

Essential Guide to Feature Engineering

Turn raw data into signal: encode, scale, transform, and pipeline with pandas and scikit-learn.

A practical guide to encoding, scaling, transforming, and piping data into model-ready features with pandas and scikit-learn.



rubió
Metabolomics

Ibon Martínez-Arranz | imartinez@labrubiocom

Data Science Manager at Rubió Metabolomics

www.rubiometabolomics.com

Ibon Martínez-Arranz, PhD in Mathematics and Statistics, holds MSc degrees in Applied Statistical Techniques and Mathematical Modeling. He has extensive expertise in statistical modeling and advanced data analysis techniques. Since 2017, he has led the Data Science area at Rubió Metabolomics, driving predictive model development and statistical computing management for metabolomics, data handling, and R&D projects. His doctoral research in Mathematics investigated and adapted genetic algorithms to improve the classification of NAFLD subtypes.

Itziar Mincholé Canals | iminchola@labrubio.com

Data Specialist at Rubió Metabolomics

www.rubiometabolomics.com

Itziar Mincholé Canals has a degree in physical sciences (1999) from the University of Zaragoza (Spain). In 1999 she began developing her professional career as a software analyst and developer, mainly working on web development projects and database programming in different sectors. In 2009 she joined OWL Metabolomics as a bioinformatician, where she has participated in several R&D projects and has developed several software tools for metabolomics. In 2016 she completed the master's degree in applied statistics with R software. After joining the Data Science department, she also supports data analysis.



Essential Guide to Feature Engineering

Ibon Martínez-Arranz

**Life
Feels
Good**



rubió
Metabolomics



Contents

Introduction and Setup	1
What is Feature Engineering?	1
Why is Feature Engineering Important?	1
Libraries Used in This Book	2
Example Datasets	3
Quick Overview of the Datasets	4
Iris Dataset	4
California Housing Dataset	6
Titanic Dataset	9
Initial Data Distributions	13
What Comes Next	17
 Encoding Categorical Variables	 19
Categorical vs Numerical Features	19
Takeaways	24
Label Encoding (with <code>sklearn.preprocessing.LabelEncoder</code>)	25
One-Hot Encoding (pandas & scikit-learn)	26
Ordinal Encoding (with <code>sklearn.preprocessing.OrdinalEncoder</code>)	31
High-Cardinality Categories	34
Practical Pitfalls: Dummy Variable Trap & Sparse Matrices	36
Dummy Variable Trap	36
Sparse Matrices	37
Putting It Together: A Mini Encoding Pipeline	38
Summary	43

Scaling Numerical Features	45
Why scaling is necessary	45
Standardization with <code>StandardScaler</code>	45
Min-Max Scaling with <code>MinMaxScaler</code>	50
Robust Scaling with <code>RobustScaler</code>	52
Comparing scaling methods on the same dataset	55
Practical tips: essential vs optional	58
What comes next	59
Transforming Variables	61
Detecting skewness and choosing a transformation	61
Log transformations for skewed data	63
Box-Cox and Yeo-Johnson with <code>PowerTransformer</code>	67
Quantile transformations for non-Gaussian distributions	69
Binarization with <code>Binarizer</code>	71
Combining transformations with pipelines	74
Summary	77
Feature Generation	79
Creating interaction features (<code>PolynomialFeatures</code>)	79
Feature crosses: combining categorical features	84
Extracting date/time features	87
Text preprocessing basics: Bag-of-Words and TF-IDF	93
Domain-driven features (California Housing)	96
Summary	104
Handling Missing Data	107
Why missing data matters (MCAR, MAR, MNAR)	107
Inspecting missingness	107
Simple imputation (mean/median/most-frequent/constant)	110
Adding missingness indicators	114
<code>KNNImputer</code> (neighbors-based imputation)	117

IterativeImputer (multivariate, model-based) 119

Visual check: distribution before/after imputation (adds value) 121

Leakage-safe evaluation: quick comparison 122

Summary 124



Introduction and Setup

What is Feature Engineering?

Feature engineering is the art and science of transforming raw data into meaningful variables (features) that a machine learning model can effectively use.

A typical dataset is rarely ready to be fed directly into an algorithm.

Categorical variables, missing values, different scales, or skewed distributions can all reduce model performance or even cause failures.

Feature engineering involves tasks such as: - Converting categories into numbers (encoding). - Scaling and normalizing numerical variables. - Applying transformations to improve distributions. - Creating new features from existing ones. - Handling missing values consistently.

In this book, we will go through these steps using practical examples in **Python**, mainly with **pandas** and **scikit-learn**.

Why is Feature Engineering Important?

Machine learning algorithms, no matter how sophisticated, cannot fully compensate for poorly prepared data.

- A strong model with bad features → poor results.

- A simple model with well-crafted features → surprisingly good results.

In practice, most machine learning projects spend more than **70% of the time on preprocessing and feature engineering**, not on model tuning.

This document will cover: - Encoding of categorical variables. - Scaling and normalization of numerical variables. - Variable transformations (log, Box-Cox, Yeo-Johnson). - Feature creation. - Missing value imputation. - Pipelines for reproducibility.

Libraries Used in This Book

We will rely on the following libraries:

- **pandas** → handling tabular data (load, clean, transform).
- **numpy** → efficient numerical operations.
- **scikit-learn** → preprocessing tools (`OneHotEncoder`, `StandardScaler`, `PowerTransformer`, `Pipeline`, etc.).
- **matplotlib** and **seaborn** → visualization of distributions and results.

Let's import them and check that everything is working.

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 from sklearn import datasets
```

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import seaborn as sns
5
6 from sklearn import datasets
```

Example Datasets

We will use some classical machine learning datasets, easily available from `scikit-learn` or `seaborn`:

- **Titanic dataset** (survival of passengers).
Useful to work with categorical variables, missing values, and encoding.
- **Iris dataset** (classification of flowers based on measurements).
Perfect for scaling, transformations, and simple pipelines.
- **California Housing dataset** (prediction of house prices).
Good for numeric features, scaling, and feature generation.

Let's load them.

```
1 # Iris dataset
2 iris = datasets.load_iris(as_frame=True)
3 iris_df = iris.frame
4
5 # California housing dataset
6 housing = datasets.fetch_california_housing(as_frame=True
7 )
8 housing_df = housing.frame
9
10 # Titanic dataset (from seaborn)
11 titanic_df = sns.load_dataset("titanic")
12
13 print("Iris shape:", iris_df.shape)
14 print("Housing shape:", housing_df.shape)
15 print("Titanic shape:", titanic_df.shape)
```

```
1 # Iris dataset
2 iris = datasets.load_iris(as_frame=True)
3 iris_df = iris.frame
4
5 # California housing dataset
6 housing = datasets.fetch_california_housing(as_frame=True
7 )
```

```
7 housing_df = housing.frame
8
9 # Titanic dataset (from seaborn)
10 titanic_df = sns.load_dataset("titanic")
11
12 print("Iris shape:", iris_df.shape)
13 print("Housing shape:", housing_df.shape)
14 print("Titanic shape:", titanic_df.shape)
```

```
1 Iris shape: (150, 5)
2 Housing shape: (20640, 9)
3 Titanic shape: (891, 15)
```

Quick Overview of the Datasets

Iris Dataset

- **150 samples**, 4 numerical features (sepal length/width, petal length/width).
- Target: species of iris flower (Setosa, Versicolor, Virginica).
- Good for classification tasks.

```
1 # Quick peek at each dataset
2 display(iris_df.head())
```

```
1 # Quick peek at each dataset
2 display(iris_df.head())
```

sepal length (cm)

sepal width (cm)

petal length (cm)

petal width (cm)

target

0

5.1

3.5

1.4

0.2

0

1

4.9

3.0

1.4

0.2

0

2

4.7

3.2

1.3

0.2

0

3

4.6

3.1

1.5

0.2

0
4
5.0
3.6
1.4
0.2
0

California Housing Dataset

- **20,640 samples**, 8 numerical features (median income, average rooms, latitude, longitude...).
- Target: median house value.
- Good for regression tasks.

```
1 # Quick peek at each dataset
2 display(housing_df.head())
```

```
1 # Quick peek at each dataset
2 display(housing_df.head())
```

MedInc

HouseAge

AveRooms

AveBedrms

Population

AveOccup

Latitude

Longitude

MedHouseVal

0

8.3252

41.0

6.984127

1.023810

322.0

2.555556

37.88

-122.23

4.526

1

8.3014

21.0

6.238137

0.971880

2401.0

2.109842

37.86

-122.22

3.585

2

7.2574

52.0

8.288136

1.073446

496.0

2.802260

37.85

-122.24

3.521

3

5.6431

52.0

5.817352

1.073059

558.0

2.547945

37.85

-122.25

3.413

4

3.8462

52.0

6.281853

1.081081

565.0

2.181467

37.85

-122.25

3.422

Titanic Dataset

- **891 passengers**, mix of categorical and numerical features (sex, age, class, embarked...).
- Target: survival (0 = died, 1 = survived).
- Good for preprocessing demonstrations (encoding, imputation).

```
1 # Quick peek at each dataset
2 display(titanic_df.head())
```

```
1 # Quick peek at each dataset
2 display(titanic_df.head())
```

survived

pclass

sex

age

sibsp

parch

fare

embarked

class
who
adult_male
deck
embark_town
alive
alone
0
0
3
male
22.0
1
0
7.2500
S
Third
man
True
NaN
Southampton
no
False

1

1

1

female

38.0

1

0

71.2833

C

First

woman

False

C

Cherbourg

yes

False

2

1

3

female

26.0

0

0

7.9250

S

Third

woman

False

NaN

Southampton

yes

True

3

1

1

female

35.0

1

0

53.1000

S

First

woman

False

C

Southampton

yes

False

4

0

3

male

35.0

0

0

8.0500

S

Third

man

True

NaN

Southampton

no

True

Initial Data Distributions

Before diving into encoding, scaling, and transformations, it is useful to visualize the **distributions of the features**.

This helps us see: - Which variables are numeric vs categorical. - If variables have different scales. - If distributions are skewed. - How datasets differ in complexity.

Let's plot some histograms for each dataset.

```
1 # Histograms for Iris dataset
2 iris_df.hist(figsize=(3*3,2.5*2), bins = 15, layout =
    (2,3), color = "#FECB00", edgecolor = '0.3',)
3 plt.suptitle("Iris Dataset - Feature Distributions",
    fontsize=14)
4 plt.show()
```

```
1 # Histograms for Iris dataset
2 iris_df.hist(figsize=(3*3,2.5*2), bins = 15, layout =
    (2,3), color = "#FECB00", edgecolor = '0.3',)
3 plt.suptitle("Iris Dataset - Feature Distributions",
    fontsize=14)
4 plt.show()
```

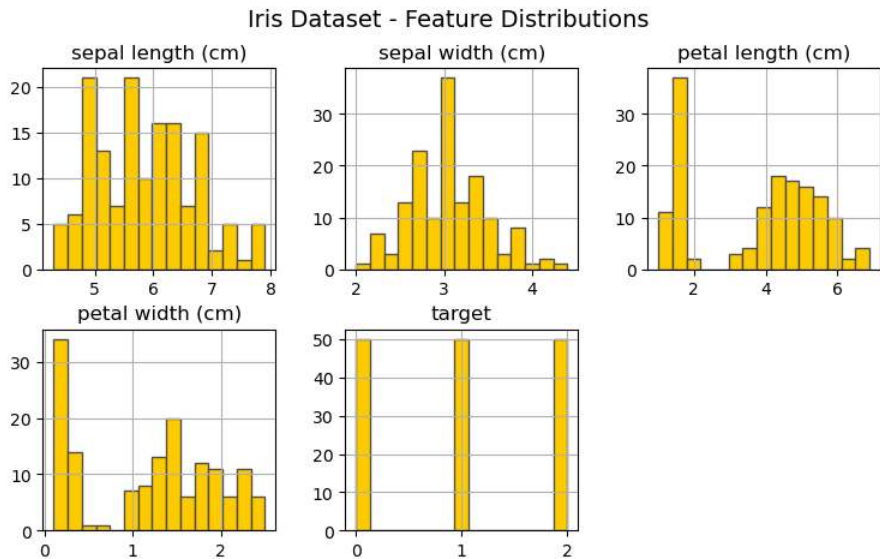


Figure 1: Image generated by the provided code.

```
1 # Histograms for California Housing dataset
2 housing_df.hist(figsize=(3*3,2.5*3), bins = 15, layout =
   (3,3), color = "#FECB00", edgecolor = '0.3',)
3 plt.suptitle("California Housing Dataset - Feature
   Distributions", fontsize=14)
4 plt.show()
```

```
1 # Histograms for California Housing dataset
2 housing_df.hist(figsize=(3*3,2.5*3), bins = 15, layout =
   (3,3), color = "#FECB00", edgecolor = '0.3',)
3 plt.suptitle("California Housing Dataset - Feature
   Distributions", fontsize=14)
4 plt.show()
```


California Housing Dataset - Feature Distributions

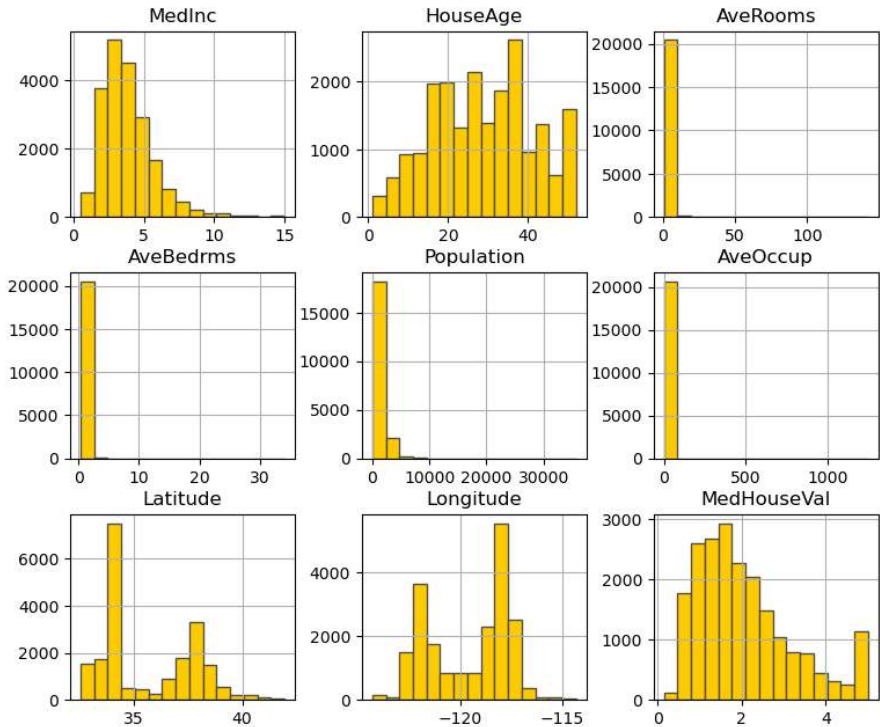


Figure 2: Image generated by the provided code.

```

1 # Histograms for Titanic dataset (numeric features only)
2 titanic_df.select_dtypes(include=["float64", "int64"]).
   hist(figsize=(3*3,2.5*2), bins = 15, layout = (2,3),
   color = "#FECB00", edgecolor = '0.3',)
3 plt.suptitle("Titanic Dataset - Numeric Feature
   Distributions", fontsize=14)
4 plt.show()

```

```

1 # Histograms for Titanic dataset (numeric features only)

```

```
2 titanic_df.select_dtypes(include=["float64", "int64"]).  
   hist(figsize=(3*3,2.5*2), bins = 15, layout = (2,3),  
       color = "#FECB00", edgecolor = '0.3',)  
3 plt.suptitle("Titanic Dataset - Numeric Feature  
   Distributions", fontsize=14)  
4 plt.show()
```

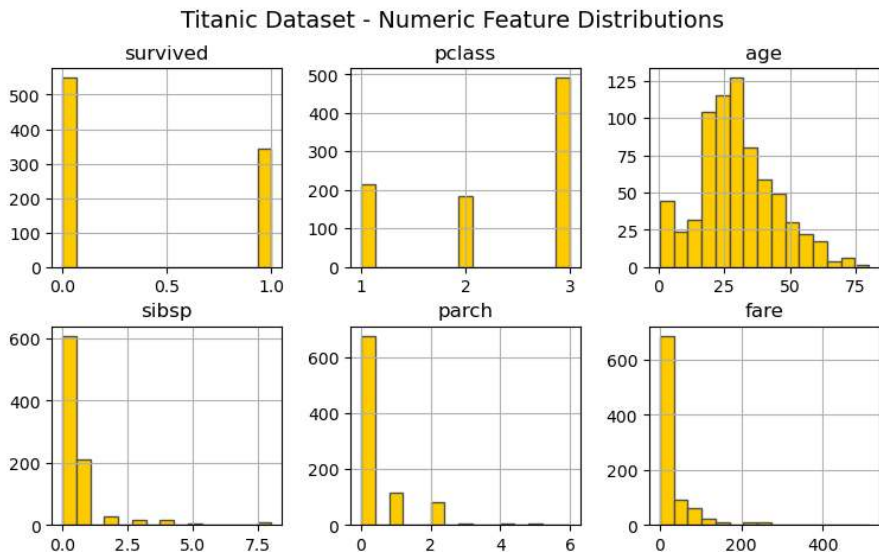


Figure 3: Image generated by the provided code.

What Comes Next

Now that the environment is set up and we have a first look at the datasets, in the next chapters we will learn to:

1. Encode categorical variables (label encoding, one-hot, ordinal).
2. Scale and normalize numerical variables.

3. Transform distributions with log, Box-Cox, and Yeo-Johnson.
4. Create new features and handle missing values.
5. Build complete preprocessing pipelines for reproducibility and production.

