

Search...

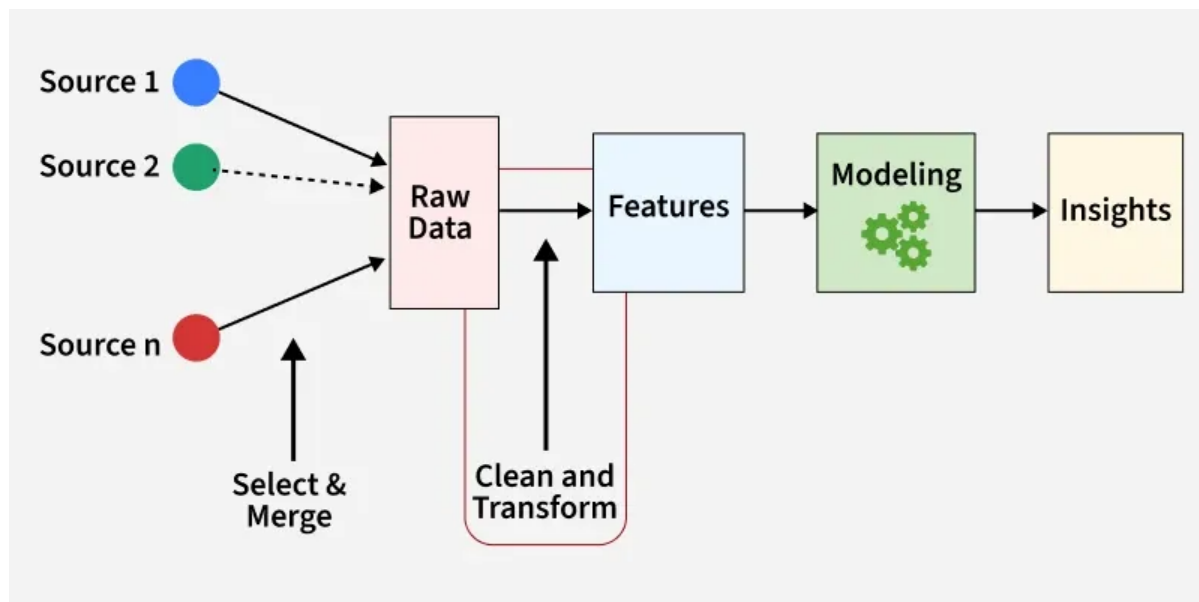
[Sign In](#)[Python for Machine Learning](#)[Machine Learning with R](#)[Machine Learning Algorit](#)[Sign In](#)

What is Feature Engineering?

Last Updated : 29 Aug, 2025

Feature engineering is the process of turning raw data into useful features that help improve the performance of machine learning models. It includes choosing, creating and adjusting data attributes to make the model's predictions more accurate. The goal is to make the model better by providing relevant and easy-to-understand information.

A feature or attribute is a measurable property of data that is used as input for machine learning algorithms. Features can be numerical, categorical or text-based representing essential data aspects which are relevant to the problem. For example in housing price prediction, features might include the number of bedrooms, location and property age.



