

Cmd 11

```
1 lr = LogisticRegression(max_iter=10000)
```

Command took 0.04 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 12

Create Feature Vectors

🤖 sklearn wants the shape of our data to be a matrix for our feature(s) and the shape of our target to be a vector. This is why you will see two square brackets around our feature - a matrix - and a single set of square brackets around our target - a vector.

Cmd 13

```
1 X = ht_agg_pandas_df[['mean_bmi']]
```

Command took 0.04 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 14

Create Target Vector

An additional step, not required when perform the linear regression, is necessary to encode our target vector when performing a logistic regression.

This has to do with the way the lifestyle labels are stored.

Cmd 15

```
1 ht_agg_pandas_df["lifestyle"].unique()
```

```
Out[11]: array(['Cardio Enthusiast', 'Athlete', 'Sedentary', 'Weight Trainer'],
              dtype=object)
```

Command took 0.04 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 16

⚠️ Each lifestyle is recorded as a string value.

sklearn models can only work on numerical values. For this reason, it is required to numerically encode our lifestyle values.

We will use an sklearn transformer to do this encoding.

An sklearn transformer is like an sklearn estimator except rather than using it to `.predict()` or `.score()`, we will use it to `.transform()`

```
estimator.fit(data)
estimator.transform(data)
```

Cmd 17

```
1 from sklearn.preprocessing import LabelEncoder
2
3 le = LabelEncoder()
4 lifestyle = ht_agg_pandas_df['lifestyle']
5 le.fit(lifestyle)
6 y = le.transform(lifestyle)
```

Command took 0.03 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 18

Fit the Model

Next, fit our model, using the same `.fit(feature, target)` pattern we learned earlier.

The model will learn the relationship between features and target, i.e. we will "train or fit the model".

Cmd 19

```
1 lr.fit(X, y)
```

```
Out[13]: LogisticRegression(max_iter=10000)
```

Command took 0.32 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 20

Evaluate the model

Finally, use the `.score()` method to evaluate the single-variable model.

Note that a classifier estimator in sklearn uses accuracy for scoring by default.

Cmd 21

```
1 lr.score(X, y)
```

```
Out[14]: 0.417
```

Command took 0.04 seconds -- by t.jamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 22

Your Turn

Exercise 1: Single-Variable Logistic Regression

Fit a single-variable logistic model for each of the remaining feature.

1. prepare a feature matrix for each of these features:

- `mean_bmi`
- `mean_active_heartrate`
- `mean_resting_heartrate`

- `mean_vo2`
- 2. fit a single-variable logistic model for each of these features
- 3. evaluate using `.score()` each of these models and print the result

Cmd 23

```
1 # ANSWER
2 X_bmi = ht_agg_pandas_df[['mean_bmi']]
3 X_active_hearttrate = ht_agg_pandas_df[['mean_active_hearttrate']]
4 X_resting_hearttrate = ht_agg_pandas_df[['mean_resting_hearttrate']]
5 X_vo2 = ht_agg_pandas_df[['mean_vo2']]
6
7 lr_bmi = LogisticRegression(max_iter=10000)
8 lr_active_hearttrate = LogisticRegression(max_iter=10000)
9 lr_resting_hearttrate = LogisticRegression(max_iter=10000)
10 lr_vo2 = LogisticRegression(max_iter=10000)
11
12 lr_bmi.fit(X_bmi, y)
13 lr_active_hearttrate.fit(X_active_hearttrate, y)
14 lr_resting_hearttrate.fit(X_resting_hearttrate, y)
15 lr_vo2.fit(X_vo2, y)
16
17 print("bmi: ", lr_bmi.score(X_bmi, y))
18 print("active_hearttrate: ", lr_active_hearttrate.score(X_active_hearttrate, y))
19 print("resting_hearttrate: ", lr_resting_hearttrate.score(X_resting_hearttrate, y))
20 print("vo2: ", lr_vo2.score(X_vo2, y))
```

bmi: 0.417
active_hearttrate: 0.576
resting_hearttrate: 0.586
vo2: 0.557

Command took 1.24 seconds -- by tjamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 24

Question: Which of these single-variable models is the best at predicting lifestyle?

