

Feature Engineering



- Transform Variables
- Extract Features
- Create New Features

Missing Data Imputation Techniques

Numerical Variables



- ☐ Mean / Median Imputation
- ☐ Arbitrary value imputation
- ☐ End of tail imputation

Categorical Variables



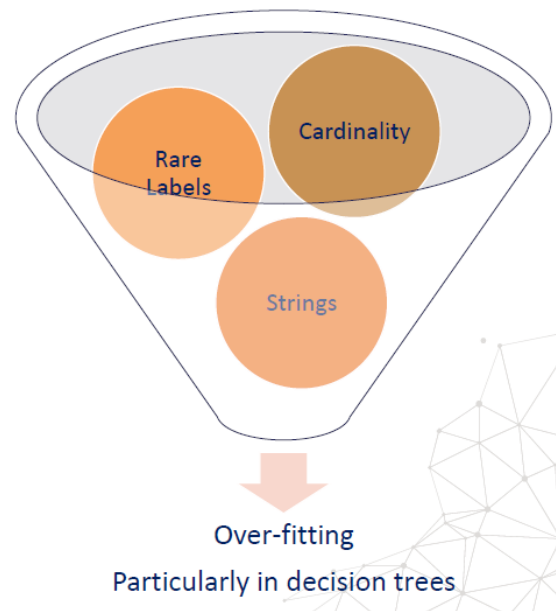
- ☐ Frequent category imputation
- ☐ Adding a "missing" category

Both

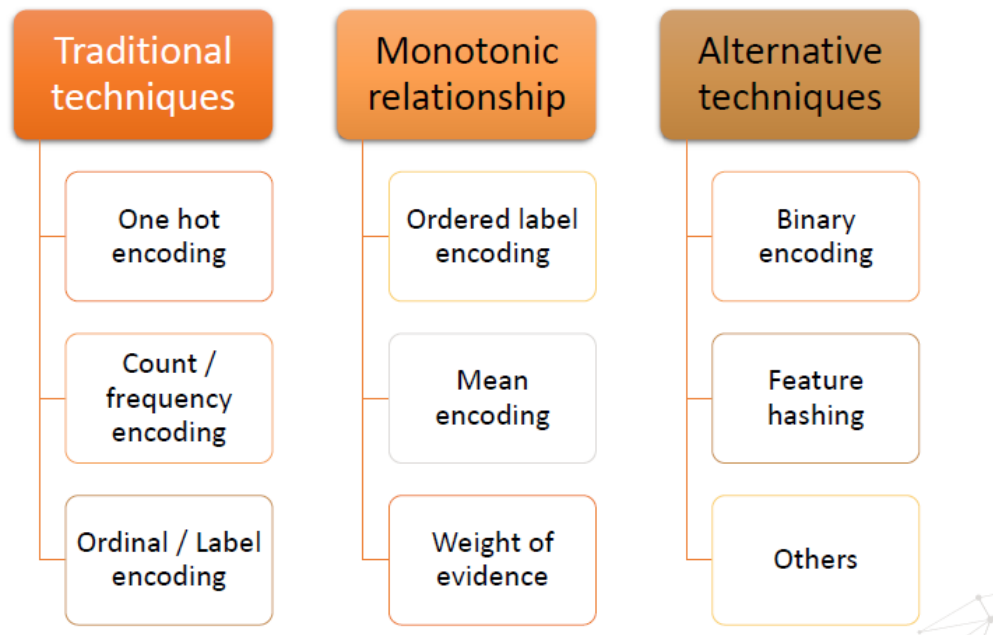


- ☐ Complete Case Analysis
- ☐ Adding a "Missing" indicator
- ☐ Random sample imputation

