

scikit-learn

Lecture 14

Dr. Colin Rundel

scikit-learn

Scikit-learn is an open source machine learning library that supports supervised and unsupervised learning. It also provides various tools for model fitting, data preprocessing, model selection, model evaluation, and many other utilities.

- Simple and efficient tools for predictive data analysis
- Accessible to everybody, and reusable in various contexts
- Built on NumPy, SciPy, and matplotlib
- Open source, commercially usable - BSD license

Submodules

The `sklearn` package contains a large number of submodules which are specialized for different tasks / models,

- `sklearn.base` - Base classes and utility functions
- `sklearn.calibration` - Probability Calibration
- `sklearn.cluster` - Clustering
- `sklearn.compose` - Composite Estimators
- `sklearn.covariance` - Covariance Estimators
- `sklearn.cross_decomposition` - Cross decomposition
- `sklearn.datasets` - Datasets
- `sklearn.decomposition` - Matrix Decomposition
- `sklearn.discriminant_analysis` - Discriminant Analysis
- `sklearn.ensemble` - Ensemble Methods
- `sklearn.exceptions` - Exceptions and warnings
- `sklearn.experimental` - Experimental
- `sklearn.feature_extraction` - Feature Extraction
- `sklearn.feature_selection` - Feature Selection
- `sklearn.gaussian_process` - Gaussian Processes
- `sklearn.impute` - Impute
- `sklearn.inspection` - Inspection
- `sklearn.isotonic` - Isotonic regression
- `sklearn.kernel_approximation` - Kernel Approximation
- `sklearn.kernel_ridge` - Kernel Ridge Regression
- `sklearn.linear_model` - Linear Models
- `sklearn.manifold` - Manifold Learning
- `sklearn.metrics` - Metrics
- `sklearn.mixture` - Gaussian Mixture Models
- `sklearn.model_selection` - Model Selection
- `sklearn.multiclass` - Multiclass classification
- `sklearn.multioutput` - Multioutput regression and classification
- `sklearn.naive_bayes` - Naive Bayes
- `sklearn.neighbors` - Nearest Neighbors
- `sklearn.neural_network` - Neural network models
- `sklearn.pipeline` - Pipeline
- `sklearn.preprocessing` - Preprocessing and Normalization
- `sklearn.random_projection` - Random projection
- `sklearn.semi_supervised` - Semi-Supervised Learning
- `sklearn.svm` - Support Vector Machines
- `sklearn.tree` - Decision Trees
- `sklearn.utils` - Utilities

Model Fitting

Sample data

To begin, we will examine a simple data set on the size and weight of a number of books. The goal is to model the weight of a book using some combination of the other features in the data.