Feature Engineering Architecture

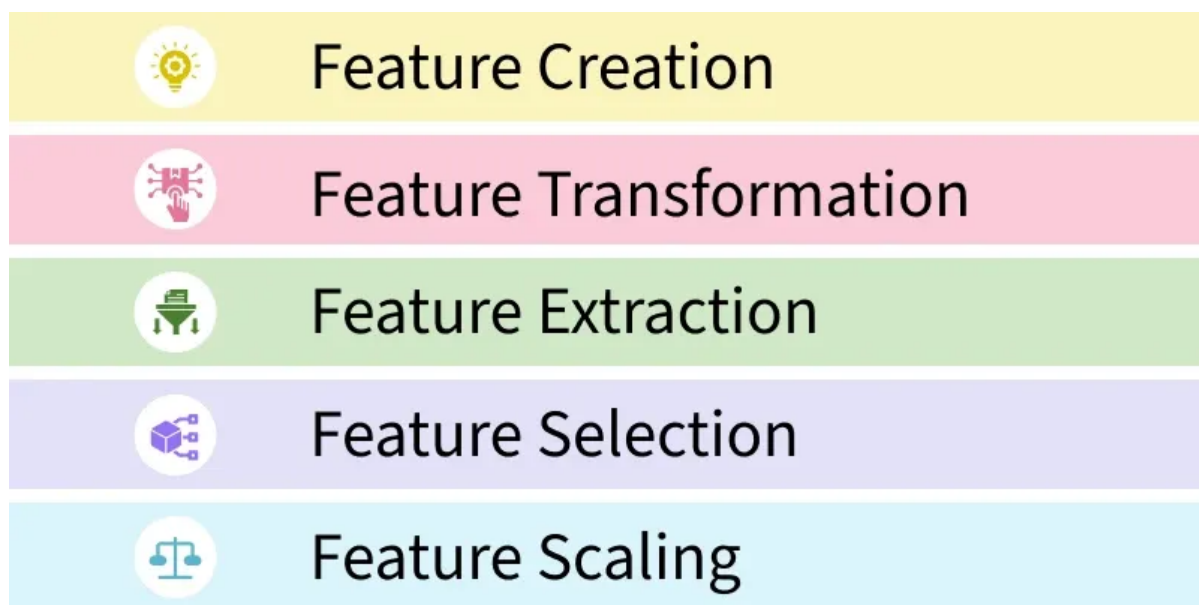
Importance of Feature Engineering

Feature engineering can significantly influence model performance. By refining features, we can:

- **Improve accuracy:** Choosing the right features helps the model learn better, leading to more accurate predictions.
- **Reduce overfitting:** Using fewer, more important features helps the model avoid memorizing the data and perform better on new data.
- **Boost interpretability:** Well-chosen features make it easier to understand how the model makes its predictions.
- **Enhance efficiency:** Focusing on key features speeds up the model's training and prediction process, saving time and resources.

Processes Involved in Feature Engineering

Lets see various features involved in feature engineering:



Processes involved in Feature Engineering

1. Feature Creation: Feature creation involves generating new features from domain knowledge or by observing patterns in the data. It can be:

1. **Domain-specific:** Created based on industry knowledge like business rules.
2. **Data-driven:** Derived by recognizing patterns in data.

3. **Synthetic:** Formed by combining existing features.

2. Feature Transformation: Transformation adjusts features to improve model learning:

1. **Normalization & Scaling:** Adjust the range of features for consistency.
2. **Encoding:** Converts categorical data to numerical form i.e one-hot encoding.
3. **Mathematical transformations:** Like logarithmic transformations for skewed data.

3. Feature Extraction: Extracting meaningful features can reduce dimensionality and improve model accuracy:

- **Dimensionality reduction:** Techniques like PCA reduce features while preserving important information.
- **Aggregation & Combination:** Summing or averaging features to simplify the model.

4. Feature Selection: Feature selection involves choosing a subset of relevant features to use:

- **Filter methods:** Based on statistical measures like correlation.
- **Wrapper methods:** Select based on model performance.
- **Embedded methods:** Feature selection integrated within model training.

5. Feature Scaling: Scaling ensures that all features contribute equally to the model:

- **Min-Max scaling:** Rescales values to a fixed range like 0 to 1.
- **Standard scaling:** Normalizes to have a mean of 0 and variance of 1.

Steps in Feature Engineering

Feature engineering can vary depending on the specific problem but the general steps are:

1. **Data Cleaning:** Identify and correct errors or inconsistencies in the dataset to ensure data quality and reliability.
2. **Data Transformation:** Transform raw data into a format suitable for modeling including scaling, normalization and encoding.
3. **Feature Extraction:** Create new features by combining or deriving information from existing ones to provide more meaningful input to the model.
4. **Feature Selection:** Choose the most relevant features for the model using techniques like correlation analysis, mutual information and stepwise regression.
5. **Feature Iteration:** Continuously refine features based on model performance by adding, removing or modifying features for improvement.

