



A CLOUD GURU

Numeric Feature Engineering



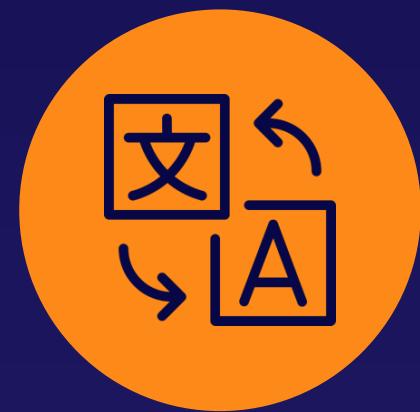
Brock Tubre

INSTRUCTOR



Numeric Feature Engineering

Transforming numeric values within our data so Machine Learning algorithms can better analyze them.



Changing numeric values in our datasets so they are easier to work with.

- **Feature Scaling**

Changes numeric values so all values are on the same scale.

- Normalization
- Standardization

- **Binning**

Changes numeric values into groups or buckets of similar values.

- Quantile Binning aims to assign the same number of features to each bin.

scaling = feature scaling = normalization



Problem

Will a home loan be approved?

ID	Type	Bedrooms	Area	Pool_Size	Price	Loan_Approved
1	condo	2	2432	S	250555	N
2	house	4	2988	L	243566	Y
3	house	3	1877	N	125700	N
4	condo	5	3876	M	345000	Y
5	apartment	2	1250	N	120900	Y



Problem

Will a home loan be approved?

ID	Type	Bedrooms	Area	Pool_Size	Price	Loan_Approved
1	condo	2	2432	S	250555	N
2	house	4	2988	L	243566	Y
3	house	3	1877	N	125700	N
4	condo	5	3876	M	345000	Y
5	apartment	2	1250	N	120900	Y

