

**PACE in machine learning: The plan and analysis stages**

- Value: Welcome to week 2 4 min
- Value: PACE in machine learning 4 min
- Value: Plan for a machine learning project 2 min
- Reading: More about planning a machine learning project 20 min
- Value: Launch: Overview, strategies and learn from your mistakes 7 min
- Value: Prepare data for a machine learning model 4 min
- Value: Introduction to feature engineering 3 min
- Reading: Explore feature engineering 10 min
- Value: Select features that come with meaningful data 4 min
- Reading: More about imbalanced datasets 20 min
- Reading: Following instructions: Feature engineering with Python 25 min
- Lab: Automated following along guide: Feature engineering with Python 25 min
- Value: Feature engineering and data science 1 min
- Lab: Deploy feature feature engineering 10 min
- Lab: Deploy: Perform feature engineering 20 min
- Practice Quiz: Test your knowledge: PACE in machine learning: The plan and analysis stages 5 questions

**PACE in machine learning: The construct and analysis stages**

**Review: Workflow for building complex models**

## Explore feature engineering

In this reading, you will learn more about what happens in the Analysis stage of PACE – namely, feature engineering. The meaning of the term “feature engineering” can vary broadly, but in this course it includes feature selection, feature transformation, and feature extraction. You will come to understand more about the considerations and process of adjusting your predictive variables to improve model performance.

### Feature Engineering

When building machine learning models, your model is only ever going to be as good as your data. Sometimes, the data you have will not be predictive of your target variable. For example, it’s unlikely that you can build a good model that predicts whether it will rain on Saturday based on stock market data. In this case, it might seem obvious, but when you’re building a model, you’ll often have features that *possibly* could be predictive of your target, but in fact are not. Other times, your model’s features might contain a predictive signal for your model, but this signal can be strengthened if you manipulate the feature in a way that makes it more detectable by the model.

**Feature engineering** is the process of using practical, statistical, and data science knowledge to select, transform, or extract characteristics, properties, and attributes from data. In this reading, you will learn more about these processes, when and why to use them, and what good feature engineering can do for your model.

#### Feature Selection

Feature selection is the process of picking variables from a dataset that will be used as predictive variables for your model. With very large datasets, there are classes of features that are not useful for your observations in the data. Taking all of the features in a dataset often doesn’t give any performance boost. In fact, it may actually hurt performance by adding complexity and noise to the model. Therefore, choosing the features to use for the model is an important part of the model development process.

Generally, there are three types of features:

1. **Predictive Features** that by themselves contain information useful to predict the target
2. **Interactive Features** that are not useful by themselves to predict the target variable, but become predictive in conjunction with other features
3. **Irrelevant Features** that don’t contain any useful information to predict the target

You want predictive features, but a predictive feature can also be a redundant feature. **Redundant features** are highly correlated with other features and therefore do not provide the model with any new information. For example, the steps you took in a day may be highly correlated with the calories you burned. The goal of feature selection is to find the predictive and interactive features and exclude redundant and irrelevant features.



The feature selection process typically occurs at multiple stages of the PACE workflow. The first place it occurs is during the Plan phase. Once you have defined your problem and decided on target variables to predict, you need to find features. Keep in mind that datasets are not always prepackaged in nice little tables ready to model. Data professionals can spend days, weeks, or even months acquiring and assembling features from many different sources.



Feature selection can happen once more during the Analysis phase. Once you do an exploratory data analysis, it might become clear that some of the features you included are not suitable for modeling. This could be for a number of reasons. For example, you might find that a feature has too many missing or clearly erroneous values, or perhaps it’s highly correlated with another feature and must be dropped as a not to violate the assumptions of your model. It’s also common that the feature is some kind of metadata, such as an ID number with no inherent meaning. Whenever the case may be, you might want to drop these features.



During the Construct phase, when you are building models, the process of improving your model might include more feature selection. At this point, the objective usually is to find the smallest set of predictive features that will result in good overall model performance. In fact, data professionals will often begin their model selection with only one score, but also on model simplicity and explainability. A model with an  $R^2$  of 0.52 and 10 features might get selected over a model with an  $R^2$  of 0.48 and 100 features. Models with fewer features are simpler, and simpler models are generally more stable and easier to understand.

When data professionals perform feature selection during the Construct phase, they typically use statistical methodologies to determine which features to keep and which to drop. It could be as simple as ranking the model’s feature importance and keeping only the top 4% of them. Another way of doing it is to keep the top features that account for 90% of the model’s predictive signal. There are many different ways of performing feature selection, but they all seek to keep the predictive features and exclude the non-predictive features.