With this, we are ready to move to **Chapter 2: Encoding Categorical Variables**.

Encoding Categorical Variables

Categorical vs Numerical Features

Numerical features are measured on a numeric scale (e.g., `age`, `fare`, `sepal_length`).

Categorical features represent discrete groups or labels (e.g., `sex`, `class`, `embarked`).

Some categorical features are **ordinal** (they have an inherent order, e.g., `low` < `medium` < `high`), while others are **nominal** (no order, e.g., `red`, `green`, `blue`).

Why this matters:

- Many Machine Learning algorithms require numeric input. Categorical variables must be encoded.
- The choice of encoding depends on whether the category is **ordinal** or **nominal**.
- For high-cardinality categories (hundreds/thousands of unique values), naive one-hot encoding can explode dimensionality.

Below we'll inspect the Titanic dataset to identify categorical vs numerical columns and look at their distributions.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
```

```
5 titanic_df = sns.load_dataset("titanic").copy()
6
7 # Basic dtype view
8 display(titanic_df.dtypes.to_frame("dtype"))
9
10 # Split by type
11 cat_cols = titanic_df.select_dtypes(include=["object", "
    category", "bool"]).columns.tolist()
12 num_cols = titanic_df.select_dtypes(include=["number"]).
    columns.tolist()
13
14 print("Categorical columns:", cat_cols)
15 print("Numerical columns:", num_cols)
16
17 # Quick cardinality check
18 cardinality = titanic_df[cat_cols].nunique().sort_values(
    ascending=False)
19 display(cardinality.to_frame("n_unique"))
```

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 titanic_df = sns.load_dataset("titanic").copy()
6
7 # Basic dtype view
8 display(titanic_df.dtypes.to_frame("dtype"))
9
10 # Split by type
11 cat_cols = titanic_df.select_dtypes(include=["object", "
    category", "bool"]).columns.tolist()
12 num_cols = titanic_df.select_dtypes(include=["number"]).
    columns.tolist()
13
14 print("Categorical columns:", cat_cols)
15 print("Numerical columns:", num_cols)
16
17 # Quick cardinality check
18 cardinality = titanic_df[cat_cols].nunique().sort_values(
    ascending=False)
19 display(cardinality.to_frame("n_unique"))
```

dtype

survived

int64

pclass

int64

sex

object

age

float64

sibsp

int64

parch

int64

fare

float64

embarked

object

class

category

who

object

adult_male

bool

deck

category

embark_town

object

alive

object

alone

bool

```
1 Categorical columns: ['sex', 'embarked', 'class', 'who',  
    'adult_male', 'deck', 'embark_town', 'alive', 'alone']  
2 Numerical columns: ['survived', 'pclass', 'age', 'sibsp',  
    'parch', 'fare']
```

n_unique

deck

7

embark_town

3

embarked

3

who

3

class

3

sex

2

adult_male

2

alive

2

alone

2

```
1 # Plot distributions for some common categorical
  variables
2 fig, axes = plt.subplots(1, 3, figsize=(12, 4))
3 sns.countplot(data=titanic_df, x="sex", ax=axes[0],
  palette = ['#fecb00ff', '#6b99a6ff'])
4 axes[0].set_title("sex")
5
6 sns.countplot(data=titanic_df, x="class", ax=axes[1],
  palette = ['#b7d3dbff', '#6b99a6ff', '#3a6674ff'])
7 axes[1].set_title("class")
8
9 sns.countplot(data=titanic_df, x="embarked", ax=axes[2],
  palette = ['#ffef80ff', '#fecb00ff', '#cc9b00ff'])
10 axes[2].set_title("embarked")
11 plt.tight_layout()
12 plt.show()
```

```
1 # Plot distributions for some common categorical
  variables
2 fig, axes = plt.subplots(1, 3, figsize=(12, 4))
3 sns.countplot(data=titanic_df, x="sex", ax=axes[0],
  palette = ['#fecb00ff', '#6b99a6ff'])
4 axes[0].set_title("sex")
5
```



```
6 sns.countplot(data=titanic_df, x="class", ax=axes[1],  
    palette = ['#b7d3dbff', '#6b99a6ff', '#3a6674ff'])  
7 axes[1].set_title("class")  
8  
9 sns.countplot(data=titanic_df, x="embarked", ax=axes[2],  
    palette = ['#ffef80ff', '#fecb00ff', '#cc9b00ff'])  
10 axes[2].set_title("embarked")  
11 plt.tight_layout()  
12 plt.show()
```

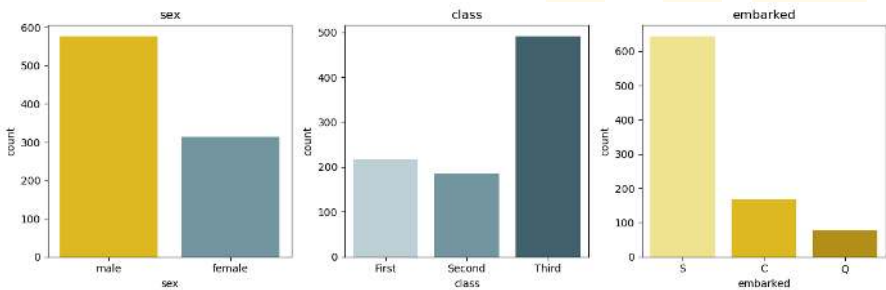


Figure 1: Image generated by the provided code.

Takeaways

- `sex`, `class`, and `embarked` are typical **nominal**/mixed categorical features.
- `class` is **ordinal** on Titanic (`Third` < `Second` < `First`).
- High-cardinality columns can cause problems for naive one-hot encoding.

Label Encoding (with `sklearn.preprocessing.LabelEncoder`)

Label Encoding converts labels to integers `{class_1: 0, class_2: 1, ...}`.

Best use: encoding the **target variable (y)** in classification problems, or truly **ordinal** features if you understand the implied ordering.

Pitfall: Using `LabelEncoder` on **nominal feature columns** injects a fake order (e.g., `male=0`, `female=1`), which can mislead linear and distance-based models.

We'll demonstrate it on the **Iris** dataset's target (species names).

```
1 from sklearn.preprocessing import LabelEncoder
2
3 iris_df = sns.load_dataset("iris").copy() # species is
   string labels
4 y_str = iris_df["species"]
5
6 le = LabelEncoder()
7 y_encoded = le.fit_transform(y_str)
8
9 print("Classes:", list(le.classes_))
10 print("Encoded sample:", y_str.head().tolist(), "->",
      y_encoded[:5].tolist())
```

```
1 from sklearn.preprocessing import LabelEncoder
2
3 iris_df = sns.load_dataset("iris").copy() # species is
   string labels
4 y_str = iris_df["species"]
5
6 le = LabelEncoder()
7 y_encoded = le.fit_transform(y_str)
8
9 print("Classes:", list(le.classes_))
10 print("Encoded sample:", y_str.head().tolist(), "->",
      y_encoded[:5].tolist())
```

```
1 Classes: ['setosa', 'versicolor', 'virginica']
2 Encoded sample: ['setosa', 'setosa', 'setosa', 'setosa',
  'setosa'] -> [0, 0, 0, 0, 0]
```

When not to use it: Do **not** use `LabelEncoder` for nominal features like `sex` or `embarked` unless your algorithm is tree-based (trees are insensitive to monotonic transformations). For general pipelines, prefer **One-Hot** or **OrdinalEncoder** (with explicit order).

One-Hot Encoding (pandas & scikit-learn)

One-Hot Encoding (OHE) creates one binary column per category (e.g., `sex_female`, `sex_male`).

- Pros: safe for nominal data, widely used.
- Cons: dimensionality blow-up for high-cardinality features.

We'll show:

- 1) `pandas.get_dummies()` — fast and convenient for EDA/quick baselines.
- 2) `sklearn.preprocessing.OneHotEncoder` — integrates with pipelines/transformers (recommended for production).

```
1 # One-Hot with pandas.get_dummies()
2 titanic_small = titanic_df[["survived", "sex", "class", "
   embarked", "age", "fare"]].copy()
3
4 # drop_first=True can avoid the "dummy variable trap" for
   linear models with intercept
5 dummies_pd = pd.get_dummies(
6     titanic_small,
7     columns=["sex", "class", "embarked"],
```

```
8     drop_first=True # set to False if you want the full
      set of dummies
9 )
10 display(dummies_pd.head())
```

```
1 # One-Hot with pandas.get_dummies()
2 titanic_small = titanic_df[["survived", "sex", "class", "
   embarked", "age", "fare"]].copy()
3
4 # drop_first=True can avoid the "dummy variable trap" for
   linear models with intercept
5 dummies_pd = pd.get_dummies(
6     titanic_small,
7     columns=["sex", "class", "embarked"],
8     drop_first=True # set to False if you want the full
      set of dummies
9 )
10 display(dummies_pd.head())
```

survived

age

fare

sex_male

class_Second

class_Third

embarked_Q

embarked_S

0

0

22.0

7.2500

1

0

1

0

1

1

1

38.0

71.2833

0

0

0

0

0

2

1

26.0

7.9250

0

0

1

0

1

3

1

35.0

53.1000

0

0

0

0

1

4

0

35.0

8.0500

1

0

1

0

1

`drop_first=True` removes one level per categorical feature to avoid perfect multicollinearity in linear models with intercept.

For tree-based or regularized models, keeping all dummies is usually fine.

```
1 # One-Hot with scikit-learn's OneHotEncoder +  
  ColumnTransformer  
2 from sklearn.compose import ColumnTransformer
```

```

3  from sklearn.preprocessing import OneHotEncoder
4
5  X = titanic_df.drop(columns=["survived"])
6  y = titanic_df["survived"]
7
8  categorical_features = ["sex", "class", "embarked"]
9  numeric_features = ["age", "fare"]
10
11 # Note: For scikit-learn >=1.2 use sparse_output=False;
    for <1.2 use sparse=False
12 ohe = OneHotEncoder(handle_unknown="ignore",
    sparse_output=False, drop=None)
13
14 ct = ColumnTransformer(
15     transformers=[
16         ("ohe", ohe, categorical_features),
17         ("passthrough", "passthrough", numeric_features),
18     ],
19     remainder="drop"
20 )
21
22 X_enc = ct.fit_transform(titanic_df)
23 feature_names = ct.get_feature_names_out()
24 print("Encoded shape:", X_enc.shape)
25 print("Sample feature names:", feature_names[:10])

```

```

1  # One-Hot with scikit-learn's OneHotEncoder +
    ColumnTransformer
2  from sklearn.compose import ColumnTransformer
3  from sklearn.preprocessing import OneHotEncoder
4
5  X = titanic_df.drop(columns=["survived"])
6  y = titanic_df["survived"]
7
8  categorical_features = ["sex", "class", "embarked"]
9  numeric_features = ["age", "fare"]
10
11 # Note: For scikit-learn >=1.2 use sparse_output=False;
    for <1.2 use sparse=False
12 ohe = OneHotEncoder(handle_unknown="ignore",
    sparse_output=False, drop=None)

```

```
13
14 ct = ColumnTransformer(
15     transformers=[
16         ("ohe", ohe, categorical_features),
17         ("passthrough", "passthrough", numeric_features),
18     ],
19     remainder="drop"
20 )
21
22 X_enc = ct.fit_transform(titanic_df)
23 feature_names = ct.get_feature_names_out()
24 print("Encoded shape:", X_enc.shape)
25 print("Sample feature names:", feature_names[:10])
```

```
1 Encoded shape: (891, 11)
2 Sample feature names: ['ohe__sex_female' 'ohe__sex_male'
3   'ohe__class_First' 'ohe__class_Second'
4   'ohe__class_Third' 'ohe__embarked_C' 'ohe__embarked_Q' '
   ohe__embarked_S'
   'ohe__embarked_nan' 'passthrough__age']
```

Notes

- `handle_unknown="ignore"` prevents errors when unseen categories appear at inference.
- Use `drop="if_binary"` (or `drop="first"`) to reduce redundancy for linear models.
- Prefer integrating OHE via `ColumnTransformer` and `Pipeline` so training and inference apply the exact same steps.

Ordinal Encoding (with `sklearn.preprocessing.OrdinalEncoder`)

Ordinal Encoding maps categories to ordered integers **based on a known order**. Use it **only** when the order is meaningful (e.g., `Third < Second < First`).

class in Titanic).

For nominal variables, do **not** use ordinal encoding.

We will encode **class** with a custom order.

```

1  from sklearn.preprocessing import OrdinalEncoder
2  from numpy import nan
3
4  ord_cols = ["class"]
5  ord_order = ["Third", "Second", "First"]
6
7  ord_enc = OrdinalEncoder(
8      categories=ord_order,
9      handle_unknown="use_encoded_value",
10     unknown_value=np.nan
11 )
12
13 # Copy to avoid modifying original
14 titanic_ord = titanic_df[["class"]].copy()
15 titanic_ord_enc = ord_enc.fit_transform(titanic_ord)
16
17 titanic_df["class_ordinal"] = titanic_ord_enc
18 display(titanic_df[["class", "class_ordinal"]].head(10))

```

- 1 ****Tip:**** For multiple ordinal features, pass a list of lists to `categories`, one per column.
- 2 For missing or unseen categories at inference, consider imputation or an explicit "Unknown" category in preprocessing.

```

1  from sklearn.preprocessing import OrdinalEncoder
2  from numpy import nan
3
4  ord_cols = ["class"]
5  ord_order = ["Third", "Second", "First"]
6
7  ord_enc = OrdinalEncoder(
8      categories=ord_order,
9      handle_unknown="use_encoded_value",
10     unknown_value=np.nan

```

```
11 )
12
13 # Copy to avoid modifying original
14 titanic_ord = titanic_df[["class"]].copy()
15 titanic_ord_enc = ord_enc.fit_transform(titanic_ord)
16
17 titanic_df["class_ordinal"] = titanic_ord_enc
18 display(titanic_df[["class", "class_ordinal"]].head(10))
```

class

class_ordinal

0

Third

0.0

1

First

2.0

2

Third

0.0

3

First

2.0

4

Third

0.0

5

Third

0.0

6

First

2.0

7

Third

0.0

8

Third

0.0

9

Second

1.0

Tip: For multiple ordinal features, pass a list of lists to `categories`, one per column. For missing or unseen categories at inference, consider imputation or an explicit “Unknown” category in preprocessing.

High-Cardinality Categories

The seaborn `titanic` dataset does not include truly high-cardinality categorical features (like `ticket` or `cabin` in the Kaggle version). To demonstrate best practices:

- We will illustrate **frequency encoding** on a low-cardinality feature (`embark_town`) to show the mechanics.

- We will also create a **synthetic high-cardinality** feature (`synthetic_id`) with a long-tail distribution, to show why frequency encoding and the **hash-ing trick** are useful when one-hot encoding is impractical.

```
1 # Load seaborn titanic
2 titanic_df = sns.load_dataset("titanic").copy()
3
4 ## Create a synthetic high-cardinality feature (e.g.,
   ~600 unique IDs, Zipf-like long tail) ---
5 rng = npy.random.default_rng(42)
6 n = len(titanic_df)
7
8 # Generate category indices with a Zipf distribution (
   skewed, many rare categories)
9 # Cap very large values to keep the number of unique
   categories reasonable
10 zipf_raw = rng.zipf(a=2.0, size=n)
11 zipf_capped = npy.minimum(zipf_raw, 600) # at most 600
   categories
12
13 titanic_df["synthetic_id"] = pd.Series(zipf_capped).map(
   lambda i: f"ID_{i:03d}")
14
15 ## Quick sanity check
16 print("Unique synthetic_id:", titanic_df["synthetic_id"].
   nunique())
17 print(titanic_df[["synthetic_id"]].head())
```

```
1 # Load seaborn titanic
2 titanic_df = sns.load_dataset("titanic").copy()
3
4 ## Create a synthetic high-cardinality feature (e.g.,
   ~600 unique IDs, Zipf-like long tail) ---
5 rng = npy.random.default_rng(42)
6 n = len(titanic_df)
7
8 # Generate category indices with a Zipf distribution (
   skewed, many rare categories)
9 # Cap very large values to keep the number of unique
   categories reasonable
```

```

10 zipf_raw = rng.zipf(a=2.0, size=n)
11 zipf_capped = npy.minimum(zipf_raw, 600) # at most 600
    categories
12
13 titanic_df["synthetic_id"] = pd.Series(zipf_capped).map(
    lambda i: f"ID_{i:03d}")
14
15 ## Quick sanity check
16 print("Unique synthetic_id:", titanic_df["synthetic_id"].
    nunique())
17 print(titanic_df[["synthetic_id"]].head())

```

```

1 Unique synthetic_id: 40
2   synthetic_id
3 0          ID_004
4 1          ID_001
5 2          ID_001
6 3          ID_001
7 4          ID_001

```

Pros of frequency encoding: keeps dimensionality to 1 column; often effective for tree-based models.

Cons: can leak target information if computed improperly (always compute frequencies on **train only**, then map to valid sets in validation/test).

Practical Pitfalls: Dummy Variable Trap & Sparse Matrices

Dummy Variable Trap

Creating one dummy column for **every** category + an intercept can yield **perfect multicollinearity** in linear models.

Solutions: - Drop one level per feature: `drop_first=True` in `get_dummies()` or `drop="first"` in `OneHotEncoder`. - Use regularization (Ridge/Lasso) or

algorithms robust to collinearity (trees).

Sparse Matrices

Encoders like `OneHotEncoder` and `FeatureHasher` often produce **sparse matrices**.

- Pros: huge memory savings when most entries are zero.
- Cons: not all libraries accept sparse inputs; sometimes you'll need `X.toarray()` (careful with memory!).