Normalization

Price



Normalization

Price



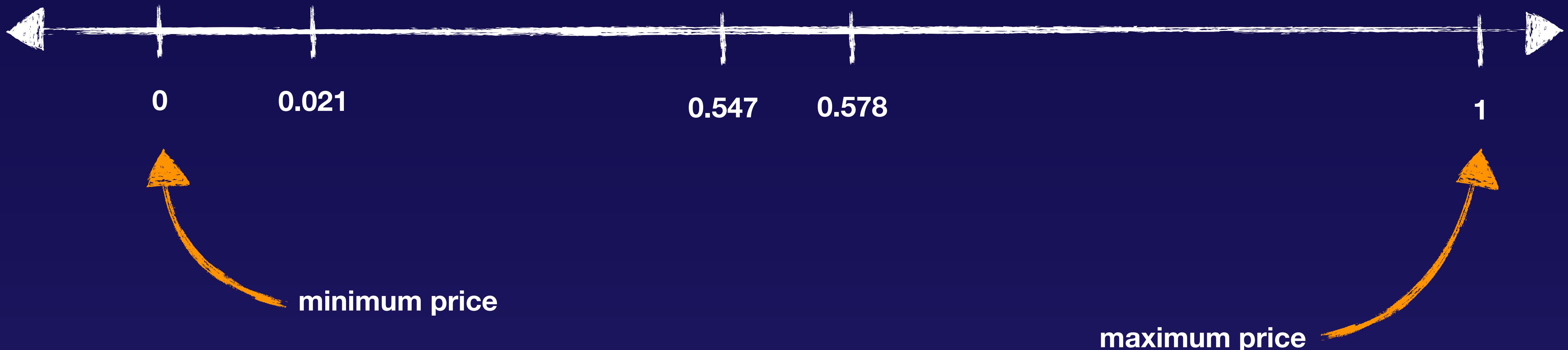
Normalization

Price



Normalization

Price



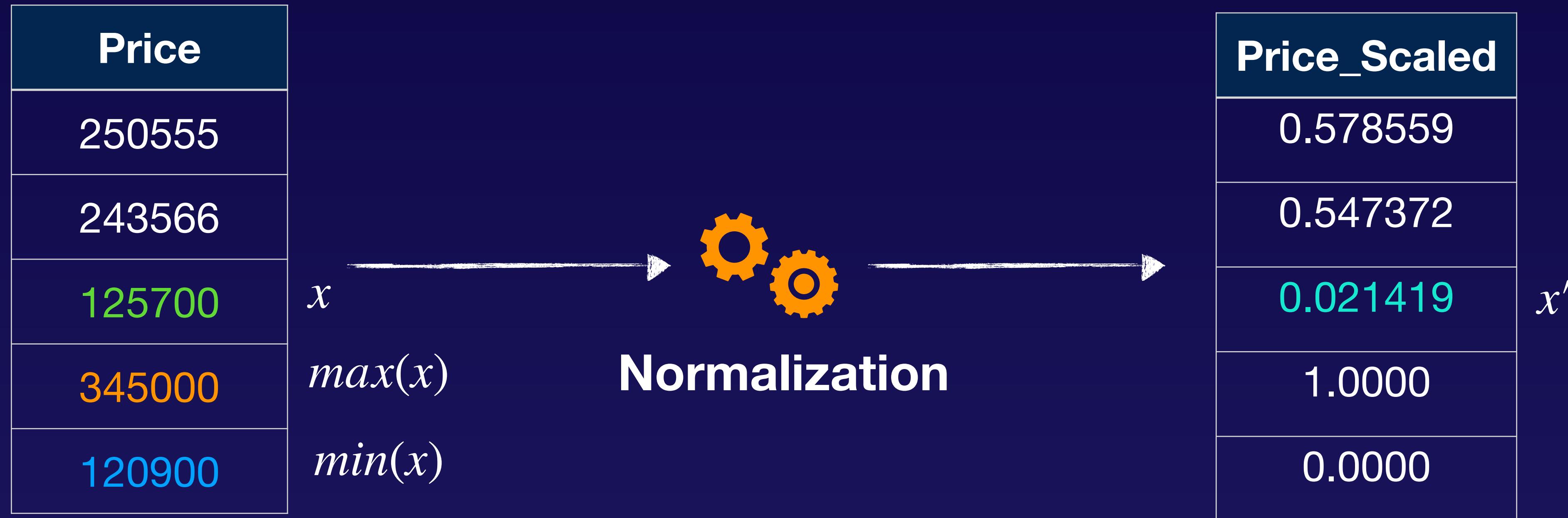
Normalization



Normalization



Normalization



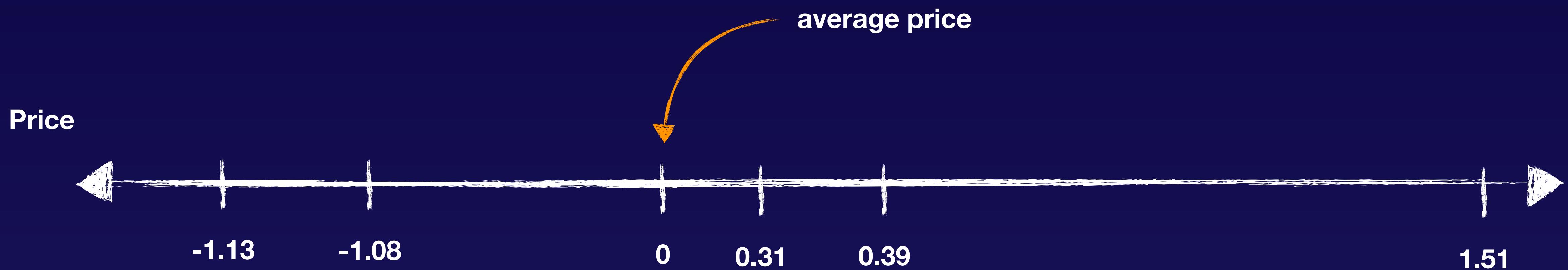
$$x' = \frac{x - \min(x)}{\max(x) - \min(x)}$$