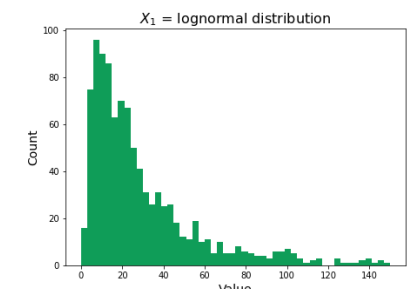
### Feature Transformation

Feature transformation is a process where you take features that already exist in the dataset, and alter them so that they’re better suited to be used for training the model. Data professionals usually perform feature transformation during the Construct phase, after they’ve analyzed the data and made decisions about how to transform it based on what they’ve learned.

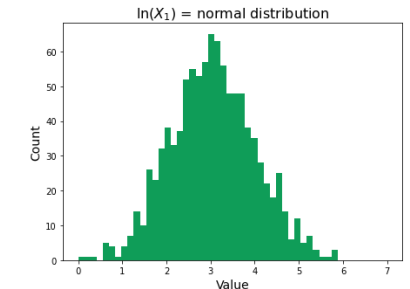
#### Log-normalization

There are various types of transformations that might be required for any given model. For example, some models do not handle continuous variables with skewed distributions very well. As a solution, you can take the log of a skewed feature, reducing the skew and making the data better for modeling. This is known as **log-normalization**.

For instance, suppose you have a feature  $X$  whose histogram demonstrated the following distribution:



This is known as a **log-normal distribution**. A log-normal distribution is a continuous distribution whose logarithm is normally distributed. In this case, the distribution shows right, but if you transform the feature by taking its natural log, it normalizes the distribution:



Normalizing a feature’s distribution is often better for training a model, and you can later verify whether or not taking the log has helped by analyzing the model’s performance.

#### Scaling

Another kind of feature transformation is **scaling**. Scaling is when you adjust the range of a feature’s values by applying a normalization function to them. Scaling helps prevent features with very large values from having more influence over a model compared to features with smaller values, but which may be equally important as predictors.

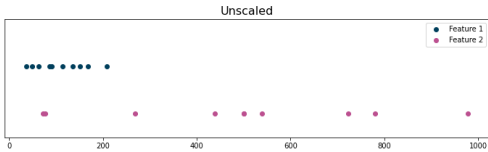
There are many scaling methodologies available. Some of the most common include:

#### Normalization

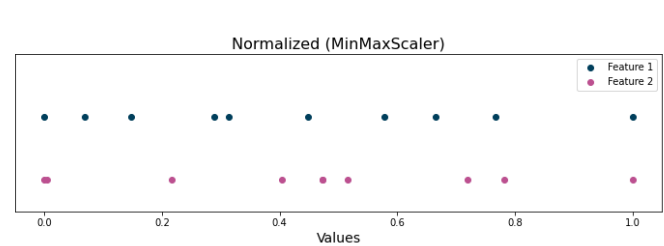
**Normalization** (e.g., **MinMaxScaler** in scikit-learn) transforms data to rescale each value to fall within the range [0, 1]. When applied to features, the feature’s minimum value becomes zero and its maximum value becomes one. All other values scale to somewhere between them. The formula for this transformation is:

$$P_{\text{normalized}} = \frac{P_{\text{value}} - P_{\text{min}}}{P_{\text{max}} - P_{\text{min}}}$$

For example, suppose you have feature 1, whose values range from 38 to 205, and feature 2, whose values range from 72 to 978.



It is apparent that these features are on different scales from one another. Features with higher magnitudes of scale will be more influential in some machine learning algorithms, like K-means, where Euclidean distances between data points are calculated with the absolute value of the feature. On large feature values have more effect, compared to small feature values. By min-max scaling/normalizing each feature, they are both reduced to the same range.

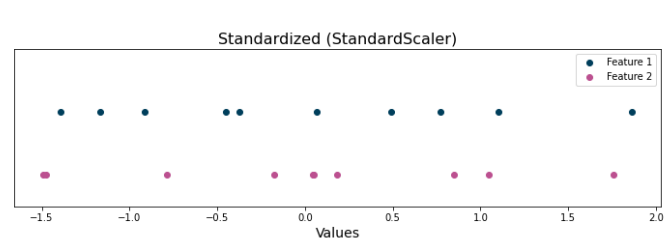


#### Standardization

Another type of scaling is called **standardization** (e.g., **StandardScaler** in scikit-learn). Standardization transforms each value within a feature so they collectively have a mean of zero and a standard deviation of one. To do this, for each value, subtract the mean of the feature and divide by the feature’s standard deviation:

$$P_{\text{standardized}} = \frac{P_{\text{value}} - \mu_{\text{feature}}}{\sigma_{\text{feature}}}$$

This method is useful because it centers the feature values at zero, which is useful for some machine learning algorithms. It also prevents outliers, since it does not place a hard cap on the range of possible values, hence it’s the same data from above after applying standardization:



Notice that the points are spatially distributed in a way that is very similar to the result of normalizing, but the values and scales are different. In this case, the values now range from -1.49 to 1.76.