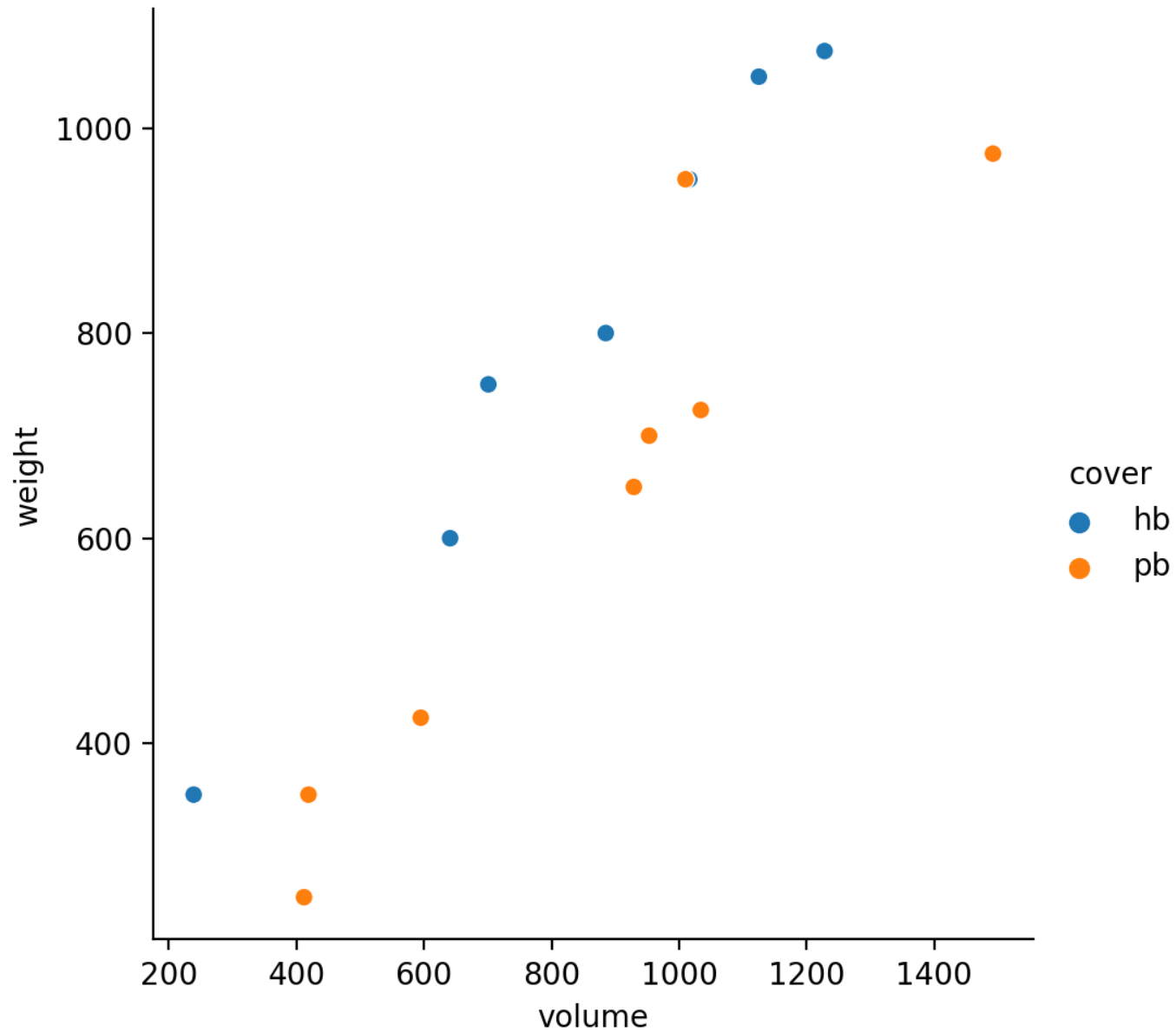
The included columns are:

- `volume` - book volumes in cubic centimeters
- `weight` - book weights in grams
- `cover` - a categorical variable with levels "hb" hardback, "pb" paperback

```
1 books = pd.read_csv("data/daag_books.csv"); book
```

	volume	weight	cover
0	885	800	hb
1	1016	950	hb
2	1125	1050	hb
3	239	350	hb
4	701	750	hb
5	641	600	hb
6	1228	1075	hb
7	412	250	pb
8	953	700	pb
9	929	650	pb
10	1492	975	pb
11	419	350	pb
12	1010	950	pb
13	595	425	pb
14	1034	725	pb

```
1 sns.relplot(data=books, x="volume", y="weight", hue="cover")
```



Linear regression

scikit-learn uses an object oriented system for implementing the various modeling approaches, the class for `LinearRegression` is part of the `linear_model` submodule.

```
1 from sklearn.linear_model import LinearRegression
```

Each modeling class needs to be constructed (potentially with options) and then the resulting object will provide attributes and methods.

```
1 lm = LinearRegression()
2
3 m = lm.fit(
4     X = books[["volume"]],
5     y = books.weight
6 )
7
8 m.coef_
```

```
array([0.70863714])
```

```
1 m.intercept_
```

```
107.679310613766
```