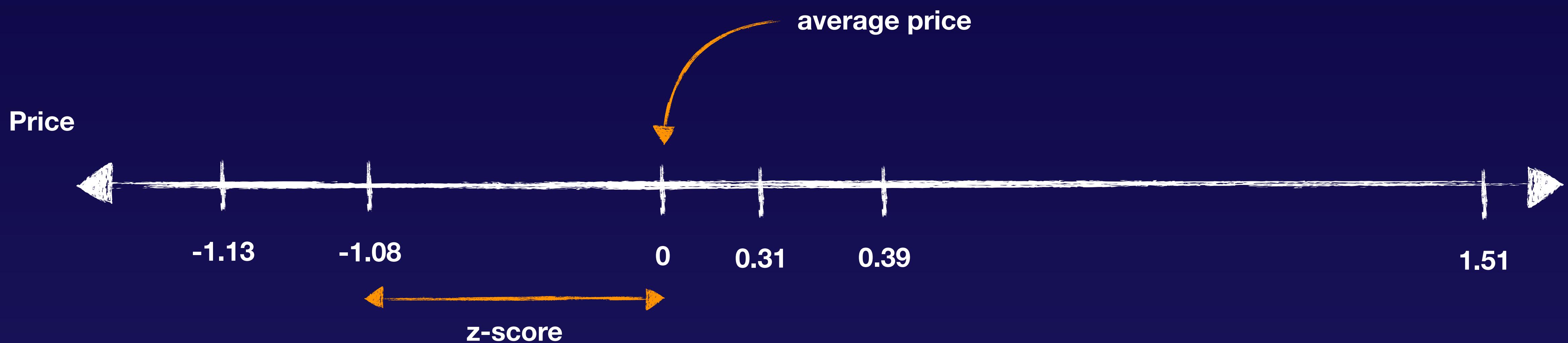
$$\frac{125700 - 120900}{345000 - 120900} = 0.021419$$

Outliers can throw off normalization!

Standardization



Standardization



Standardization



Standardization

Price
250555
243566
125700
345000
120900



Standardization

\bar{x} = Price mean: 217144.2

σ = Price std: 84600.81

$$z = \frac{x - \bar{x}}{\sigma}$$

125700 - 217144.2
84600.81 = -1.080890

Price_Scaled
0.394923
0.312311
-1.080890
1.511283
-1.137627

1

Scaling Features

Is required for many algorithms like linear/non-linear regression, clustering, neural networks, and more. Scaling features depends on the algorithm you use.

2

Normalization

Rescales values from 0 to 1 but doesn't handle outliers very well.

3

Standardization

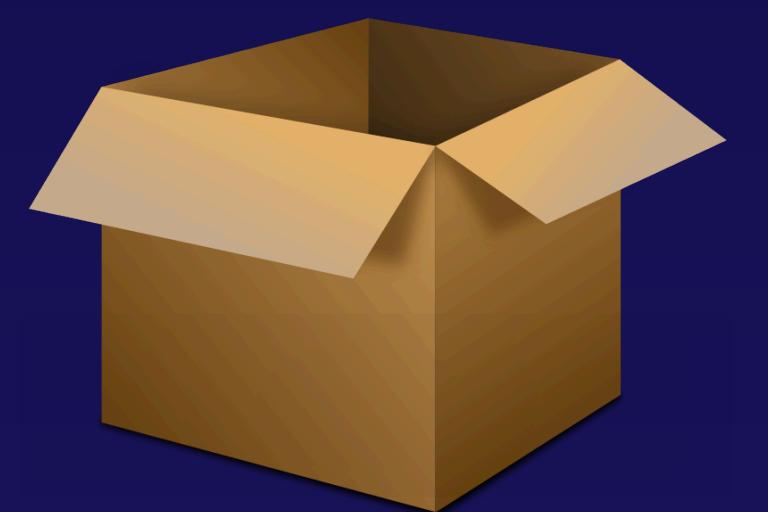
Rescales values by making the values of each feature in the data have zero mean and is much less affected by outliers.