Common Techniques in Feature Engineering

1. **One-Hot Encoding:** [One-Hot Encoding](#) converts categorical variables into binary indicators, allowing them to be used by machine learning models.

```
import pandas as pd

data = {'Color': ['Red', 'Blue', 'Green', 'Blue']}
df = pd.DataFrame(data)

df_encoded = pd.get_dummies(df, columns=['Color'], prefix='Color')

print(df_encoded)
```

Output

	Color_Blue	Color_Green	Color_Red
0	False	False	True

1	True	False	False
2	False	True	False
3	True	False	False

2. Binning: [Binning](#) transforms continuous variables into discrete bins, making them categorical for easier analysis.

```
import pandas as pd

data = {'Age': [23, 45, 18, 34, 67, 50, 21]}
df = pd.DataFrame(data)

bins = [0, 20, 40, 60, 100]
labels = ['0-20', '21-40', '41-60', '61+']

df['Age_Group'] = pd.cut(df['Age'], bins=bins, labels=labels,
right=False)

print(df)
```

Output

	Age	Age_Group
0	23	21-40
1	45	41-60
2	18	0-20
3	34	21-40
4	67	61+
5	50	41-60
6	21	21-40

3. Text Data Preprocessing: Involves removing [stop-words](#), [stemming](#) and [vectorizing](#) text data to prepare it for machine learning models.

```
import nltk
from nltk.corpus import stopwords
from nltk.stem import PorterStemmer
from sklearn.feature_extraction.text import CountVectorizer

texts = ["This is a sample sentence.", "Text data preprocessing is
important."]

stop_words = set(stopwords.words('english'))
```

```

stemmer = PorterStemmer()
vectorizer = CountVectorizer()

def preprocess_text(text):
    words = text.split()
    words = [stemmer.stem(word)
              for word in words if word.lower() not in stop_words]
    return " ".join(words)

cleaned_texts = [preprocess_text(text) for text in texts]

X = vectorizer.fit_transform(cleaned_texts)

print("Cleaned Texts:", cleaned_texts)
print("Vectorized Text:", X.toarray())

```

Output:

```

➡ Cleaned Texts: ['sampl sentence.', 'text data preprocess important.']
  Vectorized Text: [[0 0 0 1 1 0]
                    [1 1 1 0 0 1]]

```

4. Feature Splitting: Divides a single feature into multiple sub-features, uncovering valuable insights and improving model performance.

```

import pandas as pd

data = {'Full_Address': [
    '123 Elm St, Springfield, 12345', '456 Oak Rd, Shelbyville,
    67890']}
df = pd.DataFrame(data)

df[['Street', 'City', 'Zipcode']] = df['Full_Address'].str.extract(
    r'([0-9]+\s[\w\s]+),\s([\w\s]+),\s(\d+)')

print(df)

```

Output

	Full_Address	Street	City
Zipcode			
0	123 Elm St, Springfield, 12345	123 Elm St	Springfield
	12345		

1 456 Oak Rd, Shelbyville, 67890 456 Oak Rd Shelbyville
67890...

Tools for Feature Engineering

There are several tools available for feature engineering. Here are some popular ones:

- **Featuretools**: Automates feature engineering by extracting and transforming features from structured data. It integrates well with libraries like pandas and scikit-learn making it easy to create complex features without extensive coding.
- **TPOT**: Uses genetic algorithms to optimize machine learning pipelines, automating feature selection and model optimization. It visualizes the entire process, helping you identify the best combination of features and algorithms.
- **DataRobot**: Automates machine learning workflows including feature engineering, model selection and optimization. It supports time-dependent and text data and offers collaborative tools for teams to efficiently work on projects.
- **Alteryx**: Offers a visual interface for building data workflows, simplifying feature extraction, transformation and cleaning. It integrates with popular data sources and its drag-and-drop interface makes it accessible for non-programmers.
- **H2O.ai**: Provides both automated and manual feature engineering tools for a variety of data types. It includes features for scaling, imputation and encoding and offers interactive visualizations to better understand model results.

[Comment](#)[More info](#)