Note `lm` and `m` are labels for the same object,

```
1 lm.coef_
```

```
array([0.70863714])
```

```
1 lm.intercept_
```

```
107.679310613766
```

A couple of considerations

When fitting a model, scikit-learn expects `X` to be a 2d array-like object (e.g. a `np.array` or `pd.DataFrame`) and so it will not accept a `pd.Series` or 1d `np.array`.

```
1 lm.fit(  
2     X = books.volume,  
3     y = books.weight  
4 )
```

Error: ValueError: Expected 2D array, got 1D array in array=[885 1016 1125 239 701 641 1228 412 953 1034].

Reshape your data either using `array.reshape(-1, 1)` :

```
1 lm.fit(  
2     X = np.array(books.volume).reshape(-1,1),  
3     y = books.weight  
4 )
```

▼ LinearRegression

```
LinearRegression()
```

```
1 lm.fit(  
2     X = np.array(books.volume),  
3     y = books.weight  
4 )
```

Error: ValueError: Expected 2D array, got 1D array in array=[885 1016 1125 239 701 641 1228 412 953 1034].

Reshape your data either using `array.reshape(-1, 1)` :

```
1 lm.fit(  
2     X = books.drop(["weight", "cover"], axis=1),  
3     y = books.weight  
4 )
```

▼ LinearRegression

```
LinearRegression()
```


Model parameters

Depending on the model being used, there will be a number of parameters that can be configured when creating the model object or via the `set_params()` method.

```
1 lm.get_params()
```

```
{'copy_X': True, 'fit_intercept': True, 'n_jobs': None, 'positive': False}
```

```
1 lm.set_params(fit_intercept = False)
```

▼ LinearRegression

```
LinearRegression(fit_intercept=False)
```

```
1 lm = lm.fit(X = books[["volume"]], y = books.weight)
2 lm.intercept_
```

```
0.0
```

```
1 lm.coef_
```

```
array([0.81932487])
```

Model prediction

Once the model coefficients have been fit, it is possible to predict using the model via the `predict()` method, this method requires a matrix-like `X` as input and in the case of `LinearRegression` returns an array of predicted `y` values.

```
1 lm.predict(X = books[["volume"]])
```