4

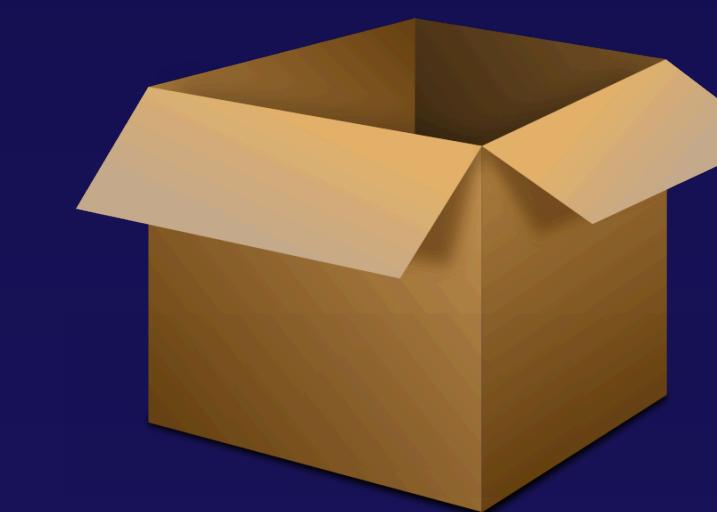
Translating Back

When you're done you can always scale back the data to the original representation.