Let's compare memory usage roughly for a dense vs sparse representation.

```
1 import sys
2 from scipy import sparse
3
4 # Example: One-hot encode three categorical cols
5 ohe = OneHotEncoder(handle_unknown="ignore") # default
      returns sparse
6 X_ohe_sparse = ohe.fit_transform(titanic_df[["sex", "
      class", "embarked"]])
7
8 # Convert to dense (for comparison ONLY)
9 X_ohe_dense = X_ohe_sparse.toarray()
10
11 # Rough memory comparison
12 sparse_mem = (X_ohe_sparse.data.nbytes + X_ohe_sparse.
      indptr.nbytes + X_ohe_sparse.indices.nbytes) / 1024
13 dense_mem = X_ohe_dense.nbytes / 1024
14
15 print(f"Sparse approx memory: {sparse_mem:,.1f} KB")
16 print(f"Dense approx memory: {dense_mem:,.1f} KB")
```

```
1 import sys
2 from scipy import sparse
3
4 # Example: One-hot encode three categorical cols
5 ohe = OneHotEncoder(handle_unknown="ignore") # default
      returns sparse
```

```

6 X_ohc_sparse = ohe.fit_transform(titanic_df[["sex", "
    class", "embarked"]])
7
8 # Convert to dense (for comparison ONLY)
9 X_ohc_dense = X_ohc_sparse.toarray()
10
11 # Rough memory comparison
12 sparse_mem = (X_ohc_sparse.data.nbytes + X_ohc_sparse.
    indptr.nbytes + X_ohc_sparse.indices.nbytes) / 1024
13 dense_mem = X_ohc_dense.nbytes / 1024
14
15 print(f"Sparse approx memory: {sparse_mem:,.1f} KB")
16 print(f"Dense approx memory: {dense_mem:,.1f} KB")

```

```

1 Sparse approx memory: 34.8 KB
2 Dense approx memory: 62.6 KB

```

Guideline: keep representations **sparse** as long as your downstream model supports it (e.g., linear models in scikit-learn do). Convert to dense only if strictly necessary.

Putting It Together: A Mini Encoding Pipeline

Below, we combine: - One-hot for nominal features, - Ordinal encoding for **class**,
- And pass through a couple of numeric features.

This illustrates how to **compose encodings cleanly** before modeling.

```

1 # Import required classes and utilities
2 from sklearn.pipeline import Pipeline
3 from sklearn.compose import ColumnTransformer
4 from sklearn.preprocessing import OneHotEncoder,
    OrdinalEncoder
5 from sklearn.linear_model import LogisticRegression
6 from sklearn.impute import SimpleImputer
7 import numpy as npy
8

```

```
9 # Use the Titanic DataFrame from seaborn as the working
   dataset
10 df = titanic_df.copy()
11
12 # Define feature matrix X and target vector y
13 X = df[["sex", "class", "embarked", "age", "fare"]]
14 y = df["survived"]
15
16 # Group columns by semantic type: nominal, ordinal,
   numeric
17 nominal = ["sex", "embarked"]
18 ordinal = ["class"]
19 numeric = ["age", "fare"]
20
21 # Specify the explicit order for the ordinal categorical
   feature
22 ordinal_order = ["Third", "Second", "First"]
23
24 # OneHotEncoder compatibility shim: prefer sparse_output=
   False, fall back to sparse=False on older scikit-
   learn
25 try:
26     ohe = OneHotEncoder(handle_unknown="ignore",
                           sparse_output=False)
27 except TypeError:
28     ohe = OneHotEncoder(handle_unknown="ignore", sparse=
        False)
29
30 # Build the ColumnTransformer so column selection happens
   on the original pandas DataFrame
31 preprocessor = ColumnTransformer(
32     transformers=[
33         # Nominal features: impute missing values with
           the most frequent category, then one-hot
           encode
34         ("nominal_ohe",
35         Pipeline(steps=[
36             ("imputer", SimpleImputer(strategy="
               most_frequent")),
37             ("ohe", ohe),
38         ]),
```



```

39         nominal),
40         # Ordinal feature: impute missing values with the
           most frequent category, then apply
           OrdinalEncoder with explicit category order
41         ("class_ord",
42         Pipeline(steps=[
43             ("imputer", SimpleImputer(strategy="
               most_frequent")),
44             ("ord", OrdinalEncoder(categories=
               ordinal_order))),
45         ]),
46         ordinal),
47         # Numeric features: impute missing values with
           the median
48         ("num",
49         Pipeline(steps=[
50             ("imputer", SimpleImputer(strategy="median")
               ),
51         ]),
52         numeric),
53     ],
54     # Drop any remaining columns and keep concise feature
       names in the output
55     remainder="drop",
56     verbose_feature_names_out=False
57 )
58
59 # Assemble the final Pipeline: preprocessing followed by
   a simple Logistic Regression classifier
60 pipe = Pipeline(steps=[
61     ("prep", preprocessor),
62     ("clf", LogisticRegression(max_iter=1000))
63 ])
64
65 # Fit the end-to-end pipeline on the training data
66 pipe.fit(X, y)
67 print("Pipeline ready. Sample prediction:", pipe.predict(
   X.head(5)))
68
69 # (Optional) Inspect the transformed feature names and
   the final feature count

```

```
70 feat_names = pipe.named_steps["prep"].  
    get_feature_names_out()  
71 print("n_features:", len(feat_names))  
72 print("first 10 features:", feat_names[:10])
```

```
1 # Import required classes and utilities  
2 from sklearn.pipeline import Pipeline  
3 from sklearn.compose import ColumnTransformer  
4 from sklearn.preprocessing import OneHotEncoder,  
    OrdinalEncoder  
5 from sklearn.linear_model import LogisticRegression  
6 from sklearn.impute import SimpleImputer  
7 import numpy as npy  
8  
9 # Use the Titanic DataFrame from seaborn as the working  
    dataset  
10 df = titanic_df.copy()  
11  
12 # Define feature matrix X and target vector y  
13 X = df[["sex", "class", "embarked", "age", "fare"]]  
14 y = df["survived"]  
15  
16 # Group columns by semantic type: nominal, ordinal,  
    numeric  
17 nominal = ["sex", "embarked"]  
18 ordinal = ["class"]  
19 numeric = ["age", "fare"]  
20  
21 # Specify the explicit order for the ordinal categorical  
    feature  
22 ordinal_order = ["Third", "Second", "First"]  
23  
24 # OneHotEncoder compatibility shim: prefer sparse_output=  
    False, fall back to sparse=False on older scikit-  
    learn  
25 try:  
26     ohe = OneHotEncoder(handle_unknown="ignore",  
        sparse_output=False)  
27 except TypeError:  
28     ohe = OneHotEncoder(handle_unknown="ignore", sparse=  
        False)
```

```

29
30 # Build the ColumnTransformer so column selection happens
    on the original pandas DataFrame
31 preprocessor = ColumnTransformer(
32     transformers=[
33         # Nominal features: impute missing values with
            the most frequent category, then one-hot
            encode
34         ("nominal_ohe",
35          Pipeline(steps=[
36              ("imputer", SimpleImputer(strategy="
37                  most_frequent")),
38              ("ohe", ohe),
39          ]),
40         nominal),
41         # Ordinal feature: impute missing values with the
            most frequent category, then apply
            OrdinalEncoder with explicit category order
42         ("class_ord",
43          Pipeline(steps=[
44              ("imputer", SimpleImputer(strategy="
45                  most_frequent")),
46              ("ord", OrdinalEncoder(categories=
47                  ordinal_order)),
48          ]),
49         ordinal),
50         # Numeric features: impute missing values with
            the median
51         ("num",
52          Pipeline(steps=[
53              ("imputer", SimpleImputer(strategy="median"
54              )
55          ]),
56         numeric),
57     ],
58     # Drop any remaining columns and keep concise feature
        names in the output
59     remainder="drop",
60     verbose_feature_names_out=False
61 )

```

```
59 # Assemble the final Pipeline: preprocessing followed by
    a simple Logistic Regression classifier
60 pipe = Pipeline(steps=[
61     ("prep", preprocessor),
62     ("clf", LogisticRegression(max_iter=1000))
63 ])
64
65 # Fit the end-to-end pipeline on the training data
66 pipe.fit(X, y)
67 print("Pipeline ready. Sample prediction:", pipe.predict(
    X.head(5)))
68
69 # (Optional) Inspect the transformed feature names and
    the final feature count
70 feat_names = pipe.named_steps["prep"].
    get_feature_names_out()
71 print("n_features:", len(feat_names))
72 print("first 10 features:", feat_names[:10])
```

```
1 Pipeline ready. Sample prediction: [0 1 1 1 0]
2 n_features: 8
3 first 10 features: ['sex_female' 'sex_male' 'embarked_C'
    'embarked_Q' 'embarked_S' 'class'
4    'age' 'fare']
```

Summary

- **Categorical vs numerical:** determine type and cardinality early.
- **LabelEncoder:** typically for **targets** or **true ordinal** features; avoid for nominal features in general.
- **One-Hot:** safe default for nominal data; watch dimensionality.
- **OrdinalEncoder:** only when order is meaningful (pass explicit category order).
- **High-cardinality:** prefer **frequency encoding** or **hashing** to limit dimensions.

- **Pitfalls:** avoid the **dummy variable trap** in linear models; keep **sparse** encodings when possible.

Next, we'll move to **Scaling Numerical Features**, comparing standardization, min-max scaling, and robust scaling with practical guidance on when each is appropriate.

Scaling Numerical Features

Why scaling is necessary

Several algorithms assume features live on comparable scales or use distances/gradients that are sensitive to scale: - **Distance-based models**: k-Nearest Neighbors, SVM with RBF/poly kernels, K-Means. - **Gradient-based models**: Linear/Logistic Regression with regularization, Neural Networks (faster, more stable convergence). - **Regularization**: L1/L2 penalties depend on coefficient magnitudes; standardized features make penalties comparable.

When scaling is often **not** needed: - **Tree-based models** (Decision Trees, Random Forests, Gradient Boosted Trees) are invariant to monotonic feature scaling. - **Naive Bayes** (in many variants) is mostly unaffected by linear rescaling.

Below we compare the main scalers and show when each one shines.

Standardization with StandardScaler

What it does: subtract the mean and divide by the standard deviation → each feature roughly **mean 0, std 1**.

When to use: default choice for many linear models and distance-based models when distributions are not heavy-tailed.

```
1 import numpy as npy
```

```
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_iris
5 from sklearn.preprocessing import StandardScaler
6
7 # Load Iris (all-numeric, clean)
8 iris = load_iris(as_frame=True)
9 X_iris = iris.frame.drop(columns=["target"]) if "target"
10     in iris.frame else iris.data
11 cols = X_iris.columns
12
13 scaler_std = StandardScaler()
14 X_std = pd.DataFrame(scaler_std.fit_transform(X_iris),
15     columns=cols)
16
17 # Quick sanity: means ~0, stds ~1
18 summary = pd.DataFrame({
19     "mean_before": X_iris.mean(),
20     "std_before": X_iris.std(),
21     "mean_after": X_std.mean(),
22     "std_after": X_std.std()
23 }).round(3)
24 display(summary)
25
26 # Plot histograms AFTER scaling (adds value: shows
27     standardized spread)
28 fig, axes = plt.subplots(2, 2, figsize=(9, 6))
29 axes = axes.ravel()
30 for ax, c in zip(axes, cols):
31     ax.hist(X_std[c], bins=20)
32     ax.set_title(f"{c} (standardized)")
33 plt.tight_layout()
34 plt.show()
```

```
1 import numpy as npy
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_iris
5 from sklearn.preprocessing import StandardScaler
6
7 # Load Iris (all-numeric, clean)
```

```
8 iris = load_iris(as_frame=True)
9 X_iris = iris.frame.drop(columns=["target"]) if "target"
   in iris.frame else iris.data
10 cols = X_iris.columns
11
12 scaler_std = StandardScaler()
13 X_std = pd.DataFrame(scaler_std.fit_transform(X_iris),
   columns=cols)
14
15 # Quick sanity: means ~0, stds ~1
16 summary = pd.DataFrame({
17     "mean_before": X_iris.mean(),
18     "std_before": X_iris.std(),
19     "mean_after": X_std.mean(),
20     "std_after": X_std.std()
21 }).round(3)
22 display(summary)
```

mean_before

std_before

mean_after

std_after

sepal length (cm)

5.843

0.828

-0.0

1.003

sepal width (cm)

3.057

0.436

-0.0

1.003

petal length (cm)

3.758

1.765

-0.0

1.003

petal width (cm)

1.199

0.762

-0.0

1.003

```
1 import numpy as npy
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.datasets import load_iris
5 from sklearn.preprocessing import StandardScaler
6
7
8 # Plot histograms AFTER scaling (adds value: shows
   standardized spread)
9 fig, axes = plt.subplots(2, 2, figsize=(9, 6))
10 axes = axes.ravel()
11 for ax, c in zip(axes, cols):
12     ax.hist(X_std[c], bins=20)
13     ax.set_title(f"{c} (standardized)")
14 plt.tight_layout()
15 plt.show()
```

```
1 # Plot histograms AFTER scaling (adds value: shows
   standardized spread)
2 fig, axes = plt.subplots(2, 2, figsize=(9, 6))
3 axes = axes.ravel()
```

```
4 for ax, c in zip(axes, cols):
5     ax.hist(X_std[c], bins=20, color = "#FECB00",
6             edgcolor = '0.3',)
7     ax.set_title(f"{c} (standardized)")
8 plt.tight_layout()
9 plt.show()
```

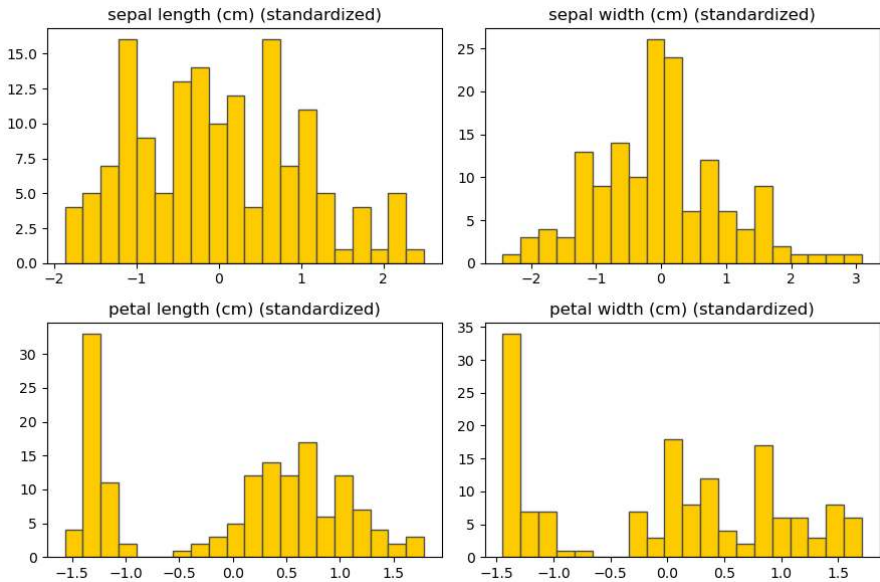


Figure 1: Image generated by the provided code.

Reading the plot/table: After standardization, each feature centers around zero with unit variance. This makes gradient steps comparable across coordinates and distances more meaningful across features.

Min-Max Scaling with MinMaxScaler

- **What it does:** linearly maps each feature to a bounded range, typically **[0, 1]**.
- **When to use:** algorithms that assume bounded inputs (e.g., some neural nets/activation setups), or when you need to preserve the original distribution shape but compress to a fixed interval.
- **Caveat:** very sensitive to **outliers** (min and max are extreme points).

```
1 from sklearn.preprocessing import MinMaxScaler
2
3 scaler_mm = MinMaxScaler(feature_range=(0, 1))
4 X_mm = pd.DataFrame(scaler_mm.fit_transform(X_iris),
5                       columns=cols)
6
7 # Show min/max after scaling
8 mm_bounds = pd.DataFrame({
9     "min_after": X_mm.min(),
10    "max_after": X_mm.max()
11 }).round(3)
12 display(mm_bounds)
```

```
1 from sklearn.preprocessing import MinMaxScaler
2
3 scaler_mm = MinMaxScaler(feature_range=(0, 1))
4 X_mm = pd.DataFrame(scaler_mm.fit_transform(X_iris),
5                       columns=cols)
6
7 # Show min/max after scaling
8 mm_bounds = pd.DataFrame({
9     "min_after": X_mm.min(),
10    "max_after": X_mm.max()
11 }).round(3)
12 display(mm_bounds)
```

min_after

max_after

sepal length (cm)

0.0

1.0

sepal width (cm)

0.0

1.0

petal length (cm)

0.0

1.0

petal width (cm)

0.0

1.0

```
1 # 2D scatter before vs after (adds value: visual bound to [0,1])
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 ax[0].scatter(X_iris[cols[0]], X_iris[cols[1]], s=12)
4 ax[0].set_title(f"Raw: {cols[0]} vs {cols[1]}")
5 ax[1].scatter(X_mm[cols[0]], X_mm[cols[1]], s=12)
6 ax[1].set_title(f"Min-Max: {cols[0]} vs {cols[1]}")
7 plt.tight_layout()
8 plt.show()
```

```
1 # 2D scatter before vs after (adds value: visual bound to [0,1])
2 fig, ax = plt.subplots(1, 2, figsize=(10, 4))
3 ax[0].scatter(X_iris[cols[0]], X_iris[cols[1]], s = 50,
4               color = '#6b99a6ff')
5 ax[1].scatter(X_mm[cols[0]], X_mm[cols[1]], s = 50, color
6               = '#6b99a6ff')
7 ax[1].set_title(f"Min-Max: {cols[0]} vs {cols[1]}")
```

```

7 plt.tight_layout()
8 plt.show()

```

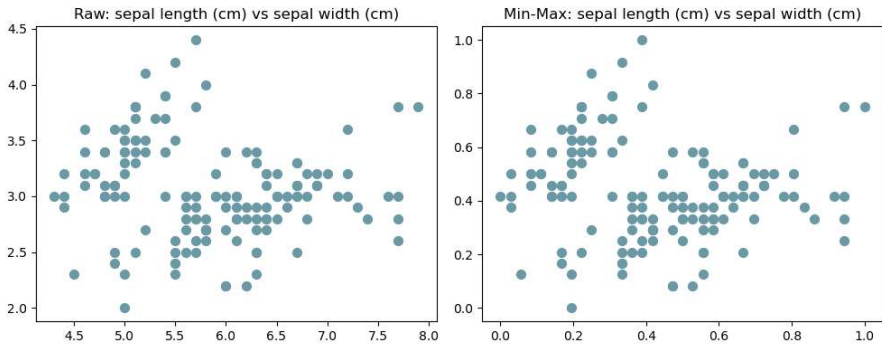


Figure 2: Image generated by the provided code.

