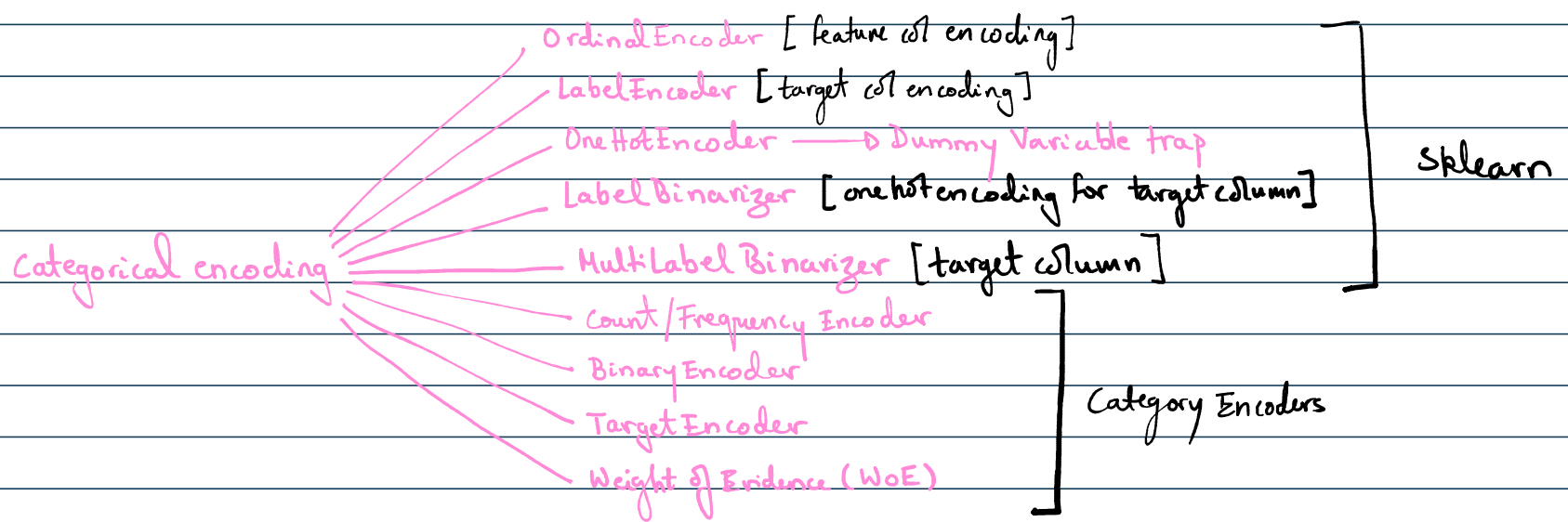


Feature Encoding

Sunday 12 May 2024 9:42 PM



Discretization → Information loss is a common problem for all binning techniques

- Custom Binning [Using Domain Knowledge]

- Uniform/equal width binning

$$\text{Bin Width} = \frac{\text{Max} - \text{Min}}{\text{\# bins}}$$

→ Tricky

KBinsDiscretizer [sklearn]

- n-bins
- encode
- strategy → uniform

→ Simple

→ Uniform coverage

→ Sensitive to outliers → Resulting in many empty bins

→ Not adaptive i.e. Doesn't study the distribution of the data
→ same algorithm regardless of underlying distribution

→ ↑ bins → noise ↓ bins → loss of info

use when

- evenly distributed data
- can act as a baseline

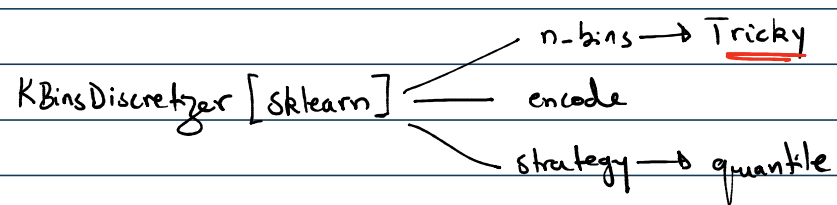
- Quantile/equal frequency binning

bins → (# bin - 1) quantiles

e.g. 4 bins → $100/4 = 25$ → 0 - 25% - 50% - 75% - 100%

e.g. 4 bins $\rightarrow 100/4 = 25 \rightarrow 0 - 25\% - 50\% - 75\% - 100\%$

\rightarrow Bin width may differ but frequency in each bin is effectively same.



\rightarrow able to handle outliers

\rightarrow Handles skewed distributions

\rightarrow Difficulty in bin interpretation

\rightarrow True info about the data distribution is lost

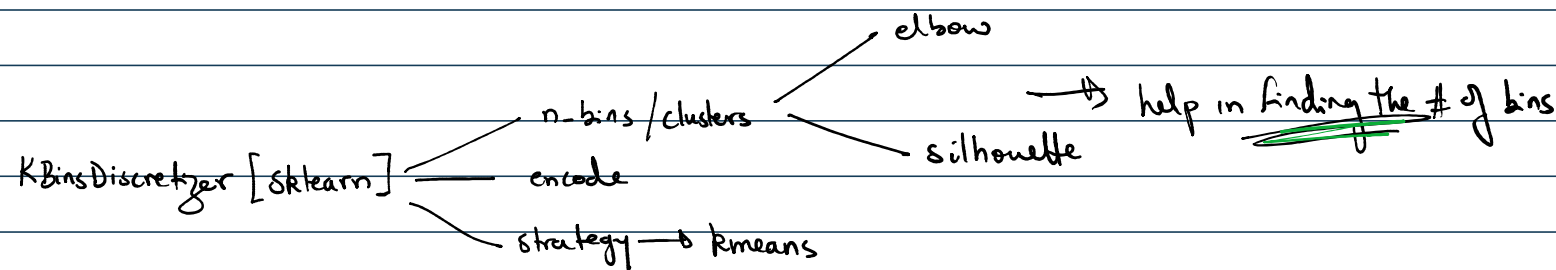
\rightarrow Computationally expensive \rightarrow sorting

K-Means Binning

\rightarrow Binning with the help of clustering

\rightarrow Works just like k-means clustering

\rightarrow Midpoints of cluster centroids act as the bin boundaries



\rightarrow Adaptive

\rightarrow Minimizes within-bin variance

\rightarrow sensitive to outliers

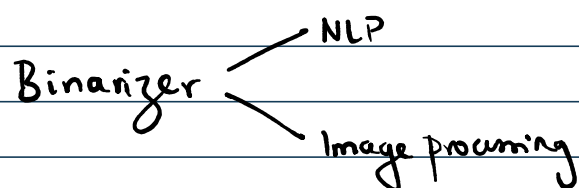
\rightarrow Computationally expensive.

\rightarrow Assumption of similar size & density of clusters

\rightarrow Interpretability

Binarization / Threshold Binning

Num \rightarrow 0 or 1 based on a threshold



\rightarrow losing a lot of info

→ losing a lot of info

→ Limited use

- Decision Tree Based Binning

→ Supervised Binning

→ feature vs target → DT → Threshold

(2) → leaf node

↳ plot tree

→ No need to spec. # of bins

→ Data leakage

Feature-engine library