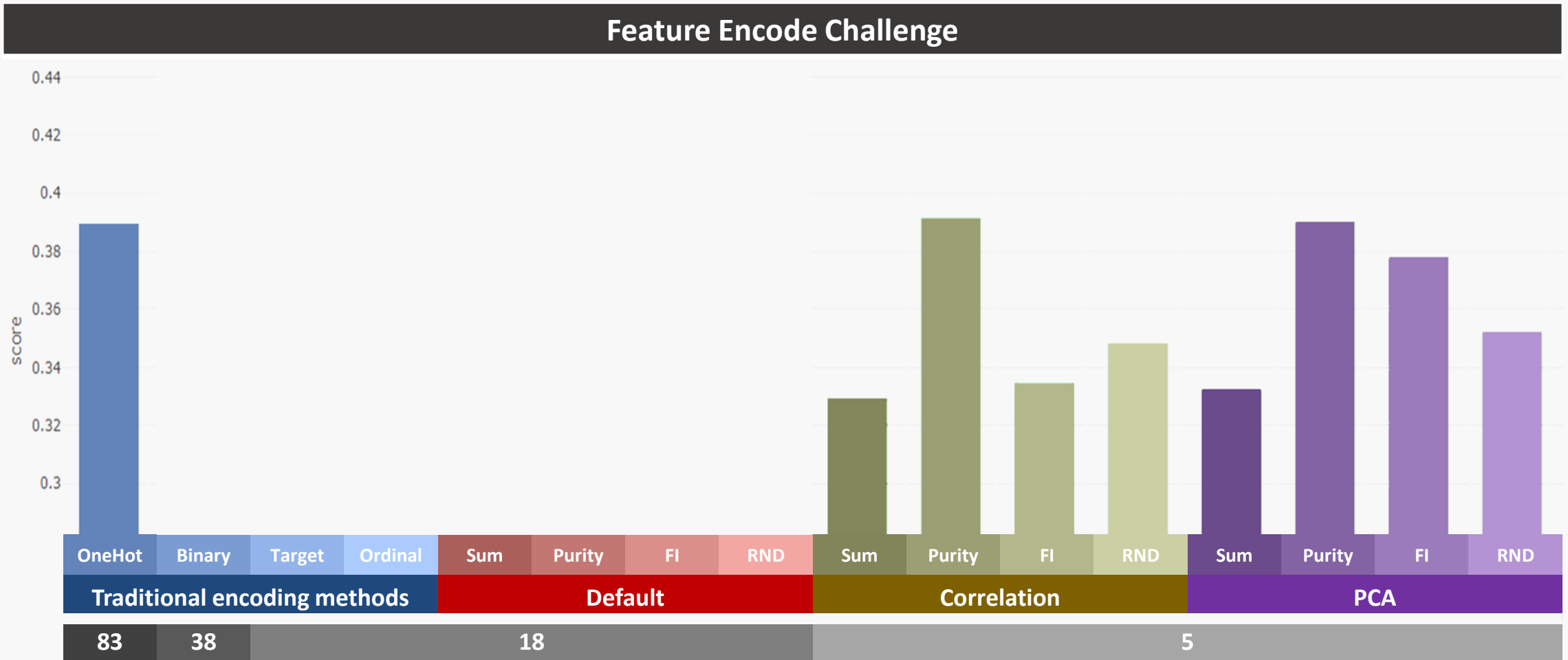
The screenshot shows the Kaggle 'Categorical Feature Encoding Challenge' page. The header includes the competition title, a description ('Binary classification, with every feature a categorical'), and the Kaggle logo with '1,338 teams · 3 years ago'. Below the header is a navigation bar with links: Overview, Data, Code, Discussion, Leaderboard, Rules, Team, Submissions, and a 'Late Submission' button. The 'Overview' section is active, showing a 'Description' tab. The description text asks 'Is there a cat in your dat?' and explains the task of encoding categorical variables. It lists features: binary, low- and high-cardinality nominal, low- and high-cardinality ordinal, and (potentially) cyclical. A photo of a cat is shown. The text continues: 'This Playground competition will give you the opportunity to try different encoding schemes for different algorithms to compare how they perform. We encourage you to share what you find with the community. If you're not sure how to get started, you can check out the [Categorical Variables](#) section of Kaggle's [Intermediate Machine Learning course](#). Have Fun!'



# Kaggle dataset - Classification



# Kaggle dataset - Classification



# Conclusion

In this research, we develop a method to encode binary feature data into numerical data, which can compress the binary features information to reduce dimension while remaining a certain level in ML model classification performance. The main advantage of this method compare to traditional encoding methods are:

1. Encoding Binary features without group knowledge to numerical data.
2. Adjust the number of features of encoded numerical data.
3. Remain ML model performance with certain sequencing method.