

**∓** Filter

**Data** 

Classification (70 min)

■ Introduction (3 min)

First steps (5 min)

**■** Binning (15 min)

■ Scrubbing (5 min)

features (5 min)

**■** Conclusion (2 min)

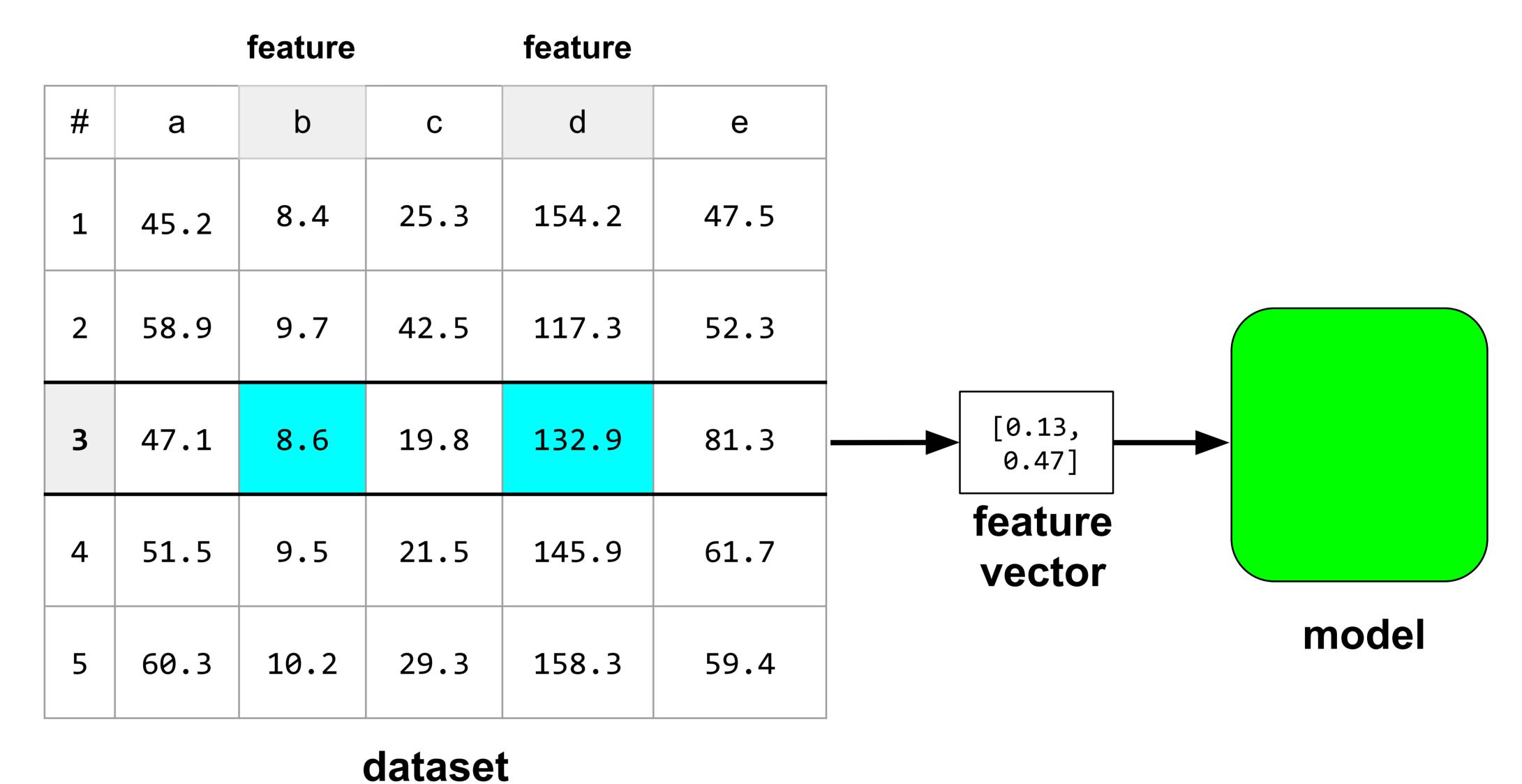
→ What's next

min)

feature vectors (5 min)

However, feature vectors seldom use the dataset's raw values. Instead, you must typically process the dataset's values into

representations that your model can better learn from. So, a more realistic feature vector might look something like this:



Wouldn't a model produce better predictions by training from the actual values in the dataset than from altered values?

Figure 3. A more realistic feature vector.

Surprisingly, the answer is no.