Reading the plot/table: All features are squeezed into [0, 1]. Relative positions are preserved, but outliers can dominate the mapping.

Robust Scaling with RobustScaler

- **What it does:** centers using the **median** and scales by the **interquartile range (IQR)**.
- **When to use:** data with **outliers** or heavy tails; keeps the bulk of the distribution comparable without being distorted by extremes.

```

1 from sklearn.preprocessing import RobustScaler
2
3 # Create a synthetic feature with clear outliers to show
  the benefit
4 rng = np.random.default_rng(7)
5 x_main = rng.normal(loc=50, scale=10, size=400)

```

```
6 x_out = rng.normal(loc=120, scale=5, size=10) #
   outliers
7 x = np.concatenate([x_main, x_out])
8
9 X_synth = pd.DataFrame({"feature": x})
10
11 scaler_std_s = StandardScaler()
12 scaler_rob = RobustScaler()
13
14 X_std_s = pd.DataFrame(scaler_std_s.fit_transform(X_synth
   ), columns=["feature"])
15 X_rob = pd.DataFrame(scaler_rob.fit_transform(X_synth),
   columns=["feature"])
16
17 # Boxplots: standard vs robust (adds value: robust
   reduces outlier distortion)
18 fig, ax = plt.subplots(1, 2, figsize=(8, 4))
19 ax[0].boxplot(X_std_s["feature"], vert=True)
20 ax[0].set_title("StandardScaler")
21 ax[1].boxplot(X_rob["feature"], vert=True)
22 ax[1].set_title("RobustScaler")
23 plt.tight_layout()
24 plt.show()
```

```
1 from sklearn.preprocessing import RobustScaler
2
3 # Create a synthetic feature with clear outliers to show
   the benefit
4 rng = npy.random.default_rng(7)
5 x_main = rng.normal(loc=50, scale=10, size=400)
6 x_out = rng.normal(loc=120, scale=5, size=10) #
   outliers
7 x = npy.concatenate([x_main, x_out])
8
9 X_synth = pd.DataFrame({"feature": x})
10
11 scaler_std_s = StandardScaler()
12 scaler_rob = RobustScaler()
13
14 X_std_s = pd.DataFrame(scaler_std_s.fit_transform(X_synth
   ), columns=["feature"])
```

```

15 X_rob = pd.DataFrame(scaler_rob.fit_transform(X_synth),
16                       columns=["feature"])
17 # Boxplots: standard vs robust (adds value: robust
18   reduces outlier distortion)
19 fig, ax = plt.subplots(1, 2, figsize=(8, 4))
20 ax[0].boxplot(X_std_s["feature"], vert=True)
21 ax[0].set_title("StandardScaler")
22 ax[1].boxplot(X_rob["feature"], vert=True)
23 ax[1].set_title("RobustScaler")
24 plt.tight_layout()
25 plt.show()

```

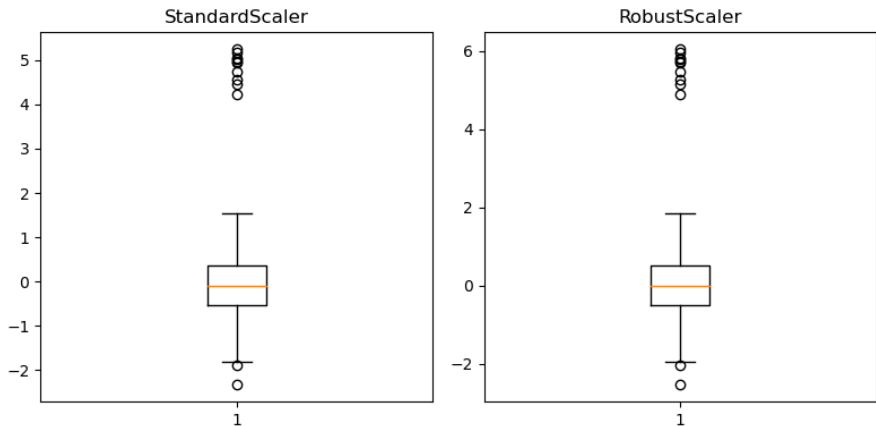


Figure 3: Image generated by the provided code.

Reading the plot: Robust scaling compresses the central mass similarly across features while preventing a few extreme values from stretching the scale.

Comparing scaling methods on the same dataset

We quantify the impact on a **distance-based model**: k-Nearest Neighbors (kNN) regression on the California Housing dataset.

We compare **no scaling**, **StandardScaler**, and **MinMaxScaler** using cross-validated RMSE.

(We subsample to keep runtime reasonable.)

```
1  from sklearn.datasets import fetch_california_housing
2  from sklearn.model_selection import cross_val_score
3  from sklearn.neighbors import KNeighborsRegressor
4  from sklearn.pipeline import Pipeline
5
6  # Load data
7  housing = fetch_california_housing(as_frame=True)
8  X_h = housing.frame.drop(columns=["MedHouseVal"])
9  y_h = housing.frame["MedHouseVal"]
10
11 # Subsample to speed up (adjust as needed)
12 sample = X_h.sample(n=5000, random_state=42)
13 target = y_h.loc[sample.index]
14
15 def rmse(scores):
16     return np.sqrt(-scores.mean())
17
18 pipelines = {
19     "No Scaling": Pipeline([("model", KNeighborsRegressor(
20         n_neighbors=5))]),
21     "StandardScaler": Pipeline([("scaler", StandardScaler(
22         )), ("model", KNeighborsRegressor(n_neighbors=5
23         ))]),
24     "MinMaxScaler": Pipeline([("scaler", MinMaxScaler()),
25         ("model", KNeighborsRegressor(n_neighbors=5))]),
26 }
27
28 results = {}
29 for name, pipe in pipelines.items():
30     scores = cross_val_score(pipe, sample, target,
31         scoring="neg_mean_squared_error", cv=5, n_jobs
```



```

    ==-1)
27     results[name] = [rmse(scores)]
28
29     results_df = pd.DataFrame(results).T
30     results_df.columns = ["RMSE"]
31     results_df = results_df.sort_values("RMSE")
32
33     # Bar plot (adds value: quick visual comparison)
34     fig, ax = plt.subplots(figsize=(6, 4))
35     ax.bar(results_df.index, results_df["RMSE"])
36     ax.set_ylabel("RMSE (lower is better)")
37     ax.set_title("kNN Regression - Effect of Scaling")
38     plt.tight_layout()
39     plt.show()

```

```

1  from sklearn.datasets import fetch_california_housing
2  from sklearn.model_selection import cross_val_score
3  from sklearn.neighbors import KNeighborsRegressor
4  from sklearn.pipeline import Pipeline
5
6  # Load data
7  housing = fetch_california_housing(as_frame=True)
8  X_h = housing.frame.drop(columns=["MedHouseVal"])
9  y_h = housing.frame["MedHouseVal"]
10
11 # Subsample to speed up (adjust as needed)
12 sample = X_h.sample(n=5000, random_state=42)
13 target = y_h.loc[sample.index]
14
15 def rmse(scores):
16     return npy.sqrt(-scores.mean())
17
18 pipelines = {
19     "No Scaling": Pipeline([("model", KNeighborsRegressor(
20         n_neighbors=5))]),
21     "StandardScaler": Pipeline([("scaler", StandardScaler(
22         )), ("model", KNeighborsRegressor(n_neighbors=5)
23         )]),
24     "MinMaxScaler": Pipeline([("scaler", MinMaxScaler()),
25         ("model", KNeighborsRegressor(n_neighbors=5))]),
26 }

```

```
23
24 results = {}
25 for name, pipe in pipelines.items():
26     scores = cross_val_score(pipe, sample, target,
27                               scoring="neg_mean_squared_error", cv=5, n_jobs
28                               =-1)
29     results[name] = [rmse(scores)]
30
31 results_df = pd.DataFrame(results).T
32 results_df.columns = ["RMSE"]
33 results_df = results_df.sort_values("RMSE")
34
35 # Bar plot (adds value: quick visual comparison)
36 fig, ax = plt.subplots(figsize=(6, 4))
37 ax.bar(results_df.index, results_df["RMSE"], color = ['#
38         b7d3dbff', '#6b99a6ff', '#3a6674ff'])
39 ax.set_ylabel("RMSE (lower is better)")
40 ax.set_title("kNN Regression - Effect of Scaling")
41 plt.tight_layout()
42 plt.show()
```

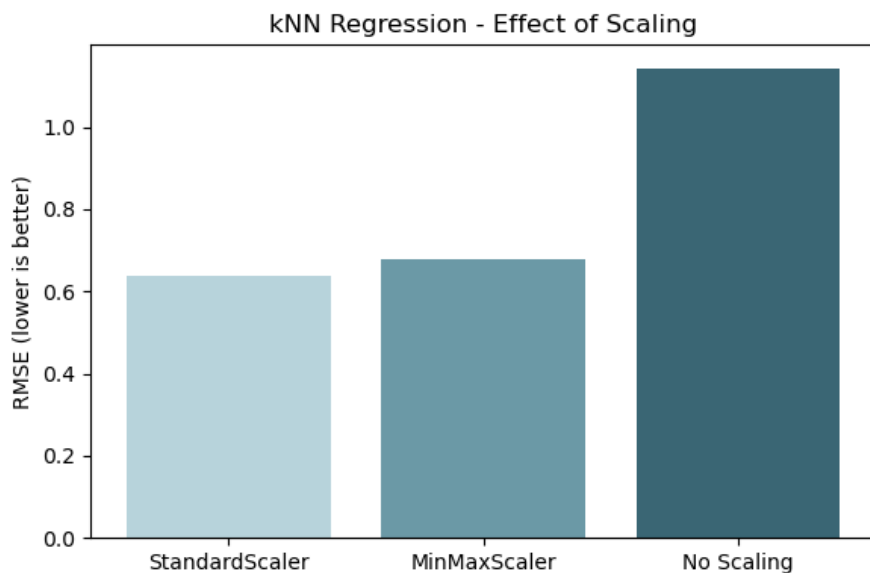


Figure 4: Image generated by the provided code.

Interpreting the numbers: Distance-based models typically benefit from scaling. Which scaler wins can depend on feature distributions; standardization is a robust default, MinMax can help if you want all features in a fixed range.

Practical tips: essential vs optional

- **Essential / highly recommended**

- kNN, K-Means, SVM with kernels, PCA/UMAP/t-SNE, Linear/Logistic Regression with regularization, Neural Networks.
- Mixed-unit features (e.g., `income` in thousands vs `rooms` count) almost always need scaling.

- **Useful but situational**

- When features have very different variances; standardization evens their influence.
- When optimization is unstable or slow to converge.

- **Usually optional**

- Tree-based methods (Random Forests, Gradient Boosted Trees, XG-Boost/LightGBM/CatBoost).
- Models that do not use distances or gradients in ways sensitive to scale.

Rule of thumb: start with **StandardScaler**; consider **RobustScaler** if outliers dominate; use **MinMaxScaler** when a bounded range is required or when feeding into methods expecting $[0, 1]$.

What comes next

We have the core scalers under our belt. Next chapter, we'll look at **transformations** to tame skewed distributions: **log transforms**, **Box-Cox**, and **Yeo-Johnson**, and how to combine them inside **pipelines**.

Transforming Variables

Detecting skewness and choosing a transformation

Many real-world features are **skewed** (long tails). Skewness can harm models that assume near-normal inputs or use distances/gradients. Typical fixes:

- **Log transform**: for strictly positive, right-skewed variables.
- **Box-Cox / Yeo-Johnson**: learn power parameters; Box-Cox needs strictly positive values; Yeo-Johnson handles zero/negative.
- **Quantile transforms**: reshape distributions to **normal** or **uniform** via rank mapping.
- **Binarization**: convert to indicator features when only a threshold matters.

We'll:

- 1) measure skewness,
- 2) apply several transforms,
- 3) compare before/after on representative features.

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 from sklearn.datasets import fetch_california_housing
5 from sklearn.preprocessing import FunctionTransformer,
   PowerTransformer, QuantileTransformer, Binarizer
6
7 # Load California Housing (all numeric, several skewed)
8 cal = fetch_california_housing(as_frame=True)
```

```

9 X = cal.frame.drop(columns=["MedHouseVal"]) # predictors
  only
10 cols = X.columns
11
12 # Compute skewness
13 skew_before = X.skew(numeric_only=True).sort_values(
    ascending=False)
14 skew_before.name = "skew_before"
15
16 display(pd.DataFrame(skew_before).T) # quick overview (
    top row sorted)

```

```

1 import pandas as pd
2 import matplotlib.pyplot as plt
3
4 from sklearn.datasets import fetch_california_housing
5 from sklearn.preprocessing import FunctionTransformer,
    PowerTransformer, QuantileTransformer, Binarizer
6
7 # Load California Housing (all numeric, several skewed)
8 cal = fetch_california_housing(as_frame=True)
9 X = cal.frame.drop(columns=["MedHouseVal"]) # predictors
    only
10 cols = X.columns
11
12 # Compute skewness
13 skew_before = X.skew(numeric_only=True).sort_values(
    ascending=False)
14 skew_before.name = "skew_before"
15
16 display(pd.DataFrame(skew_before).T) # quick overview (
    top row sorted)

```

AveOccup

AveBedrms

AveRooms

Population

MedInc

Latitude

HouseAge

Longitude

skew_before

97.639561

31.316956

20.697869

4.935858

1.646657

0.465953

0.060331

-0.297801

Log transformations for skewed data

When to use: strictly positive, right-skewed features (counts, incomes, areas). We will use `numpy.log1p(x)` ($\log(1+x)$) which is stable near zero.

We will illustrate on two commonly skewed features from this dataset:

- `MedInc` (median income)
- `Population`

We will verify positivity and then compare histograms before/after.

```
1 def plot_hist_pair(series_raw, series_tx, title_raw,
2   title_tx):
   fig, axes = plt.subplots(1, 2, figsize=(10, 4))
```



```

3     axes[0].hist(series_raw, bins = 30, color = "#FECB00",
4                   , edgecolor = '0.3',)
5     axes[0].set_title(title_raw)
6     axes[1].hist(series_tx, bins = 30, color = "#FECB00",
7                   , edgecolor = '0.3',)
8     axes[1].set_title(title_tx)
9     plt.tight_layout()
10    plt.show()
11
12    # Choose two features
13    feat_log_candidates = ["MedInc", "Population"]
14    for feat in feat_log_candidates:
15        s = X[feat].dropna()
16        # Ensure strictly positive for log
17        s_pos = s[s > 0]
18        s_log = npy.log1p(s_pos)
19
20        plot_hist_pair(
21            s_pos, s_log,
22            f"{feat} (raw)",
23            f"{feat} (log1p)"
24        )
25
26    # Skewness comparison table for selected features
27    log_rows = []
28    for feat in feat_log_candidates:
29        s = X[feat].dropna()
30        s_pos = s[s > 0]
31        s_log = npy.log1p(s_pos)
32        log_rows.append({
33            "feature": feat,
34            "skew_raw": float(s_pos.skew()),
35            "skew_log1p": float(pd.Series(s_log).skew()),
36        })

```

```

1    def plot_hist_pair(series_raw, series_tx, title_raw,
2                      title_tx):
3        fig, axes = plt.subplots(1, 2, figsize=(10, 4))
4        axes[0].hist(series_raw, bins = 30, color = "#FECB00",
5                      , edgecolor = '0.3',)
6        axes[0].set_title(title_raw)

```

```
5     axes[1].hist(series_tx, bins = 30, color = "#FECB00",
6                 edgcolor = '0.3',)
7     axes[1].set_title(title_tx)
8     plt.tight_layout()
9     plt.show()
10
11 # Choose two features
12 feat_log_candidates = ["MedInc", "Population"]
13 for feat in feat_log_candidates:
14     s = X[feat].dropna()
15     # Ensure strictly positive for log
16     s_pos = s[s > 0]
17     s_log = npy.log1p(s_pos)
18
19     plot_hist_pair(
20         s_pos, s_log,
21         f"{feat} (raw)",
22         f"{feat} (log1p)"
23     )
24
25 # Skewness comparison table for selected features
26 log_rows = []
27 for feat in feat_log_candidates:
28     s = X[feat].dropna()
29     s_pos = s[s > 0]
30     s_log = npy.log1p(s_pos)
31     log_rows.append({
32         "feature": feat,
33         "skew_raw": float(s_pos.skew()),
34         "skew_log1p": float(pd.Series(s_log).skew()),
35     })
```

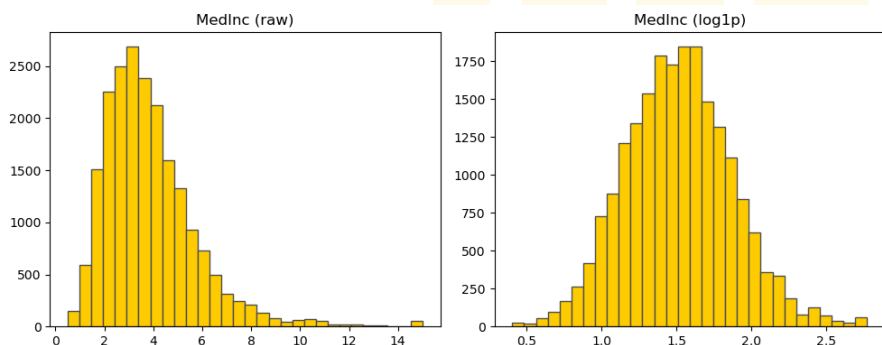


Figure 1: Image generated by the provided code.

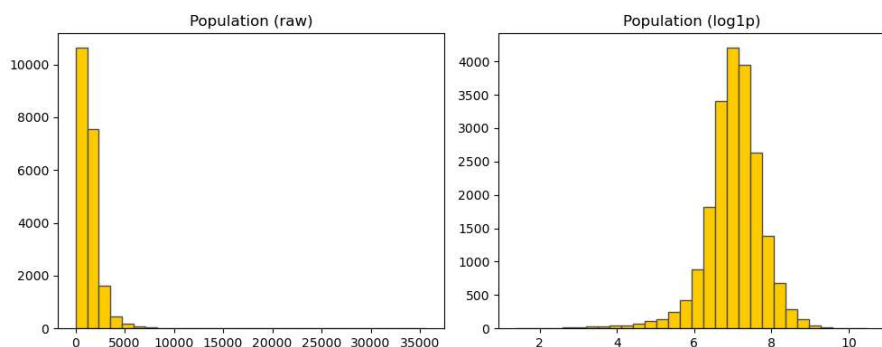


Figure 2: Image generated by the provided code.

Reading the plots/table: The log transform typically reduces right skewness and compresses long tails, making gradients and distances more stable.

Box-Cox and Yeo-Johnson with PowerTransformer

PowerTransformer estimates a power parameter per feature to make the data more Gaussian-like.

- **Box-Cox** works only for **strictly positive** data.
- **Yeo-Johnson** works for **zero/negative** as well and is often a safe default.

We will apply both where possible and compare skewness reductions.

