

1. What is Feature Engineering

1. Learning What? :

- determine which features are the most important with mutual information.
- invent new features in several real-world problem domains.
- encode high-cardinality categories into a target embedding
- create segmentation features with k-means clustering
- decompose a dataset's variant into feature with principal component analysis.

2. The Goal of Feature Engineering :

The goal of feature engineering is simply to make your data better suited to the problem at hand.

Consider "apparent temperature" measures like the heat index and the wind chill. These quantities attempt to measure the perceived temperature to humans based on an temperature, humidity, and wind speed, things which we can measure directly. You could think of an apparent temperature as the result of a kind of feature engineering, an attempt to make the observed data more relevant to what we actually care about.

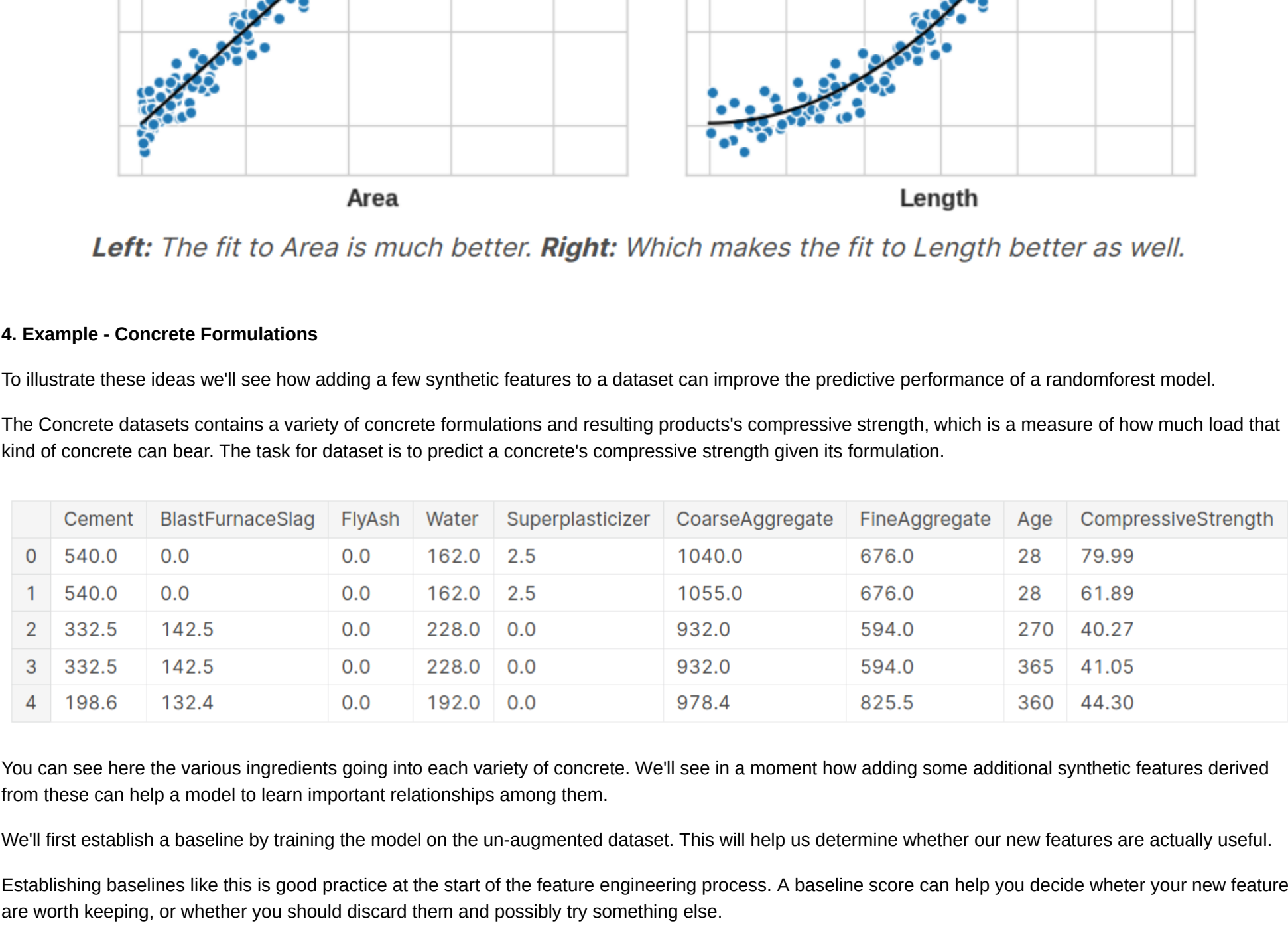
- improve a model's predictive performance
- reduce computational of data needs
- improve interpretability of the resulting model

3. A Guiding Principle of Feature Engineering :

For a feature to be useful, it must have a relationship to the target that your model is able to learn. Linear model, for instance, are only able to learn linear relationship. So, when using a linear model, your goal is to transform the features to make their relationship to the target linear.

The key idea here is that a transformation you apply to a features becomes in essence a part of the model itself. Say you were trying to predict the Price of square plots of land from the Length of one side. Fitting a linear model directly at Length gives poor results.

If we give the Length feature to get "Area", however, we create a linear relationship. Adding Area to the feature set means this linear model can now fit a parabola. Squaring a feature, in other words, gives the linear model the ability to fit squared features.