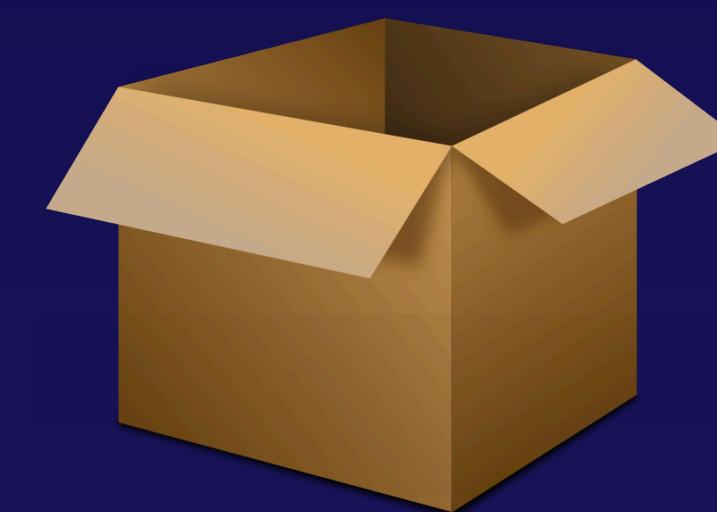




30s and below



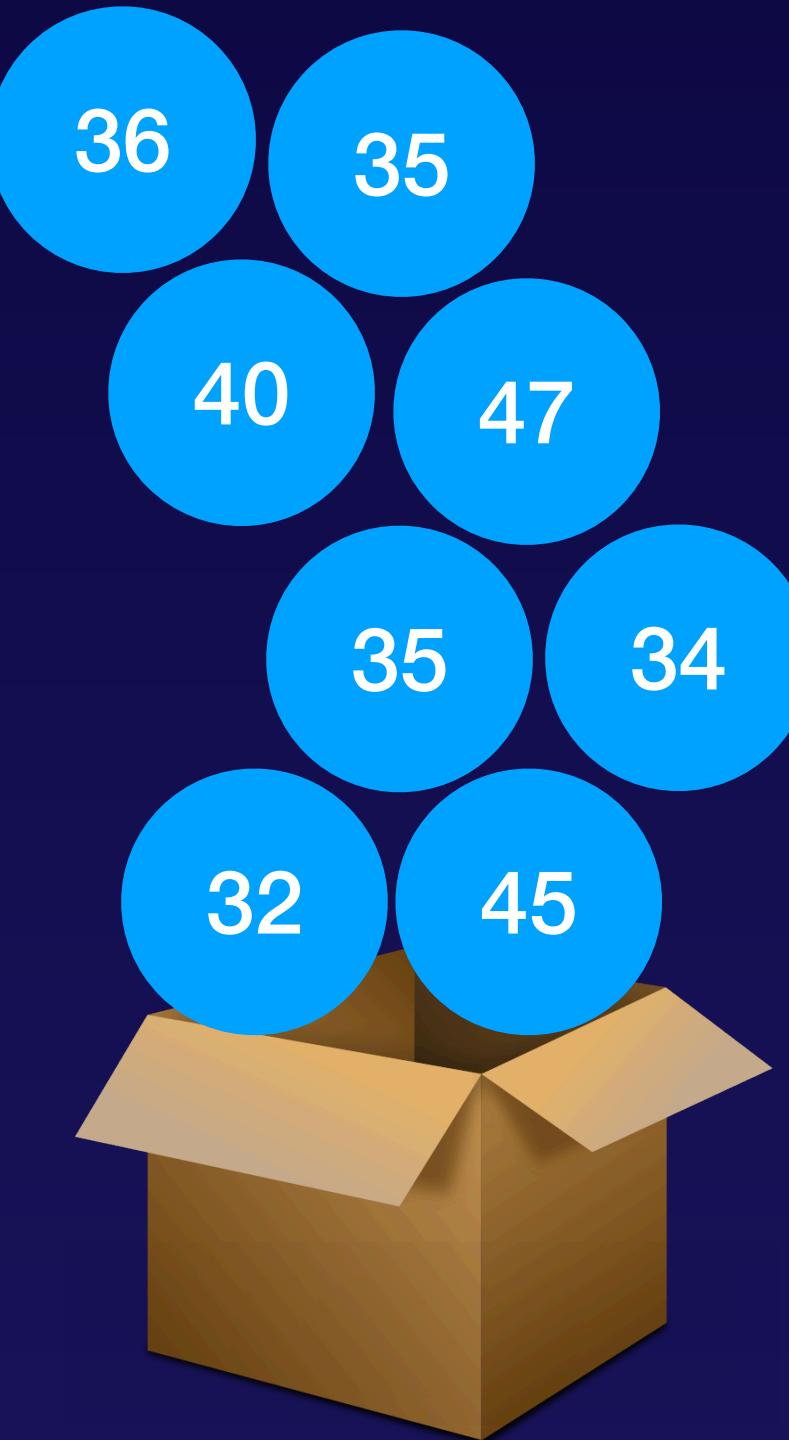
30s-50s



50 and above



30s and below



30s-50s



50 and above