```
1 from sklearn.compose import ColumnTransformer
2
3 # Identify strictly positive columns for Box-Cox (min >
  0)
4 positive_cols = [c for c in cols if (X[c] > 0).all()]
5 non_positive_cols = [c for c in cols if c not in
  positive_cols]
6
7 pt_boxcox = PowerTransformer(method="box-cox",
  standardize=True)
8 pt_yeojohn = PowerTransformer(method="yeo-johnson",
  standardize=True)
9
10 # Fit/transform separate copies for clarity
11 X_box = X[positive_cols].copy()
12 X_box_tx = pd.DataFrame(pt_boxcox.fit_transform(X_box),
  columns=positive_cols, index=X_box.index)
13
14 X_yj = X.copy()
15 X_yj_tx = pd.DataFrame(pt_yeojohn.fit_transform(X_yj),
  columns=cols, index=X.index)
16
17 # Skewness summaries (median across features as a compact
  signal)
18 summary_skew = pd.DataFrame({
19     "median_skew_raw": [X.skew().median()],
20     "median_skew_boxcox_(pos_only)": [X_box_tx.skew().
  median() if len(positive_cols) else npy.nan],
21     "median_skew_yeojohnson": [X_yj_tx.skew().median()]
22 }).round(3)
```

```

23 display(summary_skew)
24
25 # Visual example on MedInc with Yeo-Johnson (adds value
    if different from log)
26 feat = "MedInc"
27 s_raw = X[feat]
28 s_yj = X_yj_tx[feat]
29 plot_hist_pair(s_raw, s_yj, f"{feat} (raw)", f"{feat} (
    Yeo-Johnson)", "fig_medinc_yj_hist.png")

```

```

1  from sklearn.compose import ColumnTransformer
2
3  # Identify strictly positive columns for Box-Cox (min >
    0)
4  positive_cols = [c for c in cols if (X[c] > 0).all()]
5  non_positive_cols = [c for c in cols if c not in
    positive_cols]
6
7  pt_boxcox = PowerTransformer(method="box-cox",
    standardize=True)
8  pt_yeojohn = PowerTransformer(method="yeo-johnson",
    standardize=True)
9
10 # Fit/transform separate copies for clarity
11 X_box = X[positive_cols].copy()
12 X_box_tx = pd.DataFrame(pt_boxcox.fit_transform(X_box),
    columns=positive_cols, index=X_box.index)
13
14 X_yj = X.copy()
15 X_yj_tx = pd.DataFrame(pt_yeojohn.fit_transform(X_yj),
    columns=cols, index=X.index)
16
17 # Skewness summaries (median across features as a compact
    signal)
18 summary_skew = pd.DataFrame({
19     "median_skew_raw": [X.skew().median()],
20     "median_skew_boxcox_(pos_only)": [X_box_tx.skew().
        median() if len(positive_cols) else npy.nan],
21     "median_skew_yeojohnson": [X_yj_tx.skew().median()]
22 }).round(3)
23

```

```
24 # Visual example on MedInc with Yeo-Johnson (adds value
    if different from log)
25 feat = "MedInc"
26 s_raw = X[feat]
27 s_yj = X_yj_tx[feat]
28 plot_hist_pair(s_raw, s_yj, f"{feat} (raw)", f"{feat} (
    Yeo-Johnson)", )
```

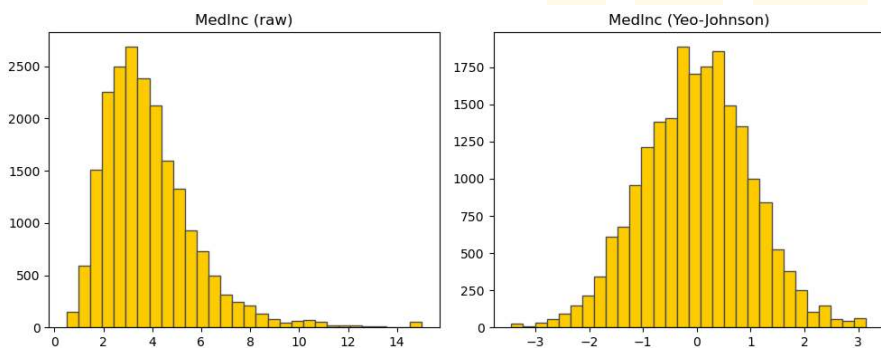


Figure 3: Image generated by the provided code.

Interpretation: Both Box-Cox and Yeo-Johnson usually reduce skewness. Yeo-Johnson is more flexible and avoids the “strictly positive” constraint.

Quantile transformations for non-Gaussian distributions

QuantileTransformer maps values to their empirical quantiles and then to a target distribution:

- `output_distribution="normal"` → approximately standard normal.

- `output_distribution="uniform"` → approximately uniform on $[0, 1]$.

This can dramatically stabilize highly non-Gaussian features, at the cost of non-linear distortions.

```

1 qt_norm = QuantileTransformer(output_distribution="normal",
2                                random_state=42)
3
4 qt_unif = QuantileTransformer(output_distribution="uniform",
5                                random_state=42)
6
7 feat = "Population" # typically heavy-tailed
8 s = X[feat].values.reshape(-1, 1)
9
10 s_norm = qt_norm.fit_transform(s).ravel()
11 s_unif = qt_unif.fit_transform(s).ravel()
12
13 # Histograms to show target distributions
14 fig, axes = plt.subplots(1, 3, figsize=(12, 4))
15 axes[0].hist(s.ravel(), bins=30)
16 axes[0].set_title(f"{feat} (raw)")
17 axes[1].hist(s_norm, bins=30)
18 axes[1].set_title(f"{feat} (Quantile -> Normal)")
19 axes[2].hist(s_unif, bins=30)
20 axes[2].set_title(f"{feat} (Quantile -> Uniform)")
21 plt.tight_layout()
22 plt.show()

```

```

1 qt_norm = QuantileTransformer(output_distribution="normal",
2                                random_state=42)
3
4 qt_unif = QuantileTransformer(output_distribution="uniform",
5                                random_state=42)
6
7 feat = "Population" # typically heavy-tailed
8 s = X[feat].values.reshape(-1, 1)
9
10 s_norm = qt_norm.fit_transform(s).ravel()
11 s_unif = qt_unif.fit_transform(s).ravel()
12
13 # Histograms to show target distributions

```

```

11 fig, axes = plt.subplots(1, 3, figsize=(12, 4))
12 axes[0].hist(s.ravel(), bins=30, color = "#FECB00",
13             edgcolor = '0.3',)
14 axes[0].set_title(f"{feat} (raw)")
15 axes[1].hist(s_norm, bins=30, color = "#FECB00",
16             edgcolor = '0.3',)
17 axes[1].set_title(f"{feat} (Quantile -> Normal)")
18 axes[2].hist(s_unif, bins=30, color = "#FECB00",
19             edgcolor = '0.3',)
20 axes[2].set_title(f"{feat} (Quantile -> Uniform)")
21 plt.tight_layout()
22 plt.show()

```

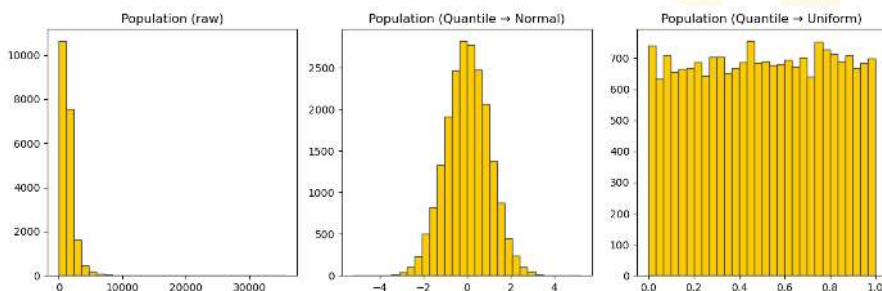


Figure 4: Image generated by the provided code.

Trade-off: Quantile transforms can greatly help algorithms needing near-Gaussian inputs but may warp relationships. Consider applying them within **pipelines** and validate their impact on your metrics.

Binarization with Binarizer

Binarizer converts a numeric feature into a binary indicator using a threshold. Use it when a **decision boundary** (e.g., above/below median or a domain-specific cutoff) carries most of the predictive signal.

We will create a binary indicator for **high income** using the median of **MedInc**.

```

1  from sklearn.preprocessing import Binarizer
2
3  threshold = float(X["MedInc"].median())
4  bin_tx = Binarizer(threshold=threshold)  # > threshold
      becomes 1.0
5  medinc_bin = bin_tx.fit_transform(X[["MedInc"]]).ravel()
6
7  # Show proportion of 1s
8  prop_high = medinc_bin.mean()
9  print(f"Threshold used: {threshold:.3f}; share above
      threshold: {prop_high:.2%}")
10
11 # Simple before/after visualization (adds value: shows
      the discretization)
12 fig, axes = plt.subplots(1, 2, figsize=(9, 4))
13 axes[0].hist(X["MedInc"], bins=30, color = "#FECB00",
      edgecolor = '0.3',)
14 axes[0].set_title("MedInc (raw)")
15 axes[1].hist(medinc_bin, bins=2, color = "#FECB00",
      edgecolor = '0.3',)
16 axes[1].set_xticks([0, 1])
17 axes[1].set_title("MedInc (binarized)")
18 plt.tight_layout()
19 plt.show()

```

```

1  from sklearn.preprocessing import Binarizer
2
3  threshold = float(X["MedInc"].median())
4  bin_tx = Binarizer(threshold=threshold)  # > threshold
      becomes 1.0
5  medinc_bin = bin_tx.fit_transform(X[["MedInc"]]).ravel()
6
7  # Show proportion of 1s
8  prop_high = medinc_bin.mean()
9  print(f"Threshold used: {threshold:.3f}; share above
      threshold: {prop_high:.2%}")
10
11 # Simple before/after visualization (adds value: shows
      the discretization)

```

```
12 fig, axes = plt.subplots(1, 2, figsize=(9, 4))
13 axes[0].hist(X["MedInc"], bins=30, color = "#FECB00",
14             edgcolor = '0.3',)
15 axes[0].set_title("MedInc (raw)")
16 axes[1].hist(medinc_bin, bins=2, color = "#FECB00",
17             edgcolor = '0.3',)
18 axes[1].set_xticks([0, 1])
19 axes[1].set_title("MedInc (binarized)")
20 plt.tight_layout()
21 plt.show()
```

1 Threshold used: 3.535; share above threshold: 50.00%

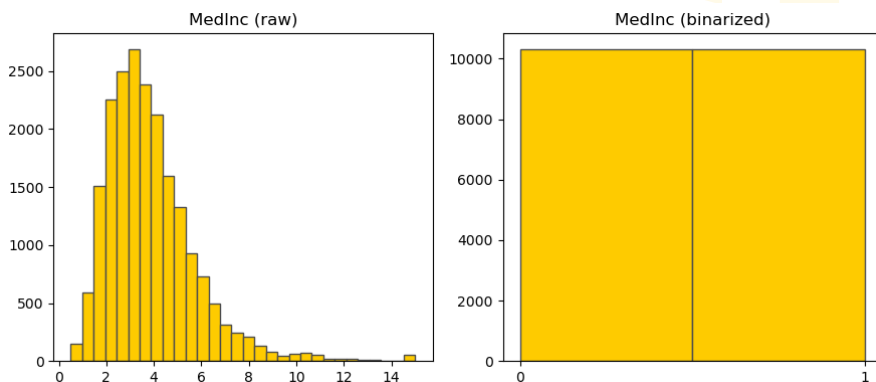


Figure 5: Image generated by the provided code.

Caution: Binarization discards within-bin variation. It can be powerful for rule-like signals but may reduce information for models that benefit from continuous structure.

Combining transformations with pipelines

We will:

- Apply **Yeo-Johnson** to skewed features,
- Leave approximately symmetric features as-is (or standardize),
- Add a **binarized** version of **MedInc** via a separate pipeline branch,
- Train a simple **Ridge** regressor on the California Housing target to illustrate end-to-end usage.

We will compare RMSE with and without transformations.

```

1  from sklearn.pipeline import Pipeline
2  from sklearn.compose import ColumnTransformer
3  from sklearn.linear_model import Ridge
4  from sklearn.model_selection import cross_val_score
5
6  target = cal.frame["MedHouseVal"]
7
8  # Heuristic: mark skewed if |skew| >= 1.0
9  skews = X.skew().abs()
10 skewed_cols = skews[skews >= 1.0].index.tolist()
11 other_cols = [c for c in cols if c not in skewed_cols]
12
13 # Pipelines for branches
14 yj_pipe = Pipeline(steps=[
15     ("yeojohn", PowerTransformer(method="yeo-johnson",
16     standardize=True))
17 ])
18 # "Passthrough" for other columns (you could add
19     StandardScaler if desired)
20 pass_pipe = "passthrough"
21 # Binarized MedInc as an extra feature (single-column
22     pipeline)
23 bin_col = ["MedInc"]
24 bin_threshold = float(X["MedInc"].median())
25 bin_pipe = Pipeline(steps=[

```

```
25     ("bin", Binarizer(threshold=bin_threshold))
26 ])
27
28 # ColumnTransformer combining branches (note: binarized
    MedInc will appear alongside raw/other features)
29 preprocessor = ColumnTransformer(
30     transformers=[
31         ("yj", yj_pipe, skewed_cols),
32         ("other", pass_pipe, other_cols),
33         ("medinc_bin", bin_pipe, bin_col),
34     ],
35     remainder="drop",
36     verbose_feature_names_out=False
37 )
38
39 model = Ridge(alpha=1.0, random_state=42)
40
41 pipe_tx = Pipeline(steps=[
42     ("prep", preprocessor),
43     ("model", model)
44 ])
45
46 pipe_raw = Pipeline(steps=[
47     ("model", Ridge(alpha=1.0, random_state=42))
48 ])
49
50 def rmse(cv_scores):
51     return npy.sqrt(-cv_scores.mean())
52
53 scores_tx = cross_val_score(pipe_tx, X, target, scoring="
    neg_mean_squared_error", cv=5, n_jobs=-1)
54 scores_raw = cross_val_score(pipe_raw, X, target, scoring
    ="neg_mean_squared_error", cv=5, n_jobs=-1)
55
56 print("RMSE (no transforms):", rmse(scores_raw))
57 print("RMSE (with transforms):", rmse(scores_tx))
```

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.compose import ColumnTransformer
3 from sklearn.linear_model import Ridge
4 from sklearn.model_selection import cross_val_score
```

```
5
6 target = cal.frame["MedHouseVal"]
7
8 # Heuristic: mark skewed if |skew| >= 1.0
9 skews = X.skew().abs()
10 skewed_cols = skews[skews >= 1.0].index.tolist()
11 other_cols = [c for c in cols if c not in skewed_cols]
12
13 # Pipelines for branches
14 yj_pipe = Pipeline(steps=[
15     ("yeojohn", PowerTransformer(method="yeo-johnson",
16     standardize=True))
17 ])
18 # "Passthrough" for other columns (you could add
19     StandardScaler if desired)
20 pass_pipe = "passthrough"
21
22 # Binarized MedInc as an extra feature (single-column
23     pipeline)
24 bin_col = ["MedInc"]
25 bin_threshold = float(X["MedInc"].median())
26 bin_pipe = Pipeline(steps=[
27     ("bin", Binarizer(threshold=bin_threshold))
28 ])
29
30 # ColumnTransformer combining branches (note: binarized
31     MedInc will appear alongside raw/other features)
32 preprocessor = ColumnTransformer(
33     transformers=[
34         ("yj", yj_pipe, skewed_cols),
35         ("other", pass_pipe, other_cols),
36         ("medinc_bin", bin_pipe, bin_col),
37     ],
38     remainder="drop",
39     verbose_feature_names_out=False
40 )
41
42 model = Ridge(alpha=1.0, random_state=42)
43
44 pipe_tx = Pipeline(steps=[
```

```
42     ("prep", preprocessor),
43     ("model", model)
44 ])
45
46 pipe_raw = Pipeline(steps=[
47     ("model", Ridge(alpha=1.0, random_state=42))
48 ])
49
50 def rmse(cv_scores):
51     return npy.sqrt(-cv_scores.mean())
52
53 scores_tx = cross_val_score(pipe_tx, X, target, scoring="
    neg_mean_squared_error", cv=5, n_jobs=-1)
54 scores_raw = cross_val_score(pipe_raw, X, target, scoring
    ="neg_mean_squared_error", cv=5, n_jobs=-1)
55
56 print("RMSE (no transforms):", rmse(scores_raw))
57 print("RMSE (with transforms):", rmse(scores_tx))
```

```
1 RMSE (no transforms): 0.7471805175916729
2 RMSE (with transforms): 0.7315995010450792
```

How to read the comparison: If RMSE improves with transformations, the preprocessing meaningfully stabilized the feature space for this model. If not, revisit which columns you transform and consider `StandardScaler` for non-skewed columns.