4. Example - Concrete Formulations

To illustrate these ideas we'll see how adding a few synthetic features to a dataset can improve the predictive performance of a randomforest model.

The Concrete datasets contains a variety of concrete formulations and resulting products's compressive strength, which is a measure of how much load that kind of concrete can bear. The task for dataset is to predict a concrete's compressive strength given its formulation.

	Cement	BlastFurnaceSlag	FlyAsh	Water	Superplasticizer	CoarseAggregate	FineAggregate	Age	CompressiveStrength
0	540.0	0.0	0.0	162.0	2.5	1040.0	676.0	28	79.99
1	540.0	0.0	0.0	162.0	2.5	1055.0	676.0	28	61.89
2	332.5	142.5	0.0	228.0	0.0	932.0	594.0	270	40.27
3	332.5	142.5	0.0	228.0	0.0	932.0	594.0	365	41.05
4	198.6	132.4	0.0	192.0	0.0	974.0	823.5	360	44.30

You can see here the various ingredients going into each variety of concrete. We'll see in a moment how adding some additional synthetic features derived from these can help a model to learn important relationships among them.

We'll first establish a baseline by training the model on the un-augmented dataset. This will help us determine whether our new features are actually useful.

Establishing baselines like this is good practice at the start of the feature engineering process. A baseline score can help you decide where your new features are worth a varying or a good deal better than the baseline score.

```
x = df.copy()
y = X_pos['compressivestrength']

# Train and score baseline model
baseline = RandomForestRegressor(criterion='mse', random_state=0)
baseline_score = cross_val_score(
    baseline, X, y, cv=5, scoring='neg_mean_absolute_error')
print(f'Baseline Score: {1 - baseline_score.mean()}'
      f'RMSE Baseline Score: {baseline_score.mean()}')

# Create synthetic features
X['CoarseAggregate'] = (X['FineAggregate'] * X['CoarseAggregate']) / X['Cement']
X['AggregateRatio'] = (X['CoarseAggregate'] + X['FineAggregate']) / X['Cement']
X['WaterRatio'] = (X['Water']) / X['Cement']

# Train and score model on dataset with additional ratio features
model = RandomForestRegressor(criterion='mse', random_state=0)
score = cross_val_score(
    model, X, y, cv=5, scoring='neg_mean_absolute_error')
score = 1 - score.mean()
print(f'RMSE Score with Ratio Features: {score.mean()}')
```

If you ever cook at home, you might know that the ratio of ingredients in a recipe is usually a better predictor of how the recipe turns out than their absolute amounts. We might reason then that ratios of the features above would be a good predictor of CompressiveStrength.

2. Mutual Information

First encountering a new datasets can sometimes feel overwhelming. You might be presented with hundreds or thousands of features without even a description to go by. Where do you even begin?

A great first step to consider a ranking with a feature utility metric, a function measuring associations between a feature and the target. Then you can choose a smaller set of the most useful features to develop initially and have more confidence that your time will be well spent.

The metric we will use is called "mutual information". Mutual information is a top level correlation in that it measures a relationship between two quantities. The advantage of mutual information is that it can detect any kind of relationship, while correlation only detects linear relationships.

Mutual information is a great general-purpose metric, and especially useful at the start of feature development when you might not know what model you'd like to use yet. It is:

- easy to use and interpret,
- computationally efficient,
- theoretically well founded,
- resistant to overfitting, and,
- able to detect any kind of relationship

1. Mutual information and what it measures

Mutual information describes relationships in terms of uncertainty. The mutual information between two quantities is a measure of the extent to which knowledge of one quantity reduces uncertainty in the other. If you know the value of a feature, how much more confident would you be about the target?

2. Interpreting Mutual Information Scores

The least possible mutual information between quantities is 0.0. When MI is zero, the quantities are independent - neither can tell you anything about the other. Conversely, if there there's an upper bound to what MI can be. In practice though values above 2.0 or so are uncommon.

Here are some things to remember when applying mutual information:

- MI can help you to understand the relative potential of a features as a predictor of the target, considered by itself.
- It's possible for a feature to be very informative when interacting with other features, but not so informative all alone. MI can't detect interaction between features. It's a univariate metric.
- The actual usefulness of a feature depends on the model you use it with. A feature is only useful to the extent that its relationship with the target is one your model can learn. But because a feature has high MI score doesn't mean you model will be able to do anything with that information. You may need to transform the feature to expose the association.

3. What need to focus

1. The scikit-learn algorithm for MI treats discrete features differently from continuous features. Consequently, you need to tell it which are which. As a rule of thumb, anything that must have a fixed dtype is not discrete. Categorical can be treated by giving them a label encoding(or ordinal encoding).
2. Scikit-learn has two mutual information metrics in its feature_selection module: one for real-valued targets(mutual_info_regression) and one for categorical targets(mutual_info_classif).

```
In [3]: from sklearn.feature_selection import mutual_info_regression

def make_mi_scores(X, y, discrete_features):
    mi_scores = mutual_info_regression(X, y, discrete=discrete_features)
    mi_scores = pd.Series(mi_scores, name="MI Scores", index=X.columns)
    mi_scores = mi_scores.sort_values(ascending=False)
    return mi_scores

def plot_mi_scores(scores):
    scores = scores.sort_values(ascending=True)
    width = np.arange(len(scores))
    ticks = list(scores.index)
    plt.barh(width, scores)
    plt.xticks(ticks, labels)
    plt.title('Mutual Information Scores')

Data visualization is a great addition to your feature-engineering toolbox. Along with utility metrics like mutual information, visualization like these can help you discover important relationships in your data.
```

3. Creating Features

Once you've identified a set of features with some potential, it's time to start developing them. In this lesson, you'll learn a number of common transformation you can do directly in Pandas.