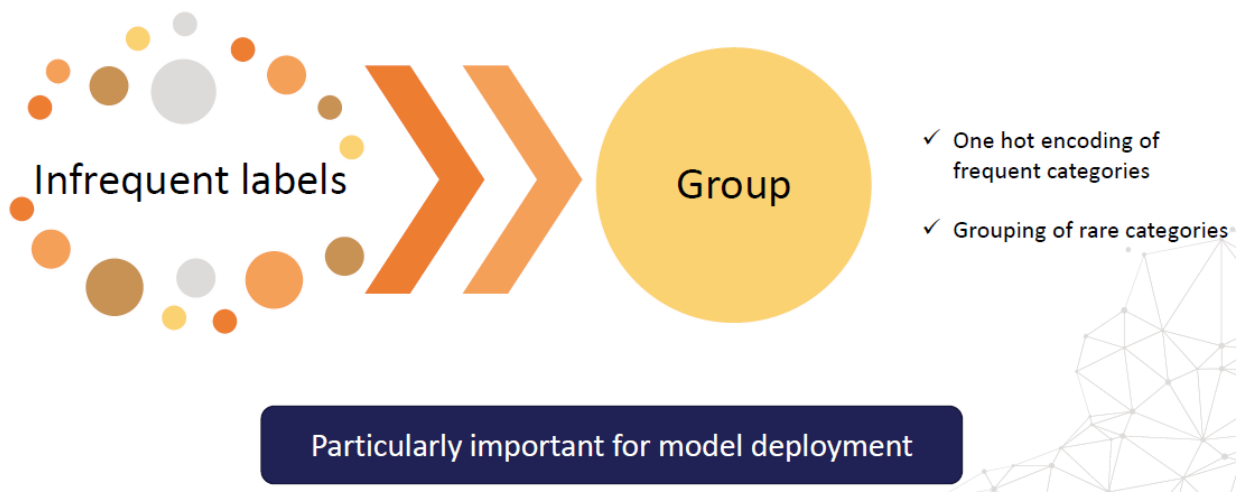
Categorical Variables



Categorical Encoding Techniques

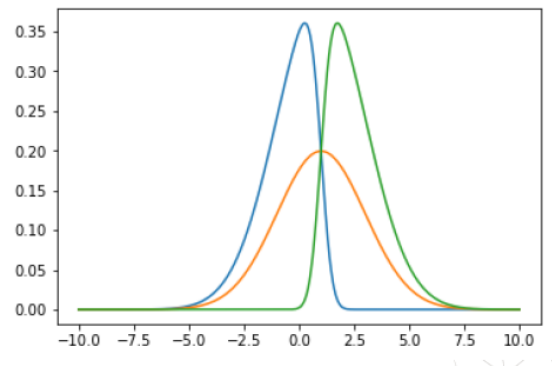


Encoding Techniques: Rare labels

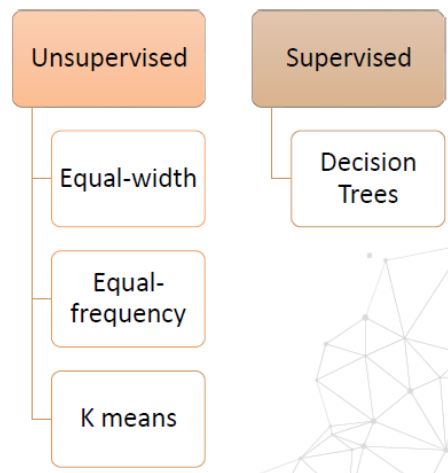
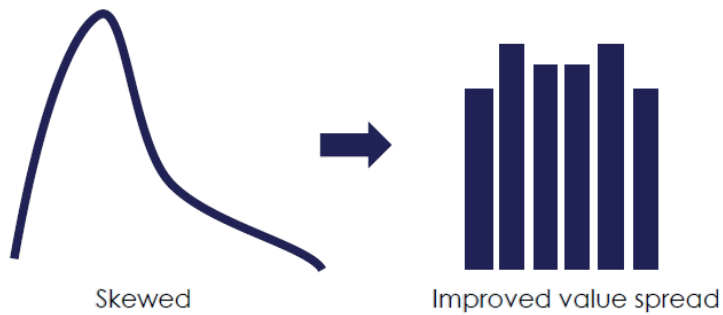