#### Encoding

Another form of feature transformation is known as **encoding**. Encoding is the process of converting categorical data to numerical data. Consider the bank churn dataset. The original data has a feature called “Geography” whose values represent each customer’s country of residence—France, Germany, or Spain. Most machine learning methodologies cannot process meaning from strings. Encoding transforms the string to numbers that can be interpreted mathematically.

The “Geography” column contains nominal values, or values that don’t have an inherent order or ranking. As such, the feature would typically be encoded into binary. This process requires that a column be added to represent any possible class contained within the feature.

Geography	is France	is Germany	is Spain
France	1	0	0
Germany	0	1	0
Spain	0	0	1
France	1	0	0

Tools commonly used to do this include `pandas.get_dummies()` and `OneHotEncoder()`. Other methods drop one of the columns to avoid having redundant information in the dataset. Note that information isn’t lost by doing this. If you have this:

Customer	is France	is Germany
Antonio Garcia	0	0

...then you know this customer is from Spain!

Keep in mind that some features may be inferred to be numerical by Python or other frameworks but still represent a category. For example, suppose you had a dataset with people assigned to different industry groups: 1, 2, and 3.

Name	Group
Rachel Stein	1
Ahmed Abadi	2
Sid Awwy	3
Ha-ri Choi	1

The “Group” column might be encoded as type `int`, but the number is really only representative of a category. Group 3 isn’t two times “greater than” group 1. The proper codified or really best labeled with values. In this case, you could first convert the column to type `string`, and then encode the strings as binary columns. This is a problem that can be solved upstream at the stage of data generation: categorize features like a group should be recorded using a number.

A different kind of encoding can be used for features that contain discrete or ordinal values. This is called **ordinal encoding**. It is used when the values do contain inherent order or ranking. For instance, consider a “Temperature” column that has values of cold, warm, and hot. In this case, ordinal encoding could rescale these classes to 1, 2, and 3.

Temperature	Temperature (Ordinal encoding)
cold	0
warm	1
hot	2

This method returns the order or ranking of the classes relative to one another.

#### Feature Extraction

Feature extraction involves producing new features from existing ones, with the goal of finding features that deliver more predictive power to your model. While there is some overlap between extraction and transformation colloquially, the main difference is that a new feature is created from one or more other features rather than simply changing one that already exists.

Consider a feature called “Date of Last Purchase,” which contains information about when a customer last purchased something from the company instead of giving the model raw data, a new feature can be extracted called “Days Since Last Purchase.” This could tell the model how long it has been since a customer has bought something from the company, going right into the likelihood that they’ll buy something again in the future. Suppose that today’s date is May 20th, something a new feature could look something like this:

Date of Last Purchase	Days Since Last Purchase
May 17th	13
May 20th	1
May 19th	16
May 21st	9

Features can also be extracted from multiple variables. For example, consider modeling if a customer will return to buy something else. In the data, there are two variables: “Days Since Last Purchase” and “Price of Last Purchase.” A new variable could be created from these by dividing the price by the number of days since the last purchase, creating a new variable altogether:

Days Since Last Purchase	Price of Last Purchase	Dollars Per Day Since Last Purchase
13	\$65	\$5.34
1	\$15	\$15.00
16	\$5	\$0.49
9	\$43	\$4.78

Sometimes, the features that you are able to generate through extraction can offer the greatest performance boosts to your model. It can be a trial and error process, but finding good features from the raw data is what will make a model stand out in industry.

### Key Takeaways

- Analyzing the features in a dataset is essential to creating a model that will produce valuable results
- Feature Selection is the process of dropping noisy and all-redundant or uncorrelated features from the dataset
- Feature Transformation is the process of refitting features into a form where they’re better for training the model
- Feature Extraction is the process of creating brand new features from other features that already exist in the dataset

#### Resources for more information

- [MinMaxScaler documentation](#) (L7 scikit-learn implementation of MinMaxScaler normalization)
- [StandardScaler documentation](#) (L7 scikit-learn implementation of StandardScaler standardization)

Mark as completed

Like Dislike Report an issue