1. Tips on Discovering New Features

- Understand the features. Refer to your dataset's data documentation, if available.
- Research the problem domain to acquire domain knowledge. If your problem is predicting house prices, do some research on real-estate for instance. Wikipedia can be a good starting point, but books and journal articles will often have the best information.
- Study previous work. Solution write-ups from past Kaggle competitions are a great resource.
- Use data visualization. Visualizations can reveal patterns in the distribution of a feature or complicated relationships that could be simplified. Be sure to visualize your dataset as you work through the feature engineering process.

2. Mathematical Transforms

Relationships among numerical features are often expressed through mathematical formulas, which you'll frequently come across as part of your domain research. In particular, you can apply arithmetic operations to columns just as if they were ordinary numbers. Data visualization can suggest transformations, often a "reshaping" of a feature through powers or logarithms.

3. Counts

Features describing the presence or absence of something often come in sets, the set of risk factor for disease, say. You can aggregate such features by creating a count. These features will be binary(1 for Present, 0 for Absent) or boolean (True or False). In Python, booleans can be added up just as if they were integers.

In traffic accidents are several features indicating whether some roadway event was near the accident. This will create a count of the total number of roadway features nearby using the sum method.

You could also use a DataFrame's built-in methods to create boolean values. In the Concrete dataset are the amounts of components in a concrete formulation. Many formulations lack one or more components. This will count how many components are in a formulation with the DataFrame's built-in greater-than-gt method.

```
# 2) Mathematical Transform
autos['displacement'] = (
    np.pi * ((8.5 * autos.bore) ** 2) * autos.stroke * autos.num_of_cylinders

# If the feature has 0.0 values, use np.log1p(log(1+x)) instead of np.log
accidents['LogWindSpeed'] = accidents.WindSpeed.apply(np.log1p)

# Plot a comparison
fig, axs = plt.subplots(1, 2, figsize=(8, 4))
sns.kdeplot(accidents.WindSpeed, shade=True, ax=axs[0])
sns.kdeplot(accidents.LogWindSpeed, shade=True, ax=axs[1]);

# 3) Counts
roadway_features = ["Aeridity", "Bump", "Crossing", "Giveway",
                   "Junction", "Junction", "Railway", "Roundabout", "Section",
                   "TrafficCalm", "TrafficSignal", "TrafficSignal"]
accidents['RoadwayFeatures'] = accidents[roadway_features].sum(axis=1)
```

4. Building-Up and Breaking-Down Features

Often you'll have complex strings that can usefully be broken into simpler pieces. Some common examples:

- ID numbers: 123-456-789
- Phone numbers: (909) 555-0232
- Street address: 3241 Maple Ln, Goose City, NV

Features like these will often have some kind of structure that you can make use of. US phone numbers, for instance that tells you the location of the caller. As always, some research can pay off here.

The str accessor lets you apply string method like split directly to columns.

5. Group Transforms

Finally we have group transforms, which aggregate information across multiple rows grouped by some category. With a group transform you can create features like "the average income of a person's state of residence" or "the proportion of movies released on a weekday, by genre". If you had discovered a category (consider, a group transform over that category could be something good to investigate).

Using a aggregation function, a group transform combines two features: a categorical features that provides the grouping and another feature whose value you wish to aggregate. For an average income by state, you would choose State for grouping feature mean for the aggregation function and income for the aggregated feature. To compute this in Pandas, we use the groupby and transform method:

```
customer['AverageIncome'] = (
    customer['AverageIncome'] * 1
    customer.groupby("State")["Income"]
    .transform("mean")
)

# Frequency of categorical value
customer["StateFreq"] = customer.groupby("State")["State"].transform("count") / customer.State.count()
```

6. Tips on Creating Features

- Linear models learn sums and differences naturally, but can't learn anything more complex.
- Ratios seem to be difficult for most model to learn. Ratio combinations often lead to some easy performance gains.
- Linear models and neural nets generally do better with normalized features. Neural nets especially need features scaled to avoid not too far from 0. Tree-based models (like random forest and XGBoost) can sometimes benefit from normalization, but usually much less so.
- The models can learn to approximate almost any combination of features, but when a combination is especially important they can still benefit from having it explicitly created, especially when data is limited.
- Counts are especially helpful for tree models, since these models don't have a natural way of aggregating information across many features at once.

3. Clustering with K-Means

This lesson and the next make use of what are known unsupervised learning algorithms. Unsupervised algorithms don't make use of a target. Instead, their purpose is to learn some property of the data, to represent the structure of the features in a certain way. In the context of feature engineering for predictions, you could think of an unsupervised algorithms as a "feature discovery" technique.

Clustering simply means the assigning of data points to group based upon how similar the points are to each other. A clustering algorithm makes "kinds of feature flock together", so to speak.