30s and below



30s-50s

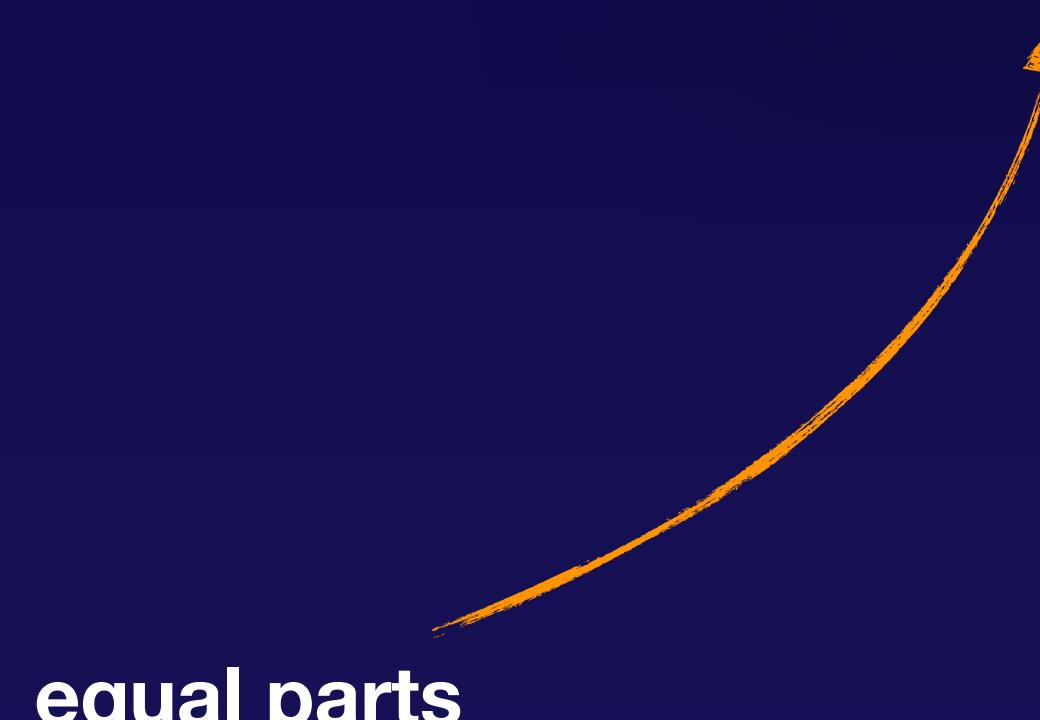


50 and above

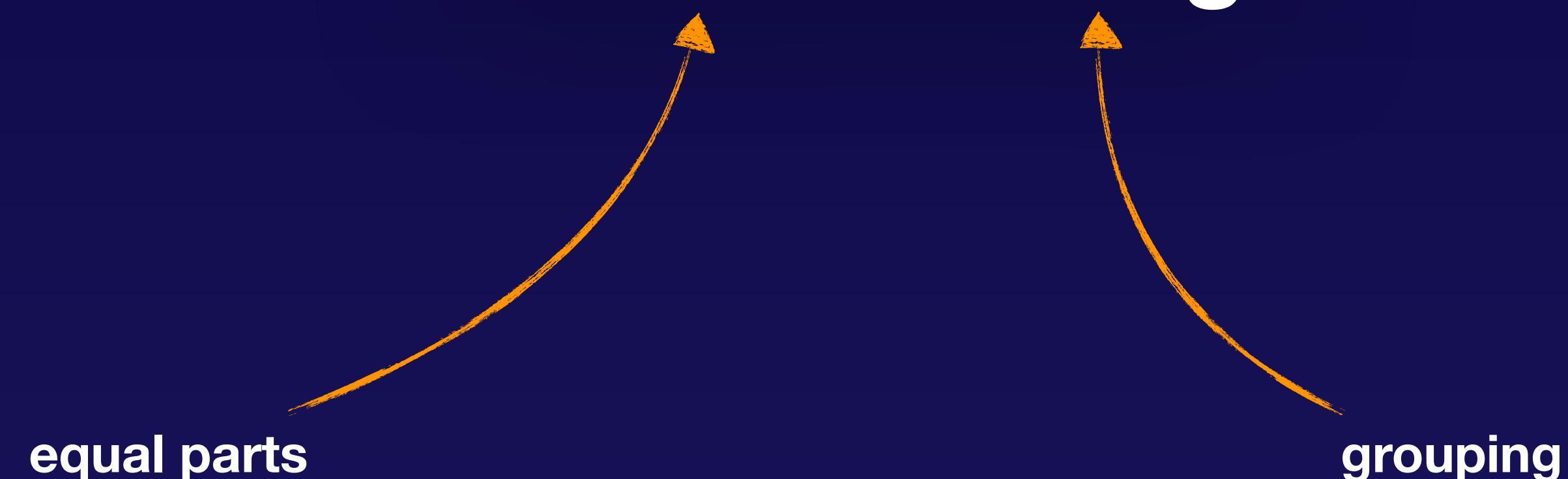
We can end up with irregular bins which are not uniform...