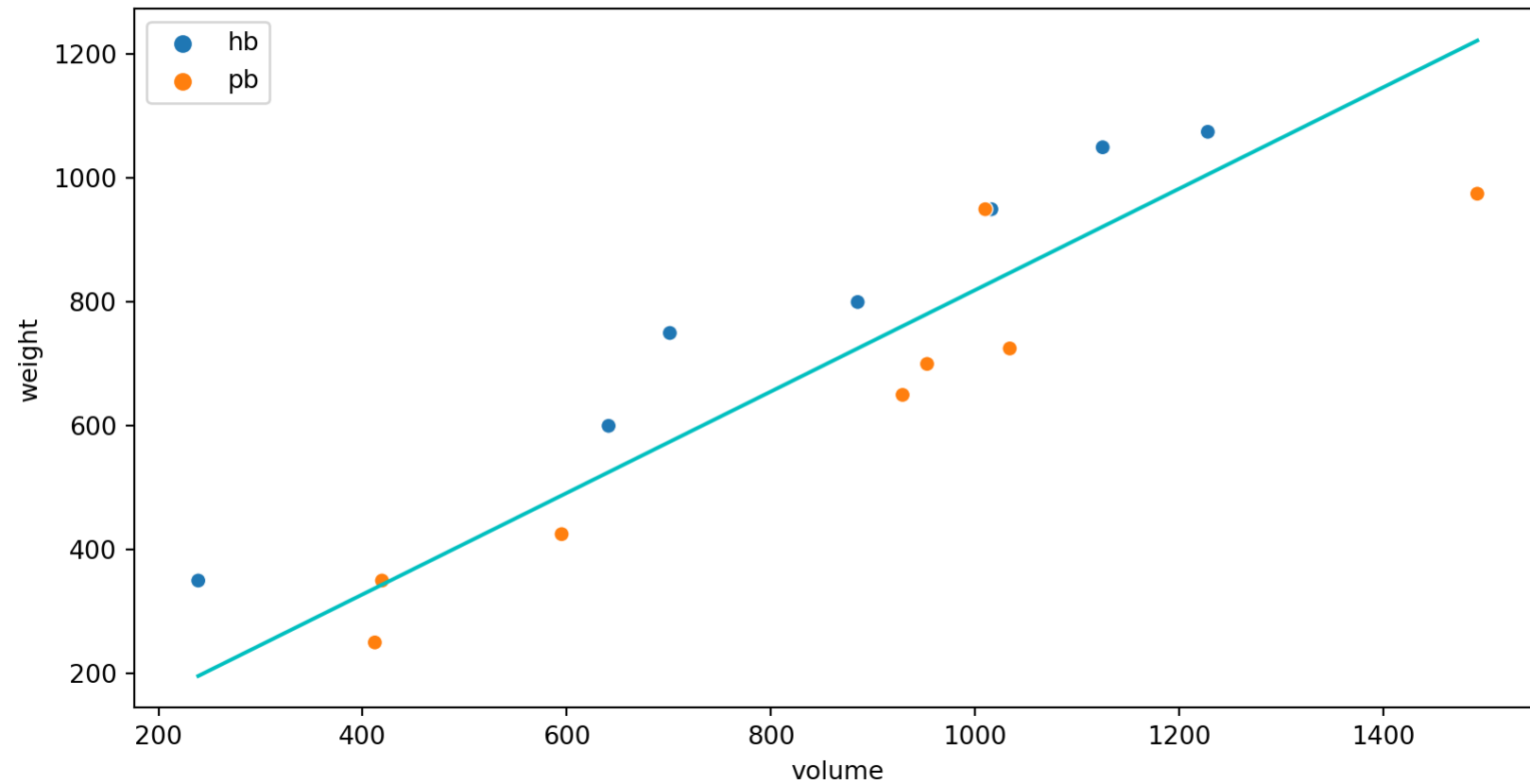
```
array([ 725.10251417,  832.43407276,  921.74048411,  195.81864507,  
       574.34673721,  525.18724472, 1006.13094621,  337.5618484 ,  
       780.81660565,  761.15280865, 1222.43271315,  343.29712253,  
       827.51812351,  487.49830048,  847.1819205  ])
```

```
1 books["weight_lm_pred"] = lm.predict(X = books[["volume"]])  
2 books
```

	volume	weight	cover	weight_lm_pred
0	885	800	hb	725.102514
1	1016	950	hb	832.434073
2	1125	1050	hb	921.740484
3	239	350	hb	195.818645
4	701	750	hb	574.346737
5	641	600	hb	525.187245
6	1228	1075	hb	1006.130946
7	412	250	pb	337.561848
8	953	700	pb	780.816606
9	929	650	pb	761.152809
10	1492	975	pb	1222.432713

11	419	350	pb	343.297123
12	1010	950	pb	827.518124
13	595	425	pb	487.498300
14	1034	725	pb	847.181921