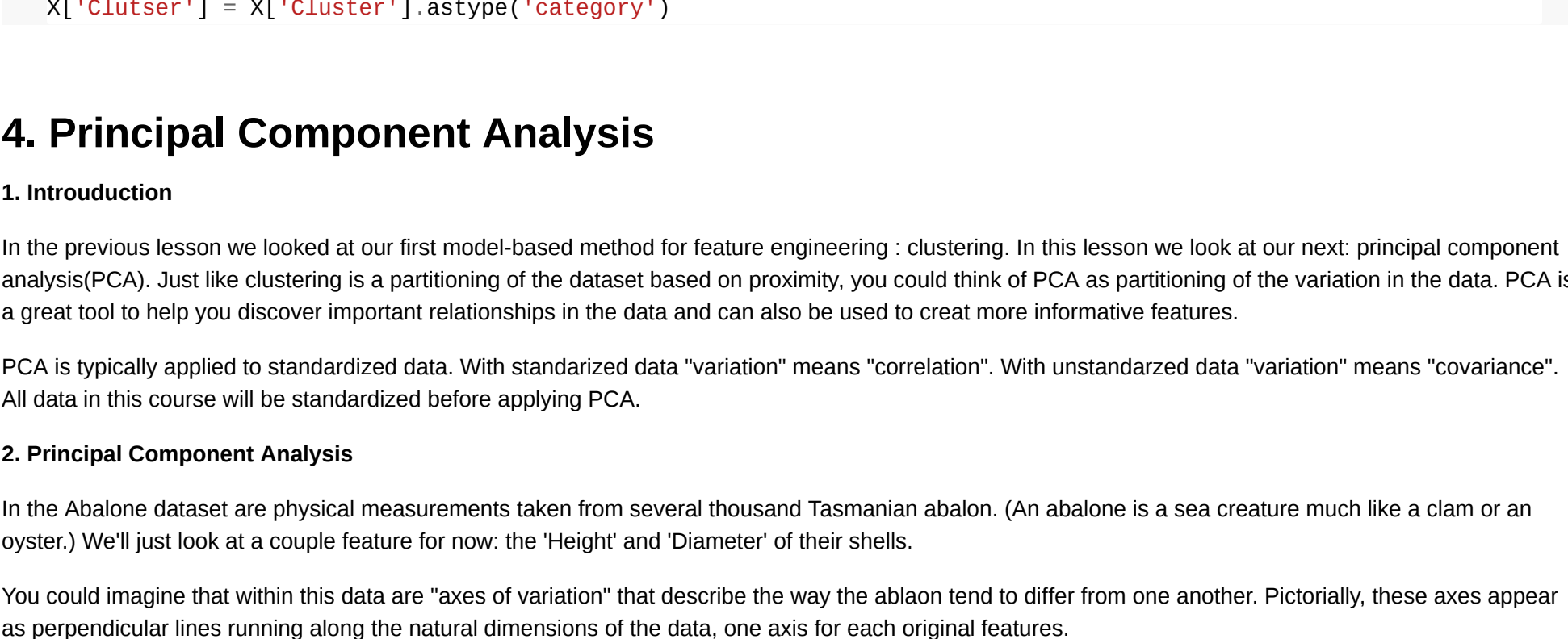
When used for feature engineering, we could attempt to discover groups of customers representing a market segment, for instance, or geographic areas that share similar weather patterns. Adding a feature of clustering labels can help machine learning models untangle complicated relationships of space or proximity.

1. Cluster Labels as a Feature

Applied to a single real-valued feature, clustering acts like a traditional "binning" or "discretization" transform. On multiple features, it's like "multi-dimensional binning".

It's important to remember that this Cluster feature is categorical. Here, it's shown with a bar encoding (same sign), but in the Shape component algorithm would produce; depending on your model, a label encoding may be more appropriate.

The motivating idea of adding cluster labels is that the clusters will break up complicated relationships across features into simpler chunks. Our model can then just learn the simpler chunks one-by-one instead having to learn the complicated whole all at once. It's a "divide and conquer" strategy.



2. K-Means Clustering

There are a great many clustering algorithms. They differ primarily in how they measure "similarity" or "proximity" and in what kinds of features they work with. The algorithm we'll use, k-means, is intuitive and easy to apply in a feature engineering context. Depending on your application another algorithm might be more appropriate.

K-means clustering measures similarity using ordinary straight-line distance (Euclidean distance, in other words). It creates clusters by placing a number of points, called centroids, inside the feature-space. Each point in the dataset is assigned to the cluster of whichever centroid it's closest to. The "K" in k-means is how many centroids (that is, clusters) it creates. You define k; you tell it.

You could imagine each centroid capturing points through a sequence of radiating circles. When sets of circles from competing centroids overlap they form a line. The result is what's called a voronoi tessellation. The tessellation shows you what clusters future data will be assigned; the tessellation is essentially what k-means learns from its training data.

Let's review how the k-means algorithm learns the clusters and what that means for feature engineering. We'll focus on three parameters from scikit-learn's implementation: n_clusters, max_iter, and n_init.

It's a simple two-step process. The algorithm starts by randomly initializing some predefined number (n_clusters) of centroids. It then iterates over these two operation:

1. assign points to the nearest cluster centroid
2. move each centroid to minimize the distance to its points

It iterates over these two steps until the centroids aren't moving anymore, or until some maximum number of iterations has passed (max_iter).