Summary

- **Log** is a simple, effective fix for positive, right-skewed data.
- **PowerTransformer** (Box-Cox / Yeo-Johnson) learns power parameters; Yeo-Johnson is more general.
- **QuantileTransformer** can enforce near-normal or uniform outputs but is highly non-linear; validate its impact.
- **Binarizer** is useful when a threshold carries most of the signal.
- Combine these steps using **pipelines** so training and inference stay consistent

and reproducible.

Feature Generation

Creating interaction features (PolynomialFeatures)

Interaction features capture multiplicative relationships between variables (e.g., $x_1 * x_2$).

They allow **linear** models to approximate non-linear boundaries without changing the model class.

We'll illustrate with two numerical predictors from California Housing and show how to create:

- *Interactions only* (no squares),
- Full *polynomial* terms up to degree 2 (including squares and interactions).

Caution: Polynomial expansion can explode dimensionality—validate with cross-validation and regularization.

```
1 import numpy as npy
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 from sklearn.datasets import fetch_california_housing
6 from sklearn.preprocessing import PolynomialFeatures
7
8 cal = fetch_california_housing(as_frame=True)
9 Xh = cal.frame.drop(columns=["MedHouseVal"])
10 y = cal.frame["MedHouseVal"]
```



```

11
12 num_pair = ["MedInc", "HouseAge"]
13 X_pair = Xh[num_pair].copy()
14
15 # Interactions only
16 poly_int = PolynomialFeatures(degree=2, include_bias=
    False, interaction_only=True)
17 X_int = pd.DataFrame(poly_int.fit_transform(X_pair),
    columns=poly_int.get_feature_names_out(num_pair))
18
19 # Full quadratic (squares + interactions)
20 poly_full = PolynomialFeatures(degree=2, include_bias=
    False, interaction_only=False)
21 X_quad = pd.DataFrame(poly_full.fit_transform(X_pair),
    columns=poly_full.get_feature_names_out(num_pair))
22
23 print("Interaction-only columns:", X_int.columns.tolist()
    )
24 print("Quadratic columns:", X_quad.columns.tolist())
25 display(X_int.head())
26 display(X_quad.head())

```

```

1 import numpy as npy
2 import pandas as pd
3 import matplotlib.pyplot as plt
4
5 from sklearn.datasets import fetch_california_housing
6 from sklearn.preprocessing import PolynomialFeatures
7
8 cal = fetch_california_housing(as_frame=True)
9 Xh = cal.frame.drop(columns=["MedHouseVal"])
10 y = cal.frame["MedHouseVal"]
11
12 num_pair = ["MedInc", "HouseAge"]
13 X_pair = Xh[num_pair].copy()
14
15 # Interactions only
16 poly_int = PolynomialFeatures(degree=2, include_bias=
    False, interaction_only=True)
17 X_int = pd.DataFrame(poly_int.fit_transform(X_pair),
    columns=poly_int.get_feature_names_out(num_pair))

```

```
18
19 # Full quadratic (squares + interactions)
20 poly_full = PolynomialFeatures(degree=2, include_bias=
    False, interaction_only=False)
21 X_quad = pd.DataFrame(poly_full.fit_transform(X_pair),
    columns=poly_full.get_feature_names_out(num_pair))
22
23 print("Interaction-only columns:", X_int.columns.tolist()
    )
24 print("Quadratic columns:", X_quad.columns.tolist())
25 display(X_int.head())
26 display(X_quad.head())
```

```
1 Interaction-only columns: ['MedInc', 'HouseAge', 'MedInc
    HouseAge']
2 Quadratic columns: ['MedInc', 'HouseAge', 'MedInc^2', '
    MedInc HouseAge', 'HouseAge^2']
```

MedInc

HouseAge

MedInc HouseAge

0

8.3252

41.0

341.3332

1

8.3014

21.0

174.3294

2

7.2574

52.0

377.3848

3

5.6431

52.0

293.4412

4

3.8462

52.0

200.0024

MedInc

HouseAge

MedInc^2

MedInc HouseAge

HouseAge^2

0

8.3252

41.0

69.308955

341.3332

1681.0

1

8.3014

21.0

68.913242

174.3294

441.0

2

7.2574

52.0

52.669855

377.3848

2704.0

3

5.6431

52.0

31.844578

293.4412

2704.0

4

3.8462

52.0

14.793254

200.0024

2704.0

When it helps: If target \approx smooth function of inputs with cross-effects (e.g., income \times house age), interactions let linear models weight those effects directly.

Feature crosses: combining categorical features

Feature crosses join categories to capture **co-occurrence patterns** (e.g., `sex` × `class` in Titanic).

Two common approaches:

- 1) **Manual cross in pandas**, then one-hot encode.
- 2) **OneHot** → **PolynomialFeatures(interaction_only=True)**: create interactions *after* one-hot (watch out for dimensionality).

We'll demonstrate the manual route for clarity and control.

```

1  import seaborn as sns
2  import numpy as npy
3  import pandas as pd
4  from sklearn.compose import ColumnTransformer
5  from sklearn.preprocessing import OneHotEncoder
6  from sklearn.linear_model import LogisticRegression
7  from sklearn.pipeline import Pipeline
8  from sklearn.model_selection import cross_val_score
9  from sklearn.impute import SimpleImputer
10
11 titanic = sns.load_dataset("titanic").copy()
12 X = titanic[["sex", "class", "embarked", "age", "fare"]]
13 y = titanic["survived"]
14
15 # Manual cross of categorical features
16 X = X.assign(sex_class=X["sex"].astype(str) + "_" + X["class"].astype(str))
17
18 # Column groups
19 nominal = ["sex", "embarked", "sex_class"]
20 numeric = ["age", "fare"]
21
22 # OneHotEncoder compatibility shim (sklearn >=1.2 uses sparse_output)
23 try:
```

```
24     ohe = OneHotEncoder(handle_unknown="ignore",
25                           sparse_output=False)
26     except TypeError:
27         ohe = OneHotEncoder(handle_unknown="ignore", sparse=
28                               False)
29
30 # Build preprocessing: impute first (train-only inside CV
31 #), then encode / passthrough
32 pre = ColumnTransformer(
33     transformers=[
34         # Nominal features: impute most frequent category
35         #, then one-hot encode
36         ("nominal", Pipeline([
37             ("imputer", SimpleImputer(strategy="
38               most_frequent")),
39             ("ohe", ohe),
40         ]), nominal),
41
42         # Numeric features: impute median (robust to
43         # outliers)
44         ("numeric", Pipeline([
45             ("imputer", SimpleImputer(strategy="median"))
46         ], numeric),
47     ],
48     verbose_feature_names_out=False
49 )
50
51 # End-to-end pipeline: preprocessing + classifier
52 pipe = Pipeline([
53     ("prep", pre),
54     ("clf", LogisticRegression(max_iter=1000))
55 ])
56
57 # Cross-validated accuracy
58 scores = cross_val_score(pipe, X, y, cv=5, scoring="
59   accuracy", n_jobs=-1)
60
61 print("Cross-validated accuracy (with crossed feature +
62   imputation):", scores.mean().round(3))
```

```
1 import seaborn as sns
```

```

2  import numpy as npy
3  import pandas as pd
4  from sklearn.compose import ColumnTransformer
5  from sklearn.preprocessing import OneHotEncoder
6  from sklearn.linear_model import LogisticRegression
7  from sklearn.pipeline import Pipeline
8  from sklearn.model_selection import cross_val_score
9  from sklearn.impute import SimpleImputer
10
11  titanic = sns.load_dataset("titanic").copy()
12  X = titanic[["sex", "class", "embarked", "age", "fare"]]
13  y = titanic["survived"]
14
15  # Manual cross of categorical features
16  X = X.assign(sex_class=X["sex"].astype(str) + "_" + X["class"].astype(str))
17
18  # Column groups
19  nominal = ["sex", "embarked", "sex_class"]
20  numeric = ["age", "fare"]
21
22  # OneHotEncoder compatibility shim (sklearn >=1.2 uses
    sparse_output)
23  try:
24      ohe = OneHotEncoder(handle_unknown="ignore",
        sparse_output=False)
25  except TypeError:
26      ohe = OneHotEncoder(handle_unknown="ignore", sparse=False)
27
28  # Build preprocessing: impute first (train-only inside CV
    ), then encode / passthrough
29  pre = ColumnTransformer(
30      transformers=[
31          # Nominal features: impute most frequent category
            , then one-hot encode
32          ("nominal", Pipeline([
33              ("imputer", SimpleImputer(strategy="
                most_frequent")),
34              ("ohe", ohe),
35          ]), nominal),

```

```
36
37     # Numeric features: impute median (robust to
38     # outliers)
39     ("numeric", Pipeline([
40         ("imputer", SimpleImputer(strategy="median"))
41     ]), numeric),
42     ],
43     verbose_feature_names_out=False
44 )
45 # End-to-end pipeline: preprocessing + classifier
46 pipe = Pipeline([
47     ("prep", pre),
48     ("clf", LogisticRegression(max_iter=1000))
49 ])
50
51 # Cross-validated accuracy
52 scores = cross_val_score(pipe, X, y, cv=5, scoring="
53     accuracy", n_jobs=-1)
54 print("Cross-validated accuracy (with crossed feature +
55     imputation):", scores.mean().round(3))
```

```
1 Cross-validated accuracy (with crossed feature +
   imputation): 0.797
```

Tip: Crosses are powerful when certain pairs of categories behave differently than each one alone (e.g., **male** in **First** class vs **Third** class).

Extracting date/time features

Timestamps encode **periodic** and **calendar** signals (seasonality, weekday vs weekend, rush hours).

Typical extractions:

- **year, month, day, weekday, hour, weekofyear.**
- **Cyclical encodings** for periodic variables (e.g., hour-of-day) using **sin/cos**.

Below we create a small example DataFrame with timestamps and extract informative features.

```

1  import pandas as pd
2  import numpy as npy
3  import matplotlib.pyplot as plt
4
5  # Synthetic event timestamps (hourly for ~2 weeks)
6  ts = pd.date_range("2023-05-01", periods=24*14, freq="H")
7  events = pd.DataFrame({
8      "timestamp": ts,
9      "value": npy.random.default_rng(0).normal(loc=0,
10         scale=1, size=len(ts))
11 })
12 # Basic calendar features
13 events["year"] = events["timestamp"].dt.year
14 events["month"] = events["timestamp"].dt.month
15 events["day"] = events["timestamp"].dt.day
16 events["weekday"] = events["timestamp"].dt.weekday # 0=
    Mon
17 events["hour"] = events["timestamp"].dt.hour
18 events["is_weekend"] = events["weekday"].isin([5, 6]).
    astype(int)
19
20 # Cyclical encodings for hour-of-day
21 events["hour_sin"] = npy.sin(2*npy.pi*events["hour"]/24)
22 events["hour_cos"] = npy.cos(2*npy.pi*events["hour"]/24)
23
24 display(events.head())
25
26 # (Adds value) Quick weekday count plot to show extracted
    structure
27 weekday_counts = events["weekday"].value_counts().
    sort_index()
28 plt.figure(figsize=(5,3))
29 plt.bar(["Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun"],
    weekday_counts)
30 plt.title("Event counts by weekday")
31 plt.tight_layout()
32 plt.savefig("fig_events_weekday_counts.png", dpi=150,

```

```
        bbox_inches="tight")
33 plt.show()

1  import pandas as pd
2  import numpy as npy
3  import matplotlib.pyplot as plt
4
5  # Synthetic event timestamps (hourly for ~2 weeks)
6  ts = pd.date_range("2023-05-01", periods=24*14, freq="H")
7  events = pd.DataFrame({
8      "timestamp": ts,
9      "value": npy.random.default_rng(0).normal(loc=0,
10         scale=1, size=len(ts))
11 })
12 # Basic calendar features
13 events["year"] = events["timestamp"].dt.year
14 events["month"] = events["timestamp"].dt.month
15 events["day"] = events["timestamp"].dt.day
16 events["weekday"] = events["timestamp"].dt.weekday # 0=
    Mon
17 events["hour"] = events["timestamp"].dt.hour
18 events["is_weekend"] = events["weekday"].isin([5, 6]).
    astype(int)
19
20 # Cyclical encodings for hour-of-day
21 events["hour_sin"] = npy.sin(2*npy.pi*events["hour"]/24)
22 events["hour_cos"] = npy.cos(2*npy.pi*events["hour"]/24)
23
24 display(events.head())
25
26 # (Adds value) Quick weekday count plot to show extracted
    structure
27 weekday_counts = events["weekday"].value_counts().
    sort_index()
28 plt.figure(figsize=(5,3))
29 plt.bar(["Mon", "Tue", "Wed", "Thu", "Fri", "Sat", "Sun"],
    weekday_counts)
30 plt.title("Event counts by weekday")
31 plt.tight_layout()
32 plt.show()
```

timestamp

value

year

month

day

weekday

hour

is_weekend

hour_sin

hour_cos

0

2023-05-01 00:00:00

0.125730

2023

5

1

0

0

0

0.000000

1.000000

1

2023-05-01 01:00:00

-0.132105

2023

5

1

0

1

0

0.258819

0.965926

2

2023-05-01 02:00:00

0.640423

2023

5

1

0

2

0

0.500000

0.866025

3

2023-05-01 03:00:00

0.104900

2023

5

1

0

3

0

0.707107

0.707107

4

2023-05-01 04:00:00

-0.535669

2023

5

1

0

4

0

0.866025

0.500000

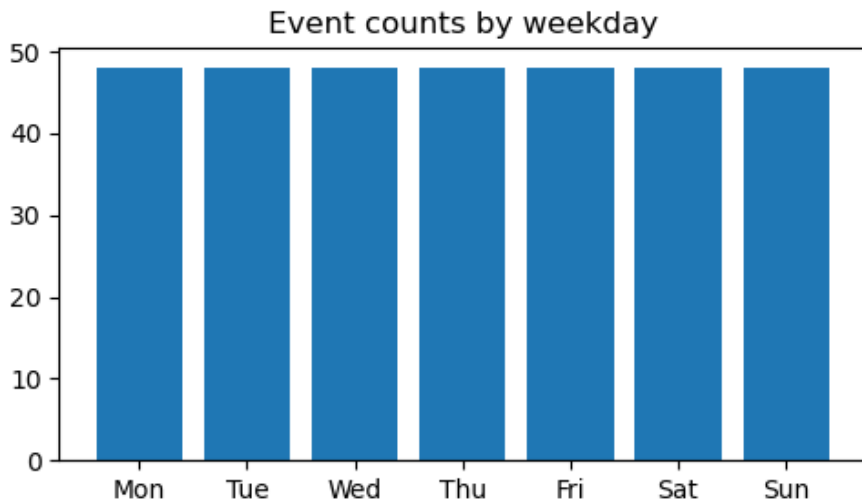


Figure 1: Image generated by the provided code.

Why cyclical features? 23 : 00 and 00 : 00 are adjacent in time but far as integers. Sine/cosine preserve circular adjacency for models that read numbers literally.

Text preprocessing basics: Bag-of-Words and TF-IDF

For short texts, two foundations:

- **Bag-of-Words (CountVectorizer):** counts token occurrences per document.
- **TF-IDF (TfidfVectorizer):** down-weights ubiquitous tokens, up-weights distinctive ones.