You must determine the best way to represent raw dataset values as trainable values in the feature vector. This process is called **feature engineering**, and it is a vital part of machine learning. The most common feature engineering techniques are:

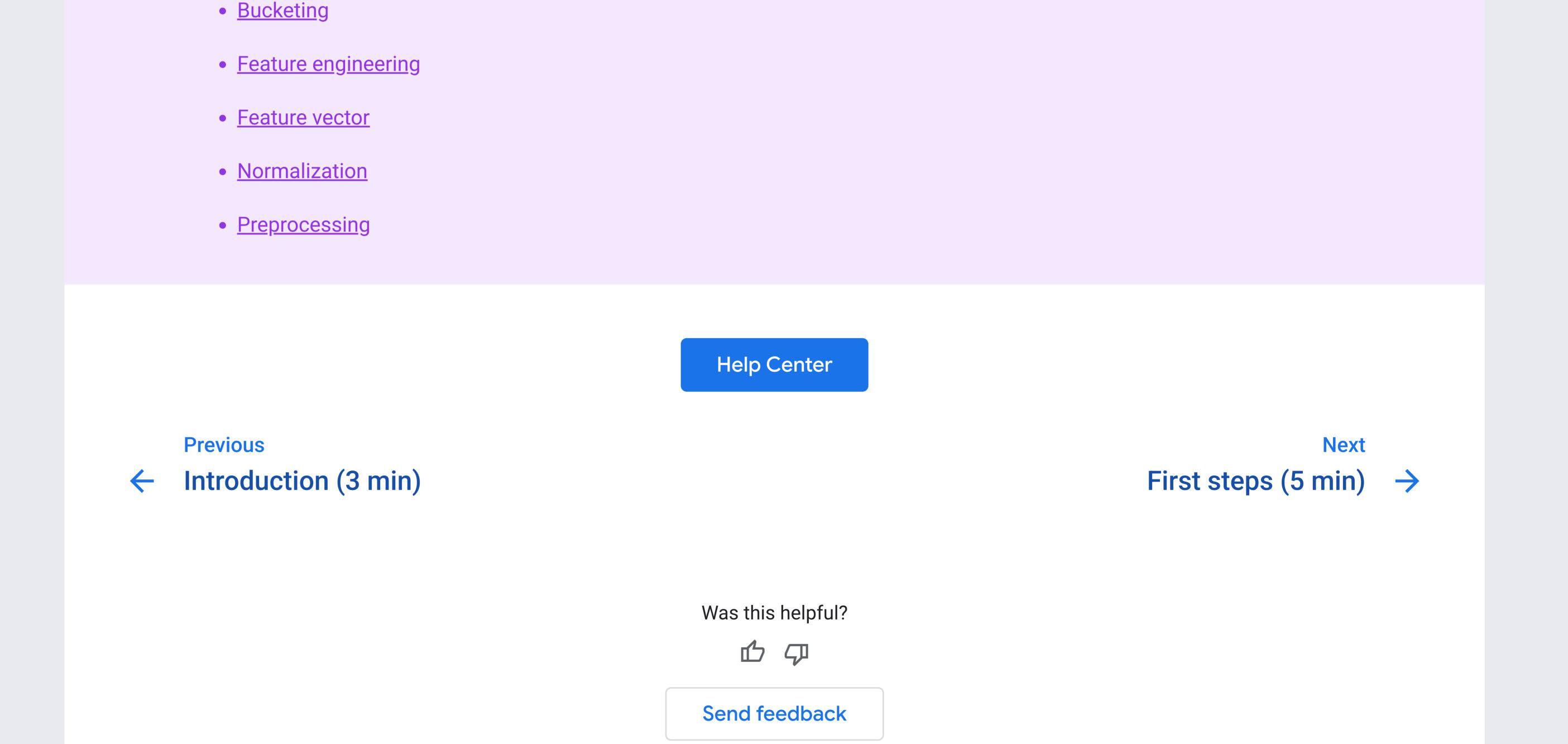
• Binning (also referred to as bucketing): Converting numerical values into buckets of ranges.

• Normalization: Converting numerical values into a standard range.

This unit covers normalizing and binning. The next unit, Working with categorical data, covers other forms of

preprocessing, such as converting non-numerical data, like strings, to floating point values.

Every value in a feature vector must be a floating-point value. However, many features are naturally strings or other nonnumerical values. Consequently, a large part of feature engineering is representing non-numerical values as numerical values. You'll see a lot of this in later modules.



2.0 License. For details, see the Google Developers Site Policies. Java is a registered trademark of Oracle and/or its affiliates. Last updated 2025-01-02 UTC.

Except as otherwise noted, the content of this page is licensed under the Creative Commons Attribution 4.0 License, and code samples are licensed under the Apache

**Key terms:** 

• <u>Binning</u>