It often happens that the initial random position of the centroids ends in a poor clustering. For this reason the algorithm repeats a number of times(n_init) and return the clustering that has the least total distance between each point and its centroid, the optimal clustering.

You may need to increase the max_iter for a large number of clusters or n_init for a complex dataset. Ordinarily though the only parameter you'll need to choose your self is n_clusters (k, max_iter). The best partitioning for a set of features depends on the model you're using and what you're trying to predict, so it's best to tune it like any hyperparameter (through cross-validation, say).

```
# Example of making new feature using K-means
kmeans = KMeans(n_clusters = 6)
X['cluster'] = kmeans.fit_predict(X)
X['cluster'] = X['cluster'].astype('category')
```

4. Principal Component Analysis

1. Introduction

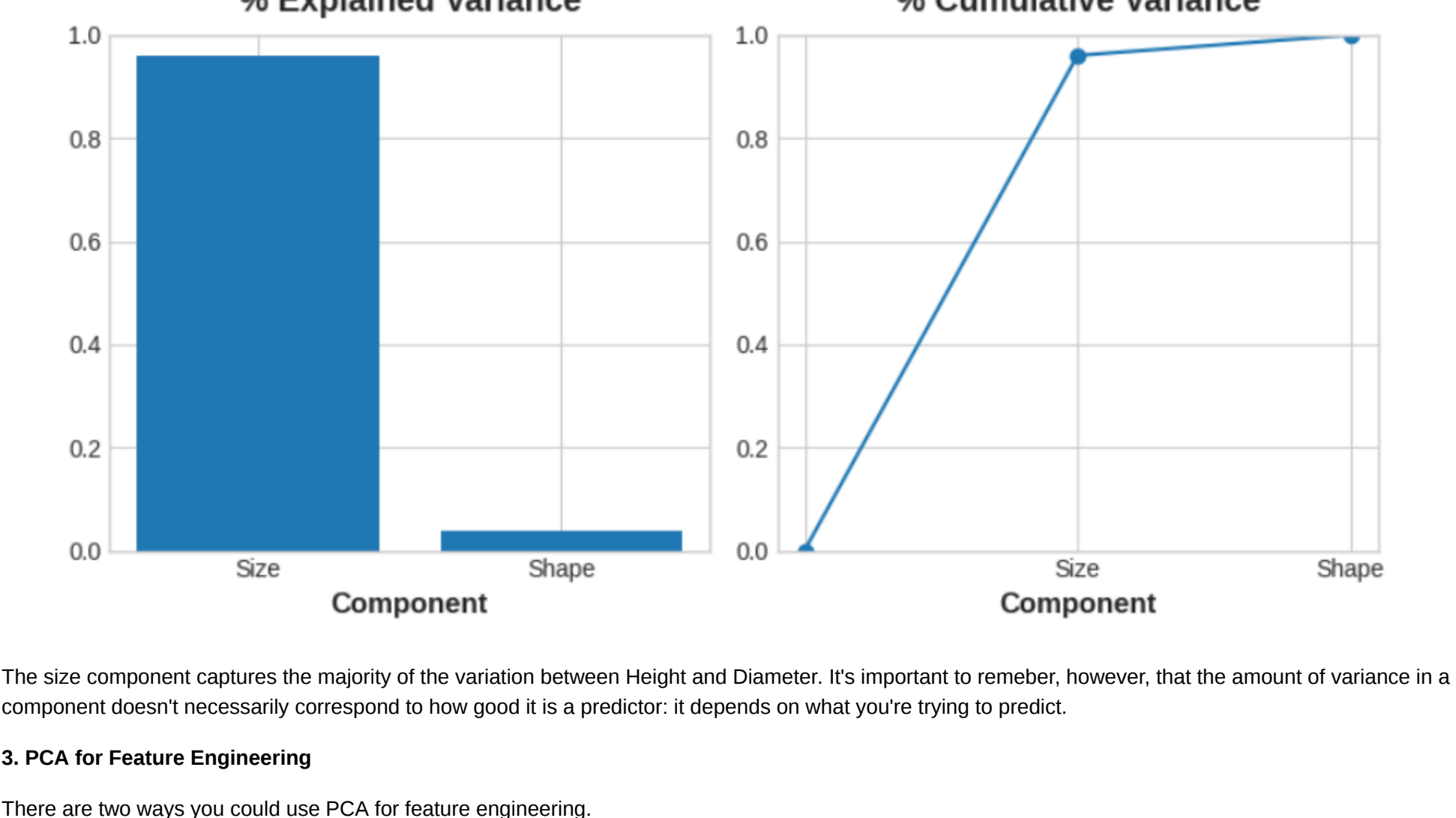
In the previous lesson we looked at our first model-based method for feature engineering - clustering. In this lesson we look at our next: principal component analysis(PCA). Just like clustering is a partitioning of the dataset based on proximity, you could think of PCA as partitioning the variance in the data. PCA is a great tool to help you discover important relationships in the data and can also be used to create more informative features.

PCA is typically applied to standardized data. With standard data "variance" means "correlation". With unstandardized data "variance" means "covariance". All data in this course will be standardized before applying PCA.