We'll use a tiny corpus and inspect the top-weighted terms.

```
1 from sklearn.feature_extraction.text import
    CountVectorizer, TfidfVectorizer
2 import numpy as npy
3
4 corpus = [
5     "Feature engineering turns data into signal",
6     "Good features help simple models compete",
7     "Pipelines keep preprocessing reproducible",
8     "Scaling and encoding are essential steps",
9 ]
10
11 # Bag of Words
12 cv = CountVectorizer(stop_words="english")
13 X_bow = cv.fit_transform(corpus)
14 vocab = npy.array(cv.get_feature_names_out())
15 print("BOW shape:", X_bow.shape)
16
17 # TF-IDF
18 tfidf = TfidfVectorizer(stop_words="english")
19 X_tfidf = tfidf.fit_transform(corpus)
20 terms = npy.array(tfidf.get_feature_names_out())
21
22 # Show top-5 terms by weight for each doc
23 for i, doc in enumerate(corpus):
24     row = X_tfidf.getrow(i).toarray().ravel()
25     top_idx = row.argsort()[-5:][::-1]
26     print(f"\nDoc {i}:", doc)
27     print("Top terms:", terms[top_idx].tolist(), "weights
        :", row[top_idx].round(3).tolist())
```

```
1 from sklearn.feature_extraction.text import
    CountVectorizer, TfidfVectorizer
2 import numpy as npy
3
4 corpus = [
5     "Feature engineering turns data into signal",
6     "Good features help simple models compete",
7     "Pipelines keep preprocessing reproducible",
8     "Scaling and encoding are essential steps",
9 ]
10
```

```
11 # Bag of Words
12 cv = CountVectorizer(stop_words="english")
13 X_bow = cv.fit_transform(corpus)
14 vocab = npy.array(cv.get_feature_names_out())
15 print("BOW shape:", X_bow.shape)
16
17 # TF-IDF
18 tfidf = TfidfVectorizer(stop_words="english")
19 X_tfidf = tfidf.fit_transform(corpus)
20 terms = npy.array(tfidf.get_feature_names_out())
21
22 # Show top-5 terms by weight for each doc
23 for i, doc in enumerate(corpus):
24     row = X_tfidf.getrow(i).toarray().ravel()
25     top_idx = row.argsort()[-5:][::-1]
26     print(f"\nDoc {i}:", doc)
27     print("Top terms:", terms[top_idx].tolist(), "weights
        :", row[top_idx].round(3).tolist())
```

```
1 BOW shape: (4, 18)
2
3 Doc 0: Feature engineering turns data into signal
4 Top terms: ['turns', 'signal', 'engineering', 'data', '
    feature'] weights: [0.447, 0.447, 0.447, 0.447,
    0.447]
5
6 Doc 1: Good features help simple models compete
7 Top terms: ['help', 'simple', 'features', 'compete', '
    models'] weights: [0.408, 0.408, 0.408, 0.408, 0.408]
8
9 Doc 2: Pipelines keep preprocessing reproducible
10 Top terms: ['preprocessing', 'reproducible', 'pipelines',
    'turns', 'steps'] weights: [0.577, 0.577, 0.577,
    0.0, 0.0]
11
12 Doc 3: Scaling and encoding are essential steps
13 Top terms: ['steps', 'scaling', 'encoding', 'essential',
    'turns'] weights: [0.5, 0.5, 0.5, 0.5, 0.0]
```

Guidance: For real projects, add tokenization/normalization (lowercasing, stem-

ming/lemmatization), n-grams, character models for misspellings, or domain vocabularies.

Domain-driven features (California Housing)

Domain knowledge often beats blind expansion. Examples for this dataset:

- **Bedrooms ratio:** `AveBedrms` / `AveRooms` (how many bedrooms per room).
- **Rooms per occupant:** `AveRooms` / `AveOccup` (space per person).
- **Income per occupant:** `MedInc` / `AveOccup` (ability to pay).
- **Coarse geospatial bins:** latitude/longitude buckets to approximate regions.

We'll implement a small transformer that **adds** these features, and show how to integrate it into a pipeline.

```

1  from sklearn.base import BaseEstimator, TransformerMixin
2
3  class AddHousingFeatures(BaseEstimator, TransformerMixin):
4      :
5      def __init__(self):
6          self.new_cols_ = ["bedrooms_ratio", "
7                          rooms_per_occup", "income_per_occup", "
8                          lat_bin", "lon_bin"]
9      def fit(self, X, y=None):
10         return self
11     def transform(self, X):
12         X = X.copy()
13         # Guard against divide-by-zero
14         eps = 1e-9
15         X["bedrooms_ratio"] = (X["AveBedrms"] / (X["
16                             AveRooms"] + eps)).clip(0, None)

```

```
13     X["rooms_per_occup"] = (X["AveRooms"] / (X["
14         AveOccup"] + eps)).clip(0, None)
15     X["income_per_occup"] = (X["MedInc"] / (X["
16         AveOccup"] + eps)).clip(0, None)
17     # Coarse bins (tune bin edges to your domain/
18         region)
19     X["lat_bin"] = pd.cut(X["Latitude"], bins=[32,
20         34, 36, 38, 40, 42], labels=False,
21         include_lowest=True)
22     X["lon_bin"] = pd.cut(X["Longitude"], bins=[-125,
23         -123, -121, -119, -117, -115], labels=False,
24         include_lowest=True)
25     return X
26
27 # Demo: transform a copy and preview
28 demo = AddHousingFeatures().fit_transform(Xh.copy())
29 display(demo[["AveRooms", "AveBedrms", "AveOccup", "MedInc",
30     "bedrooms_ratio", "rooms_per_occup", "income_per_occup",
31     "lat_bin", "lon_bin"]].head())
```

```
1  from sklearn.base import BaseEstimator, TransformerMixin
2
3  class AddHousingFeatures(BaseEstimator, TransformerMixin)
4      :
5      def __init__(self):
6          self.new_cols_ = ["bedrooms_ratio", "
7              rooms_per_occup", "income_per_occup", "
8              lat_bin", "lon_bin"]
9      def fit(self, X, y=None):
10          return self
11      def transform(self, X):
12          X = X.copy()
13          # Guard against divide-by-zero
14          eps = 1e-9
15          X["bedrooms_ratio"] = (X["AveBedrms"] / (X["
16              AveRooms"] + eps)).clip(0, None)
17          X["rooms_per_occup"] = (X["AveRooms"] / (X["
18              AveOccup"] + eps)).clip(0, None)
19          X["income_per_occup"] = (X["MedInc"] / (X["
20              AveOccup"] + eps)).clip(0, None)
21          # Coarse bins (tune bin edges to your domain/
```

```

    region)
16     X["lat_bin"] = pd.cut(X["Latitude"], bins=[32,
    34, 36, 38, 40, 42], labels=False,
    include_lowest=True)
17     X["lon_bin"] = pd.cut(X["Longitude"], bins=[-125,
    -123, -121, -119, -117, -115], labels=False,
    include_lowest=True)
18     return X
19
20 # Demo: transform a copy and preview
21 demo = AddHousingFeatures().fit_transform(Xh.copy())
22 display(demo[["AveRooms", "AveBedrms", "AveOccup", "MedInc",
    "bedrooms_ratio", "rooms_per_occup", "income_per_occup",
    "lat_bin", "lon_bin"]].head())

```

AveRooms

AveBedrms

AveOccup

MedInc

bedrooms_ratio

rooms_per_occup

income_per_occup

lat_bin

lon_bin

0

6.984127

1.023810

2.555556

8.3252

0.146591

2.732919

3.257687

2

1.0

1

6.238137

0.971880

2.109842

8.3014

0.155797

2.956685

3.934608

2

1.0

2

8.288136

1.073446

2.802260

7.2574

0.129516

2.957661

2.589838

2

1.0

3

5.817352

1.073059

2.547945

5.6431

0.184458

2.283154

2.214765

2

1.0

4

6.281853

1.081081

2.181467

3.8462

0.172096

2.879646

1.763125

2

1.0

Pipelines with generators: Place the feature generator at the **front** of your pipeline so downstream scalers/encoders see both original and derived columns consistently.

```
1 import numpy as npy
2 import pandas as pd
3 from sklearn.datasets import fetch_california_housing
4 from sklearn.pipeline import Pipeline
5 from sklearn.compose import ColumnTransformer
6 from sklearn.preprocessing import StandardScaler,
   OneHotEncoder
7 from sklearn.linear_model import Ridge
8 from sklearn.model_selection import cross_val_score
9
10 # 1) Carga y define X/y del *mismo* dataset (California
    Housing)
11 cal = fetch_california_housing(as_frame=True)
12 Xh = cal.frame.drop(columns=["MedHouseVal"]).copy()
13 y_h = cal.frame["MedHouseVal"].copy()
14
15 # 2) Columnas esperadas tras el generador
    AddHousingFeatures (definido previamente)
16 num_cols = [
17     "MedInc", "HouseAge", "AveRooms", "AveBedrms", "
        Population", "AveOccup", "Latitude", "Longitude",
18     "bedrooms_ratio", "rooms_per_occup", "income_per_occup"
19 ]
20 cat_cols = ["lat_bin", "lon_bin"]
21
22 # 3) OneHotEncoder: compatibilidad (>=1.2 usa
    sparse_output)
23 try:
24     ohe = OneHotEncoder(handle_unknown="ignore",
        sparse_output=False)
25 except TypeError:
26     ohe = OneHotEncoder(handle_unknown="ignore", sparse=
        False)
27
28 # 4) Preprocesado: escala numéricas y one-hot para bins
    geográficos
```

```

29 pre = ColumnTransformer(
30     transformers=[
31         ("num", StandardScaler(), num_cols),
32         ("cat", ohe, cat_cols),
33     ],
34     remainder="drop",
35     verbose_feature_names_out=False
36 )
37
38 # 5) Pipeline completo con el generador de *domain
39     features*
40 pipe_domain = Pipeline([
41     ("gen", AddHousingFeatures()), # añade
42         bedrooms_ratio, rooms_per_occup, income_per_occup
43         , lat_bin, lon_bin
44     ("prep", pre),
45     ("model", Ridge(alpha=1.0, random_state=42))
46 ])
47
48 # 6) Baseline en crudo (sin features de dominio ni
49     escalado)
50 baseline = Ridge(alpha=1.0, random_state=42)
51
52 # 7) Comparativa CV (métrica: RMSE)
53 def rmse(scores):
54     return npy.sqrt(-scores.mean())
55
56 scores_raw = cross_val_score(baseline, Xh, y_h, cv=5,
57     scoring="neg_mean_squared_error", n_jobs=-1)
58 scores_domain = cross_val_score(pipe_domain, Xh, y_h, cv
59     =5, scoring="neg_mean_squared_error", n_jobs=-1)
60
61 print("RMSE (raw features): ", rmse(scores_raw))
62 print("RMSE (domain features): ", rmse(scores_domain))

```

```

1 import numpy as npy
2 import pandas as pd
3 from sklearn.datasets import fetch_california_housing
4 from sklearn.pipeline import Pipeline
5 from sklearn.compose import ColumnTransformer
6 from sklearn.preprocessing import StandardScaler,

```

```
OneHotEncoder
7  from sklearn.linear_model import Ridge
8  from sklearn.model_selection import cross_val_score
9
10 # 1) Carga y define X/y del *mismo* dataset (California
    Housing)
11 cal = fetch_california_housing(as_frame=True)
12 Xh = cal.frame.drop(columns=["MedHouseVal"]).copy()
13 y_h = cal.frame["MedHouseVal"].copy()
14
15 # 2) Columnas esperadas tras el generador
    AddHousingFeatures (definido previamente)
16 num_cols = [
17     "MedInc", "HouseAge", "AveRooms", "AveBedrms", "
        Population", "AveOccup", "Latitude", "Longitude",
18     "bedrooms_ratio", "rooms_per_occup", "income_per_occup"
19 ]
20 cat_cols = ["lat_bin", "lon_bin"]
21
22 # 3) OneHotEncoder: compatibilidad (>=1.2 usa
    sparse_output)
23 try:
24     ohe = OneHotEncoder(handle_unknown="ignore",
        sparse_output=False)
25 except TypeError:
26     ohe = OneHotEncoder(handle_unknown="ignore", sparse=
        False)
27
28 # 4) Preprocesado: escala numéricas y one-hot para bins
    geográficos
29 pre = ColumnTransformer(
30     transformers=[
31         ("num", StandardScaler(), num_cols),
32         ("cat", ohe, cat_cols),
33     ],
34     remainder="drop",
35     verbose_feature_names_out=False
36 )
37
38 # 5) Pipeline completo con el generador de *domain
    features*
```



```

39 pipe_domain = Pipeline([
40     ("gen", AddHousingFeatures()), # añade
        bedrooms_ratio, rooms_per_occup, income_per_occup
        , lat_bin, lon_bin
41     ("prep", pre),
42     ("model", Ridge(alpha=1.0, random_state=42))
43 ])
44
45 # 6) Baseline en crudo (sin features de dominio ni
        escalado)
46 baseline = Ridge(alpha=1.0, random_state=42)
47
48 # 7) Comparativa CV (métrica: RMSE)
49 def rmse(scores):
50     return npy.sqrt(-scores.mean())
51
52 scores_raw = cross_val_score(baseline, Xh, y_h, cv=5,
        scoring="neg_mean_squared_error", n_jobs=-1)
53 scores_domain = cross_val_score(pipe_domain, Xh, y_h, cv
        =5, scoring="neg_mean_squared_error", n_jobs=-1)
54
55 print("RMSE (raw features):      ", rmse(scores_raw))
56 print("RMSE (domain features): ", rmse(scores_domain))

```

```

1 RMSE (raw features):      0.7471805175916729
2 RMSE (domain features):  0.697073317297543

```

Reading the comparison: If RMSE drops with domain features, you've added meaningful signal. If it doesn't, revisit definitions and consider interactions or non-linear models.

Summary

- **Interactions** (`PolynomialFeatures`) let linear models capture cross-effects; control degree to avoid feature explosion.
- **Feature crosses** combine categories to capture co-occurrence patterns; en-

code after crossing.

- **Date/time** features should include both calendar fields and **cyclical encodings** for periodic variables.
 - **Text** basics: Bag-of-Words and TF-IDF provide strong baselines; enrich with n-grams and normalization for real workloads.
 - **Domain-driven** features leverage the problem's structure and often yield outsized gains with simple models.
-

Handling Missing Data

Why missing data matters (MCAR, MAR, MNAR)

Real-world datasets almost always have missing values. Ignoring them can bias estimates or break models.

- **MCAR** (Missing Completely At Random): missingness unrelated to any variable → least harmful.
- **MAR** (Missing At Random): missingness depends on observed variables (e.g., `age` missing more often for a given `class`).
- **MNAR** (Missing Not At Random): missingness depends on the unobserved value itself (hardest case).

Key rules: - **Never** impute using information from the validation/test split → use **pipelines** so imputers learn from *train* only. - Consider adding **missingness indicators** to help the model learn patterns in what is missing.

Inspecting missingness

We will look at the seaborn `titanic` dataset and quantify missing values by column.

A simple **bar chart** of missing counts is often more actionable than a heatmap.

```
1 import pandas as pd
```

```
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 titanic = sns.load_dataset("titanic").copy()
6 missing_counts = titanic.isna().sum().sort_values(
    ascending=False)
7 missing_pct = (missing_counts / len(titanic)).round(3)
8
9 miss_tbl = pd.DataFrame({"missing_count": missing_counts,
    "missing_pct": missing_pct})
10 display(miss_tbl[miss_tbl["missing_count"] > 0])
11
12 # Bar chart (adds value: shows which cols matter)
13 plt.figure(figsize=(8,4))
14 missing_counts[missing_counts>0].plot(kind="bar")
15 plt.ylabel("count of NaN")
16 plt.title("Titanic: missing values by column")
17 plt.tight_layout()
18 plt.show()
```

```
1 import pandas as pd
2 import matplotlib.pyplot as plt
3 import seaborn as sns
4
5 titanic = sns.load_dataset("titanic").copy()
6 missing_counts = titanic.isna().sum().sort_values(
    ascending=False)
7 missing_pct = (missing_counts / len(titanic)).round(3)
8
9 miss_tbl = pd.DataFrame({"missing_count": missing_counts,
    "missing_pct": missing_pct})
10 display(miss_tbl[miss_tbl["missing_count"] > 0])
11
12 # Bar chart (adds value: shows which cols matter)
13 plt.figure(figsize=(8,4))
14 missing_counts[missing_counts>0].plot(kind="bar")
15 plt.ylabel("count of NaN")
16 plt.title("Titanic: missing values by column")
17 plt.tight_layout()
18 plt.show()
```

missing_count

missing_pct

deck

688

0.772

age

177

0.199

embarked

2

0.002

embark_town

2

0.002

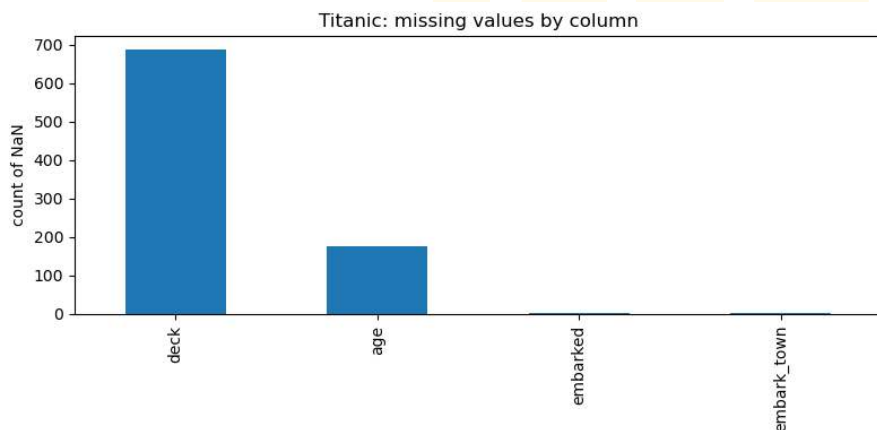


Figure 1: Image generated by the provided code.

Reading the table/plot: `age`, `embarked`, and especially `deck` often have missing values. We'll focus on features used in earlier chapters (`sex`, `class`, `embarked`, `age`, `fare`) to keep the pipeline consistent.

Simple imputation (mean/median/most-frequent/constant)

SimpleImputer is a strong baseline:

- **Numeric:** `median` is robust to outliers; `mean` if distribution is symmetric.
- **Categorical:** `most_frequent` or a **constant** like `"Unknown"`.