Distributions

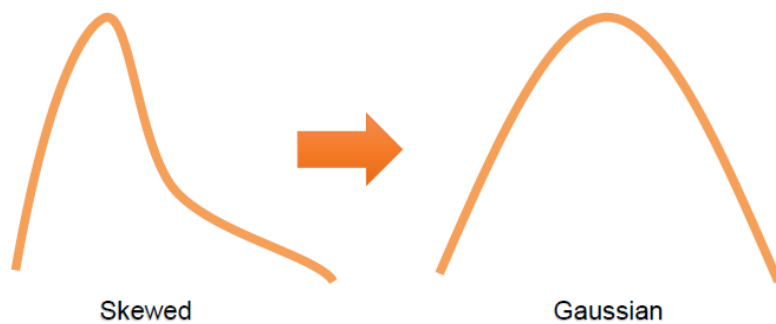
- Some models make assumptions on the variable distributions



Discretisation



Mathematical transformations

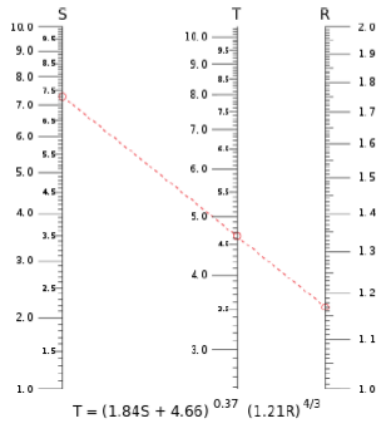


Variable transformation

- Logarithmic
- Exponential
- Reciprocal
- Box-Cox
- Yeo-Johnson



Variable Magnitude

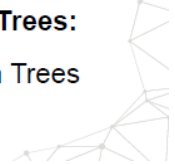


The machine learning models affected by the magnitude of the feature:

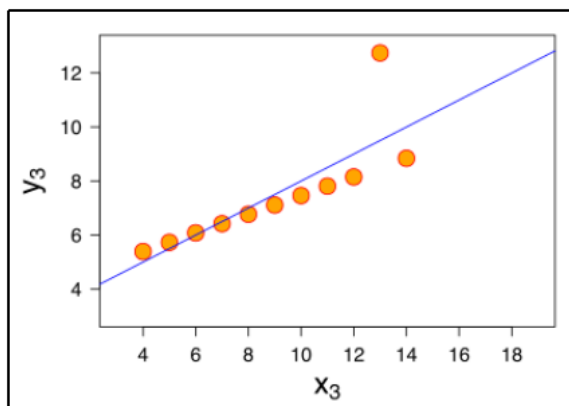
- Linear and Logistic Regression
- Neural Networks
- Support Vector Machines
- KNN
- K-means clustering
- Linear Discriminant Analysis (LDA)
- Principal Component Analysis (PCA)

Machine learning models insensitive to feature magnitude are the ones based on Trees:

- Classification and Regression Trees
- Random Forests
- Gradient Boosted Trees



Outliers



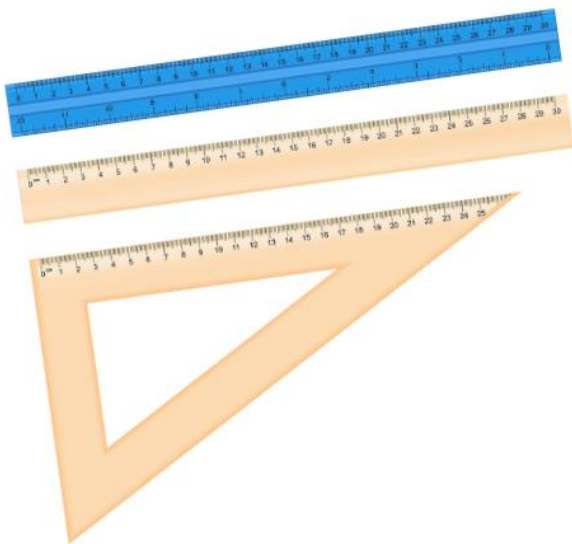
- Discretisation
- Capping / Censoring
- Truncation

Datetime Variables



- Day, Month, semester, year
- Hour, min, sec
- Elapsed Time
 - Time between transactions
 - Age