2. Principal Component Analysis

In the Abalone dataset are physical measurements taken from several thousand Tasmanian abalone. (An abalone is a sea creature much like a clam or an oyster). We'll just look at a couple feature for now: the 'Height' and 'Diameter' of their shells.

You could imagine that within this data are "axes of variation" that describe the way the abalone tend to differ from one another. Pictorially, these axes appear as perpendicular lines running along the natural dimensions of the data, one axis for each dimension.



The new features PCA constructs are actually just linear combinations (weighted sums) of the original features:

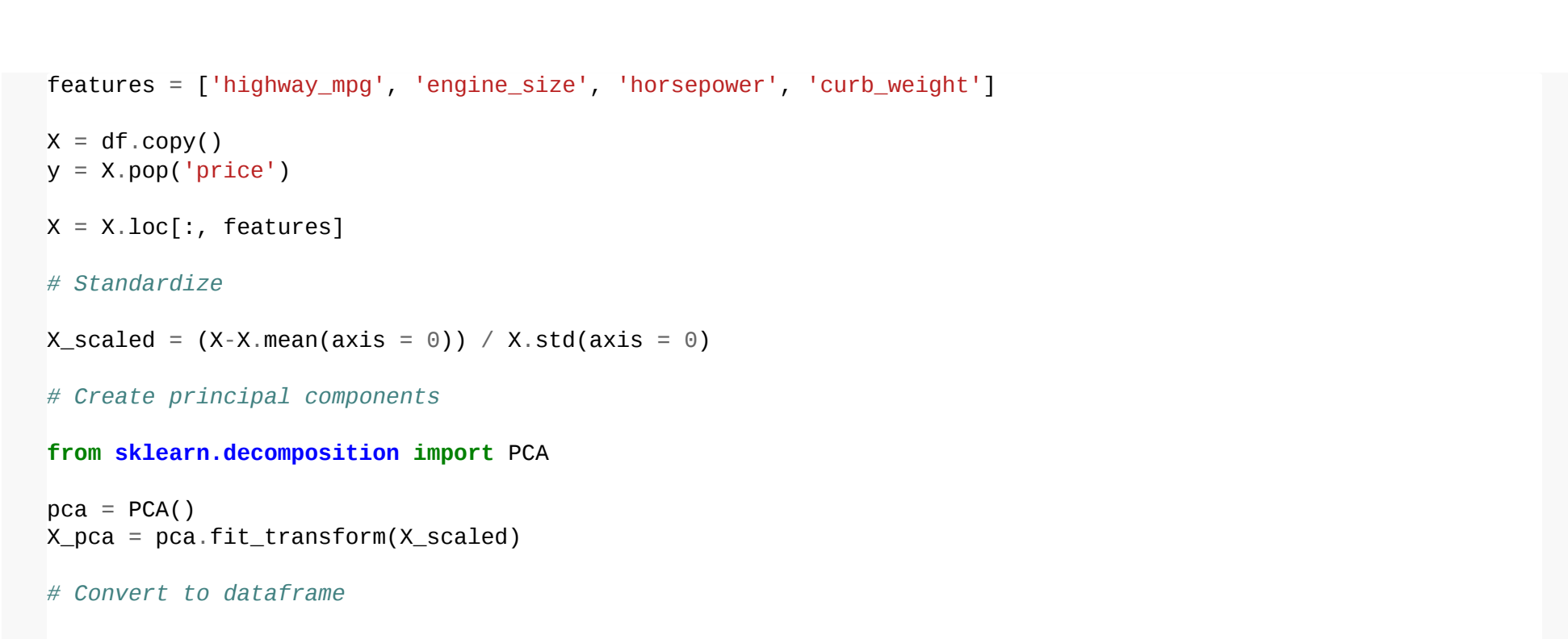
- $Size = 0.707 \times Height + 0.707 \times Diameter$
- $Shape = 0.707 \times Height - 0.707 \times Diameter$

These new features are called the principal components of the data. The weights themselves are called loadings. There will be as many principal components as there are features in the original dataset. If we had used ten features instead of two, we would have ended up with ten components.

Features	Components	Size(PC1)	Shape(PC2)
Height	0.707	0.707	
Diameter	0.707	-0.707	

This table of loadings is telling us that in the size component, Height and Diameter vary in the same direction (same sign), but in the Shape component the vary in opposite directions (opposite sign). In each component, the loadings are all of the same magnitude are so the features contribute equally in both.

PCA also tells us the amount of variation in each component. We can see from the figures that there is more variation in the data along the Size component than along the Shape component. PCA makes this precise through each component's percent of explained variance.



The size component captures the majority of the variation between Height and Diameter. It's important to remember, however, that the amount of variance in a component doesn't necessarily correspond to how good it is a predictor; it depends on what you're trying to predict.

3. PCA for Feature Engineering

There are two ways you could use PCA for feature engineering.

The first way is to use it as a descriptive technique. Since the components tell you about the variation, you could compute the MI scores for the components and see what kind of variation is most predictive of your target. That could give you ideas for kinds of feature to create - a product of 'Height' and 'Diameter' if 'Size' is important, say, or a ratio of 'Height' and 'Diameter' if 'Shape' is important. You could even try clustering on one or more of the high-scoring components.

