

Data Processing and Feature Engineering:

Encoding Categorical Variables

Categorical (also called discrete),
- Attribute can take a finite set of values

Examples

color \in (red, green, blue)

is Fraud \in (true, false)

Pandas

supports special dtype = "category"

Ordinal - categories are ordered

e.g. size - small, medium, large

Nominal - categories are unordered

e.g. color

Issue: Categoricals are often represented as text but many algorithms require numericals as in plt
- we need special encoding to convert categorical into numerical representation

loan_approval ~ encode Y or N as 1 & 0

garden_size - ordinal (S, M, L) use
pandas map function

N - no garden 0, S - S, M - 10,
L - 20

use label-encoder from sklearn

wrong solution when no relationship between

If no ordering, or no relative categories
sizing, conversion / mapping will
not give us a good result.