We will build a pipeline for Titanic classification with **imputation + encoding + logistic regression**.

```
1 from sklearn.compose import ColumnTransformer
2 from sklearn.preprocessing import OneHotEncoder
```

```

3  from sklearn.impute import SimpleImputer
4  from sklearn.pipeline import Pipeline
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import cross_val_score
7
8  X = titanic[["sex", "class", "embarked", "age", "fare"]].
    copy()
9  y = titanic["survived"].copy()
10
11  nominal = ["sex", "embarked"]
12  ordinal = ["class"]
13  numeric = ["age", "fare"]
14
15  # OneHotEncoder version-compatibility shim
16  try:
17      ohe = OneHotEncoder(handle_unknown="ignore",
                           sparse_output=False)
18  except TypeError:
19      ohe = OneHotEncoder(handle_unknown="ignore",
                           sparse_output=False)
20
21  pre_simple = ColumnTransformer(
22      transformers=[
23          ("nominal", Pipeline([
24              ("imp", SimpleImputer(strategy="most_frequent")),
25              ("ohe", ohe),
26          ]), nominal),
27          ("ordinal", Pipeline([
28              ("imp", SimpleImputer(strategy="most_frequent")),
29              # encode 'class' as ordinal  Third < Second < First
30              ("ord", OneHotEncoder(handle_unknown="ignore",
                                     drop=None, sparse_output=False) #
                                     simple OHE for class
31              if hasattr(OneHotEncoder, "sparse_output")
32              else OneHotEncoder(handle_unknown="ignore",
                                     drop=None, sparse_output=False)),
33          ]), ordinal),
34          ("numeric", Pipeline([

```



```

35         ("imp", SimpleImputer(strategy="median")),
36     ]), numeric),
37 ],
38     verbose_feature_names_out=False
39 )
40
41 pipe_simple = Pipeline([
42     ("prep", pre_simple),
43     ("clf", LogisticRegression(max_iter=1000))
44 ])
45
46 scores = cross_val_score(pipe_simple, X, y, cv=5, scoring
47     ="accuracy", n_jobs=-1)
48 print("Accuracy with simple imputation:", scores.mean().
49     round(3))

```

```

1  from sklearn.compose import ColumnTransformer
2  from sklearn.preprocessing import OneHotEncoder
3  from sklearn.impute import SimpleImputer
4  from sklearn.pipeline import Pipeline
5  from sklearn.linear_model import LogisticRegression
6  from sklearn.model_selection import cross_val_score
7
8  X = titanic[["sex", "class", "embarked", "age", "fare"]].
9      copy()
10 y = titanic["survived"].copy()
11
12 nominal = ["sex", "embarked"]
13 ordinal = ["class"]
14 numeric = ["age", "fare"]
15
16 # OneHotEncoder version-compatibility shim
17 try:
18     ohe = OneHotEncoder(handle_unknown="ignore",
19         sparse_output=False)
20 except TypeError:
21     ohe = OneHotEncoder(handle_unknown="ignore",
22         sparse_output=False)
23
24 pre_simple = ColumnTransformer(
25     transformers=[

```

```

23         ("nominal", Pipeline([
24             ("imp", SimpleImputer(strategy="most_frequent")),
25             ("ohe", ohe),
26         ]), nominal),
27         ("ordinal", Pipeline([
28             ("imp", SimpleImputer(strategy="most_frequent")),
29             # encode 'class' as ordinal Third < Second < First
30             ("ord", OneHotEncoder(handle_unknown="ignore",
31                                     drop=None, sparse_output=False) #
32                                     simple OHE for class
33                                     if hasattr(OneHotEncoder, "sparse_output")
34                                     else OneHotEncoder(handle_unknown="ignore",
35                                                         drop=None, sparse_output=False)),
36         ]), ordinal),
37         ("numeric", Pipeline([
38             ("imp", SimpleImputer(strategy="median")),
39         ]), numeric),
40     ],
41     verbose_feature_names_out=False
42 )
43
44 pipe_simple = Pipeline([
45     ("prep", pre_simple),
46     ("clf", LogisticRegression(max_iter=1000))
47 ])
48
49 scores = cross_val_score(pipe_simple, X, y, cv=5, scoring
50                           ="accuracy", n_jobs=-1)
51
52 print("Accuracy with simple imputation:", scores.mean().
53       round(3))

```

```
1 Accuracy with simple imputation: 0.789
```

Adding missingness indicators

Sometimes **whether a value is missing is predictive** (e.g., people without recorded `age` differ systematically).

We can append indicator columns. The simplest path is `SimpleImputer(add_indicator=True)` (if supported) or `MissingIndicator`.

We'll attempt `add_indicator=True` and fall back to `MissingIndicator` if needed.

```
1  from sklearn.impute import MissingIndicator
2
3  def nominal_block_with_indicator():
4      try:
5          # SimpleImputer can append indicators in recent
5              scikit-learn
6          return Pipeline([
7              ("imp", SimpleImputer(strategy="most_frequent",
7                  add_indicator=True)),
8              ("ohe", ohe),
9          ])
10     except TypeError:
11         # Fallback: separate indicator and stack via
11             ColumnTransformer
12         return Pipeline([
13             ("imp", SimpleImputer(strategy="most_frequent",
13                 add_indicator=True)),
14             ("ohe", ohe),
15         ])
16
17  def numeric_block_with_indicator():
18      try:
19          return Pipeline([
20              ("imp", SimpleImputer(strategy="median",
20                  add_indicator=True)),
21          ])
22     except TypeError:
23         return Pipeline([
24             ("imp", SimpleImputer(strategy="median")),
```

```

25         ])
26
27     pre_with_ind = ColumnTransformer(
28         transformers=[
29             ("nominal", nominal_block_with_indicator(),
30              nominal),
31             ("ordinal", Pipeline([
32                 ("imp", SimpleImputer(strategy="most_frequent
33                 ")),
34                 ("ohe", ohe),
35             ]), ordinal),
36             ("numeric", numeric_block_with_indicator(),
37              numeric),
38         ],
39         verbose_feature_names_out=False
40     )
41
42     pipe_with_ind = Pipeline([
43         ("prep", pre_with_ind),
44         ("clf", LogisticRegression(max_iter=1000))
45     ])
46
47     scores_ind = cross_val_score(pipe_with_ind, X, y, cv=5,
48                                  scoring="accuracy", n_jobs=-1)
49     print("Accuracy with indicators:", scores_ind.mean().
50           round(3))

```

```

1  from sklearn.impute import MissingIndicator
2
3  def nominal_block_with_indicator():
4      try:
5          # SimpleImputer can append indicators in recent
6          # scikit-learn
7          return Pipeline([
8              ("imp", SimpleImputer(strategy="most_frequent
9              ", add_indicator=True)),
10             ("ohe", ohe),
11         ])
12     except TypeError:
13         # Fallback: separate indicator and stack via
14         # ColumnTransformer

```

```

12         return Pipeline([
13             ("imp", SimpleImputer(strategy="most_frequent")),
14             ("ohe", ohe),
15         ])
16
17 def numeric_block_with_indicator():
18     try:
19         return Pipeline([
20             ("imp", SimpleImputer(strategy="median",
21                                   add_indicator=True)),
22         ])
23     except TypeError:
24         return Pipeline([
25             ("imp", SimpleImputer(strategy="median")),
26         ])
27
28 pre_with_ind = ColumnTransformer(
29     transformers=[
30         ("nominal", nominal_block_with_indicator(),
31          nominal),
32         ("ordinal", Pipeline([
33             ("imp", SimpleImputer(strategy="most_frequent")),
34             ("ohe", ohe),
35         ]), ordinal),
36         ("numeric", numeric_block_with_indicator(),
37          numeric),
38     ],
39     verbose_feature_names_out=False
40 )
41
42 pipe_with_ind = Pipeline([
43     ("prep", pre_with_ind),
44     ("clf", LogisticRegression(max_iter=1000))
45 ])
46
47 scores_ind = cross_val_score(pipe_with_ind, X, y, cv=5,
48                               scoring="accuracy", n_jobs=-1)
49 print("Accuracy with indicators:", scores_ind.mean().
50       round(3))

```

```
1 Accuracy with indicators: 0.799
```

Interpreting the result: If accuracy improves, the model is benefiting from the information contained in the pattern of missingness.

KNNImputer (neighbors-based imputation)

KNNImputer imputes a value from **similar rows** (nearest neighbors).

Pros: uses multivariate structure. Cons: slower, sensitive to scale and sparse categories.

We will:

- One-hot encode categoricals first (so distance is meaningful),
- Scale numeric columns (optional but recommended),
- Then apply KNNImputer on the **combined** feature space.

```
1 from sklearn.impute import KNNImputer
2 from sklearn.preprocessing import StandardScaler
3
4 pre_knn = ColumnTransformer(
5     transformers=[
6         ("nominal", Pipeline([
7             ("imp", SimpleImputer(strategy="most_frequent")),
8             ("ohe", ohe),
9         ]), nominal),
10        ("ordinal", Pipeline([
11            ("imp", SimpleImputer(strategy="most_frequent")),
12            ("ohe", ohe),
13        ]), ordinal),
14        ("numeric", Pipeline([
15            ("imp", SimpleImputer(strategy="median")),
16            ("sc", StandardScaler()),
17        ]), numeric),
```

```

18     ],
19     verbose_feature_names_out=False
20 )
21
22 pipe_knn = Pipeline([
23     ("prep", pre_knn),
24     ("knnimp", KNNImputer(n_neighbors=5, weights="
25         distance")),
26     ("clf", LogisticRegression(max_iter=1000))
27 ])
28 scores_knn = cross_val_score(pipe_knn, X, y, cv=5,
29     scoring="accuracy", n_jobs=-1)
30 print("Accuracy with KNNImputer:", scores_knn.mean().
31     round(3))

```

```

1 from sklearn.impute import KNNImputer
2 from sklearn.preprocessing import StandardScaler
3
4 pre_knn = ColumnTransformer(
5     transformers=[
6         ("nominal", Pipeline([
7             ("imp", SimpleImputer(strategy="most_frequent
8                 ")),
9             ("ohe", OneHotEncoder()),
10            ]), nominal),
11        ("ordinal", Pipeline([
12            ("imp", SimpleImputer(strategy="most_frequent
13                ")),
14            ("ohe", OneHotEncoder()),
15            ]), ordinal),
16        ("numeric", Pipeline([
17            ("imp", SimpleImputer(strategy="median")),
18            ("sc", StandardScaler()),
19            ]), numeric),
20    ],
21    verbose_feature_names_out=False
22 )
23
24 pipe_knn = Pipeline([
25     ("prep", pre_knn),

```

```
24     ("knnimp", KNNImputer(n_neighbors=5, weights="
25         distance")),
26 ] )
27
28 scores_knn = cross_val_score(pipe_knn, X, y, cv=5,
29     scoring="accuracy", n_jobs=-1)
29 print("Accuracy with KNNImputer:", scores_knn.mean().
    round(3))
```

```
1 Accuracy with KNNImputer: 0.789
```

Notes:

- KNNImputer occurs **after** encoding/scaling so distances reflect all features.
- Tune `n_neighbors`; `weights="distance"` can help when neighbors vary in similarity.

IterativeImputer (multivariate, model-based)

IterativeImputer models each feature with missing values as a function of other features (round-robin).

Pros: powerful when relationships are approximately linear/non-linear depending on estimator.

Cons: heavier, more hyperparameters.

We'll use a linear estimator (default BayesianRidge) for speed and stability.

```
1 # Enable IterativeImputer if needed in older versions
2 try:
3     from sklearn.impute import IterativeImputer
4 except Exception:
5     from sklearn.experimental import
6         enable_iterative_imputer # noqa: F401
7     from sklearn.impute import IterativeImputer
```



```

8  pre_iter = pre_knn # reuse encoding/scaling block from
   above
9
10 pipe_iter = Pipeline([
11     ("prep", pre_iter),
12     ("iterimp", IterativeImputer(max_iter=10,
13     random_state=42, sample_posterior=False)),
13     ("clf", LogisticRegression(max_iter=1000))
14 ])
15
16 scores_iter = cross_val_score(pipe_iter, X, y, cv=5,
17     scoring="accuracy", n_jobs=-1)
17 print("Accuracy with IterativeImputer:", scores_iter.mean
18     ().round(3))

```

```

1  # Enable IterativeImputer if needed in older versions
2  try:
3      from sklearn.impute import IterativeImputer
4  except Exception:
5      from sklearn.experimental import
6      enable_iterative_imputer # noqa: F401
7      from sklearn.impute import IterativeImputer
8
9  pre_iter = pre_knn # reuse encoding/scaling block from
   above
10
11 pipe_iter = Pipeline([
12     ("prep", pre_iter),
13     ("iterimp", IterativeImputer(max_iter=10,
14     random_state=42, sample_posterior=False)),
15     ("clf", LogisticRegression(max_iter=1000))
16 ])
17
18 scores_iter = cross_val_score(pipe_iter, X, y, cv=5,
19     scoring="accuracy", n_jobs=-1)
19 print("Accuracy with IterativeImputer:", scores_iter.mean
20     ().round(3))

```

```
1 Accuracy with IterativeImputer: 0.789
```

Visual check: distribution before/after imputation (adds value)

We will compare the **age** distribution before imputation vs after simple median imputation.

This helps verify that imputation did not create unrealistic spikes or distortions.

```
1 # Extract 'age' before and after a simple median
  imputation for plotting
2 age_raw = X["age"]
3
4 median_imp = SimpleImputer(strategy="median")
5 age_imp = pd.Series(median_imp.fit_transform(X[["age"]]).
  ravel(), name="age_imputed")
6
7 plt.figure(figsize=(8,4))
8 plt.hist(age_raw.dropna(), bins=20, alpha=0.6, label="age
  (raw)")
9 plt.hist(age_imp, bins=20, alpha=0.6, label="age (median-
  imputed)")
10 plt.legend()
11 plt.title("Age distribution: raw vs median-imputed")
12 plt.tight_layout()
13 plt.show()
```

```
1 # Extract 'age' before and after a simple median
  imputation for plotting
2 age_raw = X["age"]
3
4 median_imp = SimpleImputer(strategy="median")
5 age_imp = pd.Series(median_imp.fit_transform(X[["age"]]).
  ravel(), name="age_imputed")
6
7 plt.figure(figsize=(8,4))
8 plt.hist(age_raw.dropna(), bins=20, alpha=0.6, label="age
  (raw)")
9 plt.hist(age_imp, bins=20, alpha=0.6, label="age (median-
  imputed)")
```

```
10 plt.legend()
11 plt.title("Age distribution: raw vs median-imputed")
12 plt.tight_layout()
13 plt.show()
```

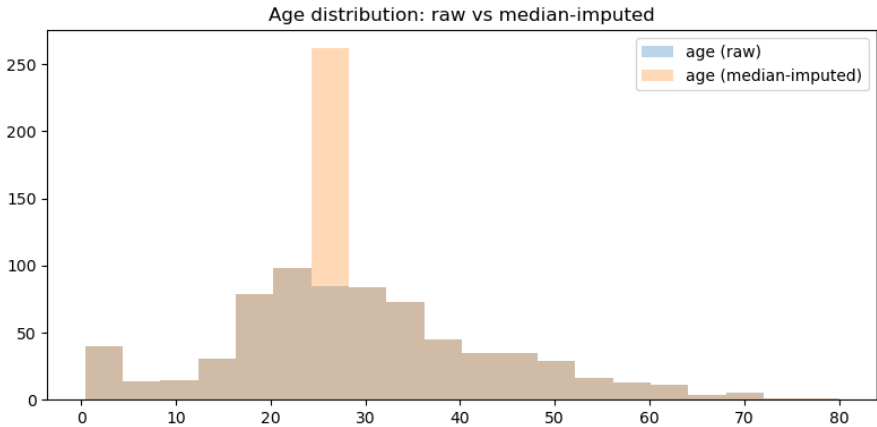


Figure 2: Image generated by the provided code.

Reading the plot: Median imputation fills gaps without changing the central tendency, but can **inflate** the median bin. Advanced imputers (KNN/Iterative) can preserve shape better—validate with metrics.

Leakage-safe evaluation: quick comparison

Let's summarize the **CV accuracy** for the different strategies we tried. Because each imputer lives *inside* the pipeline, there is **no leakage** from validation folds.

```
1 import pandas as pd
2
```

```
3 acc_table = pd.DataFrame({
4     "SimpleImputer": [scores.mean() for scores in [scores
5         ]][0],
6     "WithIndicators": [scores_ind.mean() for scores_ind
7         in [scores_ind]][0],
8     "KNNImputer": [scores_knn.mean() for scores_knn in [
9         scores_knn]][0],
10    "IterativeImputer": [scores_iter.mean() for
11        scores_iter in [scores_iter]][0],
12 }, index=["CV Accuracy"]).T.sort_values(by="CV Accuracy",
13     ascending=False)
14 display(acc_table.round(3))
```

```
1 import pandas as pd
2
3 acc_table = pd.DataFrame({
4     "SimpleImputer": [scores.mean() for scores in [scores
5         ]][0],
6     "WithIndicators": [scores_ind.mean() for scores_ind
7         in [scores_ind]][0],
8     "KNNImputer": [scores_knn.mean() for scores_knn in [
9         scores_knn]][0],
10    "IterativeImputer": [scores_iter.mean() for
11        scores_iter in [scores_iter]][0],
12 }, index=["CV Accuracy"]).T.sort_values(by="CV Accuracy",
13     ascending=False)
14 display(acc_table.round(3))
```

CV Accuracy

WithIndicators

0.799

SimpleImputer

0.789

KNNImputer

0.789

IterativeImputer

0.789

Summary

- Start with **SimpleImputer** (median for numeric, most-frequent for categorical). It's robust and fast.
 - If performance suffers and you suspect structure in missingness, try **missingness indicators**.
 - Consider **KNNImputer** or **IterativeImputer** when variables have strong relationships and simple strategies underperform.
 - Keep imputers **inside pipelines** to avoid leakage. Tune imputer hyperparameters with cross-validation as you would any model.
 - For **MNAR**, modeling missingness explicitly (two-stage models, selection models) or collecting more data may be necessary.
-