Cmd 25

Demonstration

Multiple-Variable Logistic Regression

Our next set of models will use more than one feature and but still have a single target.

Cmd 26

Display results from previous models

Before we train this new model, let's display the results from the previous models for comparison.

Cmd 27

```
1 print("bmi: ", lr_bmi.score(X_bmi, y))
2 print("active_hearttrate: ", lr_active_hearttrate.score(X_active_hearttrate, y))
3 print("resting_hearttrate: ", lr_resting_hearttrate.score(X_resting_hearttrate, y))
4 print("vo2: ", lr_vo2.score(X_vo2, y))
```

bmi: 0.417
active_hearttrate: 0.576
resting_hearttrate: 0.586
vo2: 0.557

Command took 0.05 seconds -- by tjamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 28

```
1 X_bmi_act_hr = ht_agg_pandas_df[['mean_bmi', 'mean_active_hearttrate']]
2 lr_bmi_act_hr = LogisticRegression(max_iter=10000)
3 lr_bmi_act_hr.fit(X_bmi_act_hr, y)
4 print("bmi_act_hr: ", lr_bmi_act_hr.score(X_bmi_act_hr, y))
```

bmi_act_hr: 0.5826666666666667

Command took 0.58 seconds -- by tjamesbu@gmail.com at 4/2/2021, 12:10:36 PM on My Cluster

Cmd 29

Your Turn

Exercise 2: Multi-Variable Logistic Regression

👨‍🔬 Note that this two feature model performs better than any of the single feature models.

Fit four multiple-variable logistic models.

1. prepare a feature matrix
2. fit a logistic model for each of feature matrix
3. evaluate each model using `.score()` and print the result

👨‍🔬 Did you try any models with more than two features? Multiple-variable logistic regression models can use any or all of the features.

Cmd 30

```
1 # ANSWER
2 X_1 = ht_agg_pandas_df[['mean_active_hearttrate', 'mean_resting_hearttrate']]
3 X_2 = ht_agg_pandas_df[['mean_active_hearttrate', 'mean_vo2']]
4 X_3 = ht_agg_pandas_df[['mean_active_hearttrate', 'mean_bmi', 'mean_vo2']]
5 X_4 = ht_agg_pandas_df[['mean_active_hearttrate', 'mean_bmi', 'mean_vo2', 'mean_resting_hearttrate']]
6
7 lr_1 = LogisticRegression(max_iter=10000)
8 lr_2 = LogisticRegression(max_iter=10000)
9 lr_3 = LogisticRegression(max_iter=10000)
10 lr_4 = LogisticRegression(max_iter=10000)
```

```
11
12 lr_1.fit(X_1, y)
13 lr_2.fit(X_2, y)
14 lr_3.fit(X_3, y)
15 lr_4.fit(X_4, y)
16
17 print("model 1: ", lr_1.score(X_1, y))
18 print("model 2: ", lr_2.score(X_2, y))
19 print("model 3: ", lr_3.score(X_3, y))
20 print("model 4: ", lr_4.score(X_4, y))

model 1:  0.605
model 2:  0.5963333333333334
model 3:  0.6006666666666667
model 4:  0.6043333333333333

Command took 6.86 seconds -- by tjamesbu@gmail.com at 4/2/2021, 12:10:37 PM on My Cluster
```

Cmd 31

Which of these models is the best at predicting lifestyle?

Cmd 32

© 2021 Databricks, Inc. All rights reserved.
Apache, Apache Spark, Spark and the Spark logo are trademarks of the [Apache Software Foundation](#).

[Privacy Policy](#) | [Terms of Use](#) | [Support](#)

Shift+Enter to run