Quantile Binning

Quantile Binning



Quantile Binning



Quantile Binning



25 and below



26-35

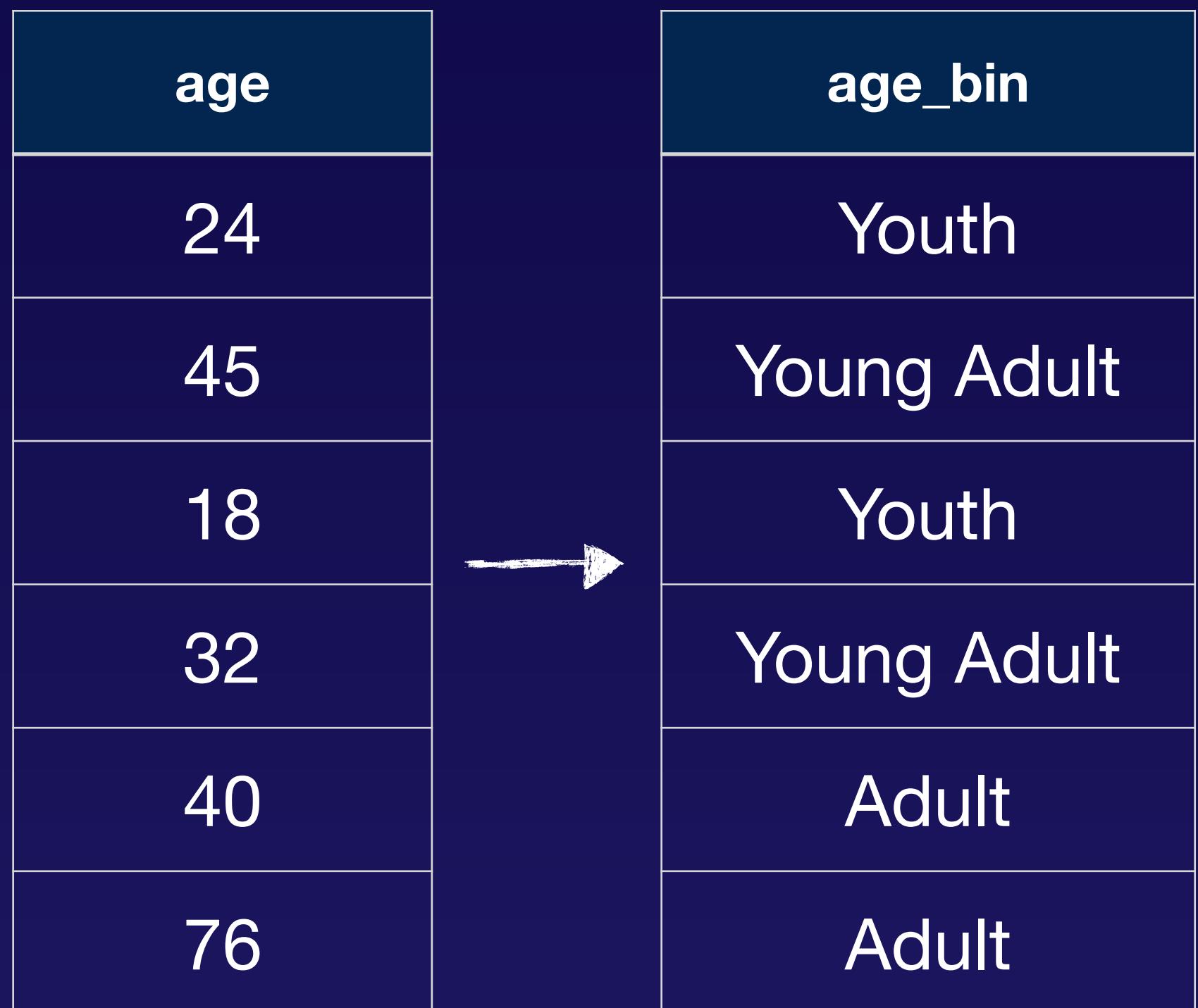


36 and above

Quantile Binning



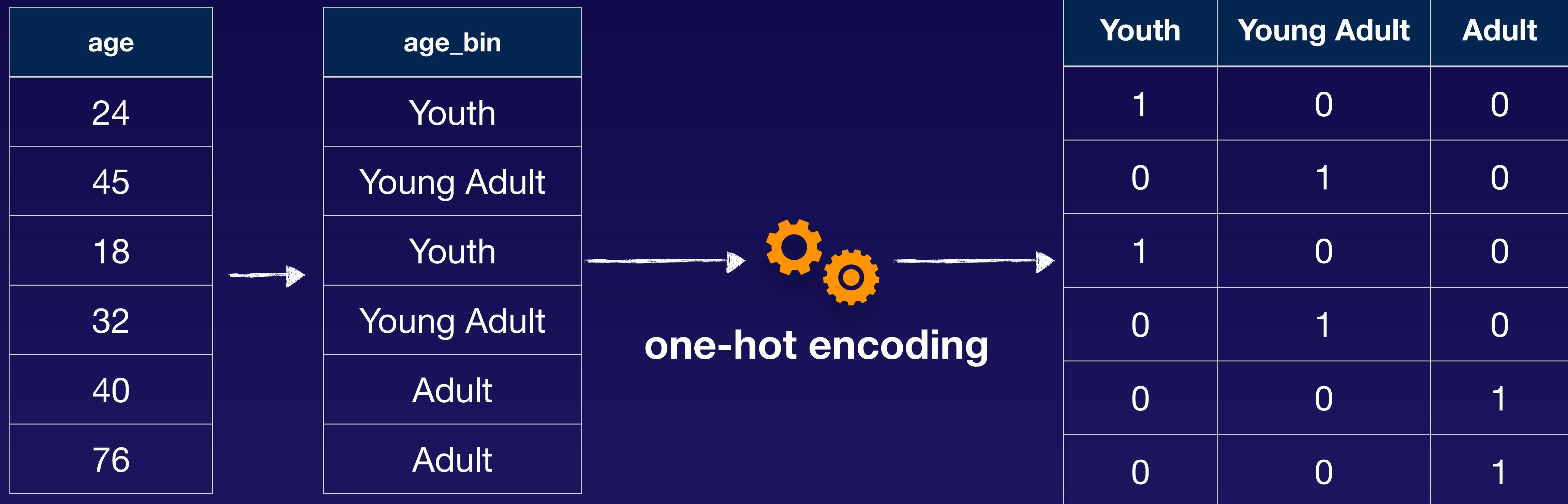
Quantile Binning



The diagram illustrates the process of quantile binning. On the left, a table labeled "age" contains seven numerical values: 24, 45, 18, 32, 40, and 76. An arrow points from this table to another table on the right, labeled "age_bin". The "age_bin" table maps these ages into categorical bins: "Youth" for ages 24, 45, and 18; "Young Adult" for ages 32 and 40; and "Adult" for the age 76.

age	age_bin
24	Youth
45	Young Adult
18	Youth
32	Young Adult
40	Adult
76	Adult

Quantile Binning



Quantile Binning Summary

1

Binning

Is used to group together values to reduce the effects of minor observation errors.

2

Quantile Binning

Bins values into equal number of bins.

3

Optimum Number of Bins

Depends on the characteristics of the variables and its relationship to the target. This is best determined through experimentation.

Numeric Feature Engineering Summary

Technique	Function
Normalization	From 0 to 1 0 - minimum value 1 - maximum value
Standardization	0 is the average value is the z-score
Quantile Binning	creates equal number of bins