The second way is to use the components themselves as features. Because the components expose the variation structure of the data directly, they can often be more informative than the original features. Here are some use-cases:

- Dimensionality reduction: When your features are highly redundant (multicollinear, specifically), PCA will partition out the redundancy into one or more near-zero variance components, which you can then drop since they will contain little or no information.
- Anomaly detection: Unusual variation, not captured from the original features, will often show up in the low-variance components. These components could be highly informative in an anomaly or outlier detection task.
- Noise reduction: A collection of sensor readings will often share some common background noise. PCA can sometimes collect the (informative) signal into a small number of features while leaving the noise alone, thus boosting the signal-to-noise ratio.
- Decorrelation: Some ML algorithms struggle with highly-correlated features. PCA transforms correlated features into uncorrelated components, which could be easier for your algorithm to work with.

PCA basically gives you direct access to the correlational structure of your data. You'll no doubt come up with application of your own!

• PCA Best Practices

- PCA only works with numeric features, like continuous quantities or counts.
- PCA is sensitive to scale. It's good practice to standardize your data before applying PCA, unless you know you have good reason not to.
- Consider removing or constraining outliers. Since they can have an undue influence on the results.

```
features = ['highway_mpg', 'engine_size', 'horsepower', 'curb_weight']

X = df.copy()
y = X_pos['price']

X = X.loc[:, features]

# Standardize
X_scaled = (X - X.mean(axis = 0)) / X.std(axis = 0)

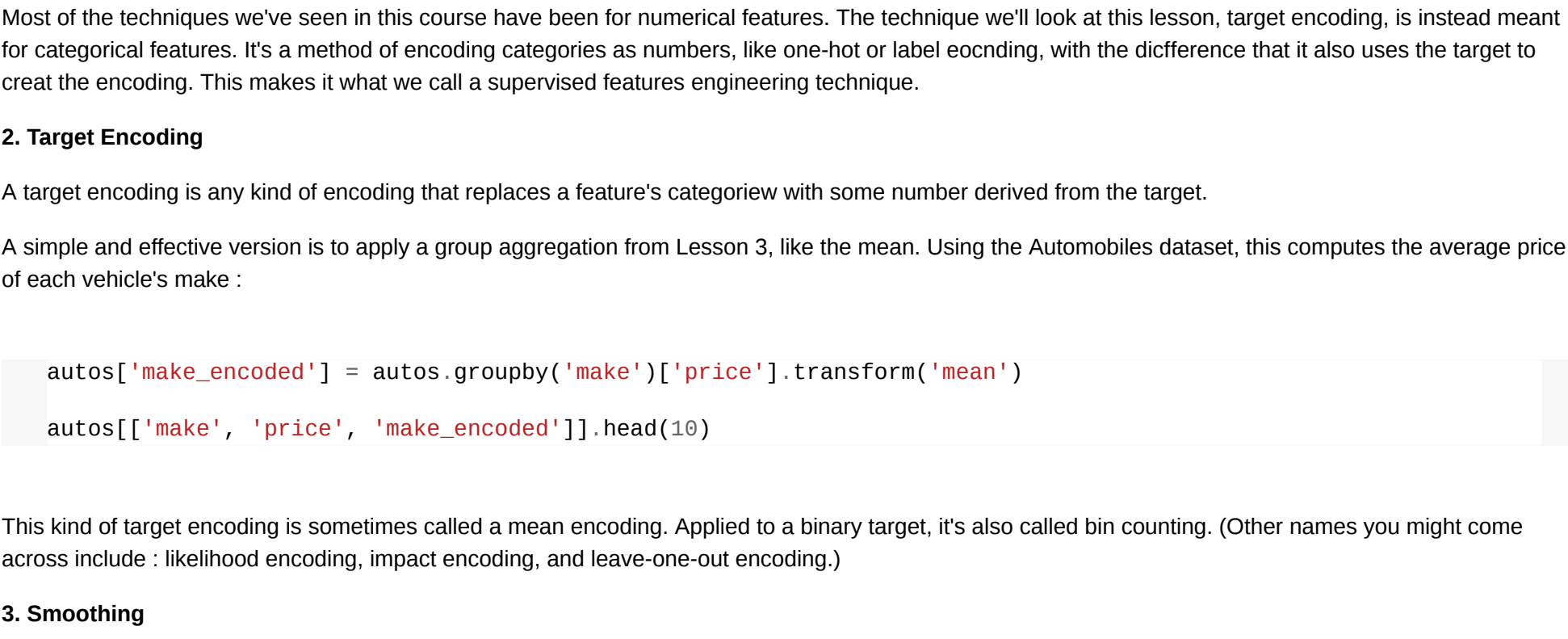
# Create principal components
from sklearn.decomposition import PCA
pca = PCA()
X_pca = pca.fit_transform(X_scaled)

# Convert to dataframe
component_names = [f'PC{i+1}' for i in range(X_pca.shape[1])]
X_pca = pd.DataFrame(X_pca, columns = component_names)
```

After fitting, the PCA instance contains the loadings in its components_ attribute.