Feature scaling methods



Scaling methods

- **Standardisation**
- Mean normalisation
- **Scaling to maximum and minimum**
- Scaling to absolute maximum
- Scaling to median and quantiles
- Scaling to unit norm



Transactions and Time Series



Aggregate data

- Number of payments in last 3, 6, 12 months
- Time since last transaction
- Total spending in last month



Text

An insurance claim
A formal request to an insur
payment based on the te
Insurance claims are re
out to the insurer

- Characters, words, unique words
- Lexical diversity
- Sentences, paragraphs
- Bag of Words
- TFIDF



Geo Data



- Distances

Feature Combination



- **Ratio:** Total debt with income → Debt to income ratio
- **Sum:** Debt in different credit cards → total debt
- **Subtraction:** Income without expenses → disposable income



Open-source for Feature engineering



🏠 Category Encoders



Feature-engine



Scikit-Learn Transformers

- Missing Data Imputation
 - SimpleImputer
 - IterativeImputer
- Categorical Variable Encoding
 - OneHotEncoder
 - OrdinalEncoder
- Scalers
 - Standard Scaler
 - MinMaxScaler
 - Robust Scaler
 - A few others
- Discretisation
 - KBinsDiscretizer
- Variable Transformation
 - PowerTransformer
 - FunctionTransformer
- Variable Combination
 - Polynomial Features
- Text
 - Word Count
 - TFIDF

Feature Engine Transformers

Discretisation methods

- EqualFrequencyDiscretiser
- EqualWidthDiscretiser
- DecisionTreeDiscretiser
- ArbitraryDiscretiser

Variable Transformation methods

- LogTransformer
- ReciprocalTransformer
- PowerTransformer
- BoxCoxTransformer
- YeoJohnsonTransformer

Scikit-learn Wrapper:

- SklearnTransformerWrapper

Variable Combinations:

- MathematicalCombination
- CombineWithReferenceFeature

Imputing Methods

- MeanMedianImputer
- RandomSampleImputer
- EndTailImputer
- AddMissingIndicator
- CategoricalImputer
- ArbitraryNumberImputer
- DropMissingData

Encoding Methods

- OneHotEncoder
- OrdinalEncoder
- CountFrequencyEncoder
- MeanEncoder
- WoEEncoder
- PRatioEncoder
- RareLabelEncoder
- DecisionTreeEncoder

Outlier Handling methods

- Winsorizer
- ArbitraryOutlierCapper
- OutlierTrimmer

Category Encoders

```
import category_encoders as ce

encoder = ce.BackwardDifferenceEncoder(cols=[...])
encoder = ce.BaseNEncoder(cols=[...])
encoder = ce.BinaryEncoder(cols=[...])
encoder = ce.CatBoostEncoder(cols=[...])
encoder = ce.HashingEncoder(cols=[...])
encoder = ce.HelmertEncoder(cols=[...])
encoder = ce.JamesSteinEncoder(cols=[...])
encoder = ce.LeaveOneOutEncoder(cols=[...])
encoder = ce.MEstimateEncoder(cols=[...])
encoder = ce.OneHotEncoder(cols=[...])
encoder = ce.OrdinalEncoder(cols=[...])
encoder = ce.SumEncoder(cols=[...])
encoder = ce.PolynomialEncoder(cols=[...])
encoder = ce.TargetEncoder(cols=[...])
encoder = ce.WoEEncoder(cols=[...])
```

The screenshot shows the official documentation for Scikit-learn's Category Encoders. The page title is "Category Encoders" and it includes a search bar. A sidebar on the left lists various encoders: Backward Difference Coding, BaseN, Binary, CatBoost Encoder, Hashing, Helmert Coding, James Stein Encoder, Leave One Out, M-estimate, One Hot, Ordinal, Polynomial Coding, Sum Coding, Target Encoder, and Weight of Evidence. The main content area describes the library as a set of scikit-learn-style transformers for encoding categorical variables. It lists key properties: first-class support for pandas dataframes, explicit configuration of columns, low variance based on training set, portability, and full compatibility with sklearn pipelines. A "Usage" section provides installation instructions: `pip install category_encoders` or `conda install -c conda-forge category_encoders`.