	PC1	PC2	PC3	PC4
highway_mpg	-0.492347	0.770892	0.070142	-0.397996
engine_size	0.503859	0.626709	0.019960	0.594107
horsepower	0.500448	0.013788	0.731093	-0.463534
curb_weight	0.503262	0.113008	-0.678369	-0.523232

Recall that the signs and magnitudes of a component's loadings tell us what kind of variation it's captured. The first component (PC1) shows a contrast between large, powerful vehicles with poor gas mileage, and smaller, more economical vehicles with good gas mileage. We might call this the "Luxury/Economy" axis. The next figures show that our four chosen features mostly vary along the Luxury/Economy axis.



Let's also look at the MI scores of the components. Not surprisingly, PC1 is highly informative, though the remaining components, despite their small variance, still have significant relationships with price. Examining those components could be worthwhile to find relationships not captured by the main Luxury/Economy axis.

```
PC1    1.013880
PC2    0.379448
PC3    0.386592
PC4    0.284447
Name: MI Scores, dtype: float64
```

5. Target Encoding

1. Introduction

Most of the techniques we've seen in this course have been for numerical features. The technique we'll look at this lesson, target encoding, is instead meant for categorical features. It's a method of encoding categories as numbers, like one-hot or label encoding, with the difference that it also uses the target to create the encoding. This makes it what we call a supervised features engineering technique.

2. Target Encoding

A target encoding is any kind of encoding that replaces a feature's categories with some number derived from the target.

A simple and effective version is to apply a group aggregation from Lesson 3, like the mean. Using the Automobiles dataset, this computes the average price of each vehicle's make:

```
autos['make_encoded'] = autos.groupby('make')['price'].transform('mean')

autos['make', 'price', 'make_encoded'].head(10)
```

This kind of target encoding is sometimes called a mean encoding. Applied to a binary target, it's also called bin counting. (Other names you might come across include: likelihood encoding, impact encoding, and leave-one-out encoding.)

3. Smoothing

An encoding like this presents a couple of problems, however. First are unknown categories. Target encodings create a special risk of overfitting, which means they need to be trained on an "independent" encoding split. When you bin the encoding to future splits, Pandas will fill in missing values for the categories not present in the encoding split. These missing values you would have to impute somehow.

Second are rare categories. When a category only occurs a few times in the dataset, any statistics calculated on its group are unlikely to be accurate. In the Automobiles dataset, the mercury make only occurs once. The "mean" price we calculated is just the price of that one vehicle, which might not be very representative of any Mercury car might sell in the future. Target encoding rare categories can make overfitting more likely.

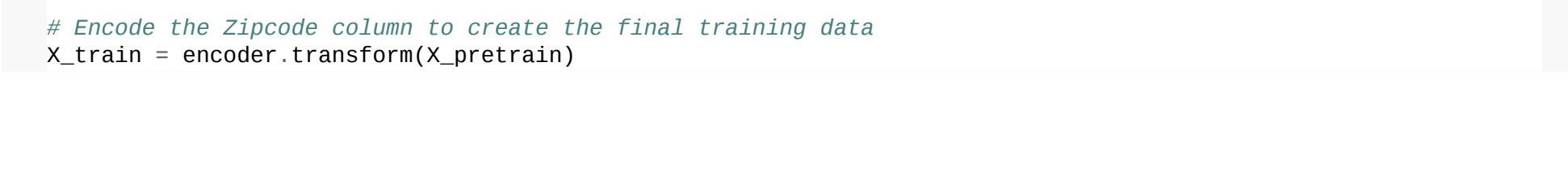
A solution to these problems is to add smoothing. The idea is to blend the in-category average with the overall average. Rare categories get less weight on their category average, while missing categories just get the overall average.

- In pseudocode: encoding = weight in_category + (1 - weight) overall

where weight is a value between 0 and 1 calculated from the category frequency. An easy way to determine the value for weight is to compute an m estimate:

$$weight = n / (n + m)$$

where n is the total number of times that category occurs in the data. The parameter m determines the "smoothing factor". Larger value of m put more weight on the overall estimate.



In the Automobiles dataset there are three cars with the make 'volvo'. If you chose m = 2.0, then the chevrolet category would be encoded with 60% of the average Chevrolet price plus 40% of the overall average price. When choosing a value of m, consider how risky you expect the categories to be. Does the price of a vehicle vary a great deal within each make? Would you need a lot of data to get good estimates? If so, it could be better to choose a larger value for m. If the average price for each make were relatively stable, a smaller value could be okay.

- Use cases for Target Encoding
- High-cardinality features: A feature with a large number of categories can be troublesome to encode. A one-hot encoding would generate too many features and alternatives, like label encoding, might not be appropriate for that feature. A target encoding derives numbers for the categories using the feature's most important property: its relationship with the target.
- Domain-motivated features: From your experience, you might suspect that a categorical features should be important even if it scored poorly with a feature metric. A target encoding can help reveal a feature's true informativeness.

```
# Split datasets
X = df.copy()
y = X_pos['Rating']

X_encode = X.sample(Frac=0.25)
y_encode = y[X_encode.index]
X_pretrain = X.drop(X_encode.index)
y_train = y[X_pretrain.index]

from category_encoder import MEstimateEncoder

# Create the encoder instance. Choose m to control noise.
encoder = MEstimateEncoder(cols=['zipcode'], m=5.0)

# Fit the encoder on the encoding split.
encoder.fit(X_encode, y_encode)

# Encode the zipcodes to create the final training data
X_train = encoder.transform(X_pretrain)
```