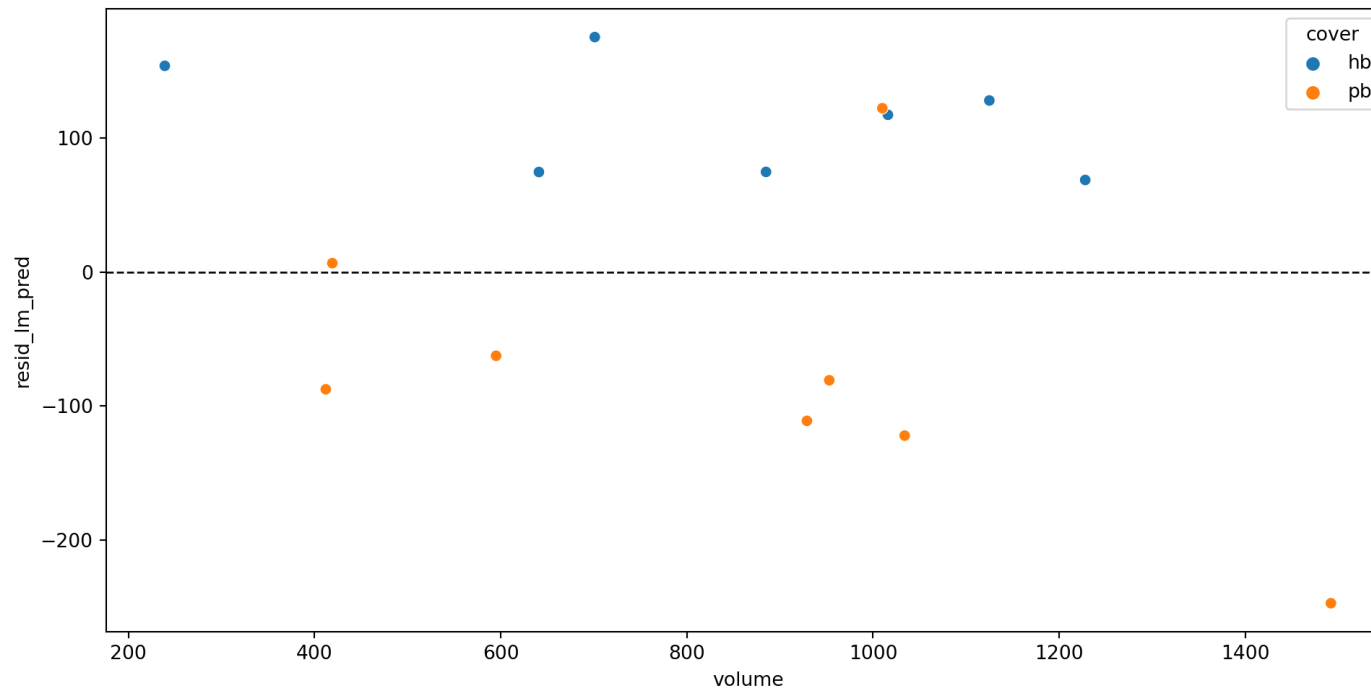
```
1 plt.figure()
2 sns.scatterplot(data=books, x="volume", y="weight", hue="cover")
3 sns.lineplot(data=books, x="volume", y="weight_lm_pred", color="c")
4 plt.show()
```



Residuals?

There is no built in functionality for calculating residuals, so this needs to be done by hand.

```
1 books["resid_lm_pred"] = books["weight"] - books["weight_lm_pred"]
2
3 plt.figure(layout="constrained")
4 ax = sns.scatterplot(data=books, x="volume", y="resid_lm_pred", hue="cover")
5 ax.axhline(c="k", ls="--", lw=1)
6 plt.show()
```



Categorical variables?

Scikit-learn expects that the model matrix be numeric before fitting,

```
1 lm = lm.fit(  
2     X = books[["volume", "cover"]],  
3     y = books.weight  
4 )
```

Error: ValueError: could not convert string to float: 'hb'

the solution here is to dummy code the categorical variables - this can be done with pandas via `pd.get_dummies()` or with a scikit-learn preprocessor.

```
1 pd.get_dummies(books[["volume", "cover"]])
```

	volume	cover_hb	cover_pb
0	885	1	0
1	1016	1	0
2	1125	1	0
3	239	1	0
4	701	1	0
5	641	1	0
6	1228	1	0
7	412	0	1
8	953	0	1
9	929	0	1
10	1492	0	1

11	419	0	1
12	1010	0	1
13	595	0	1
14	1034	0	1

What went wrong?

Do the following results look reasonable? What went wrong?

```
1 lm = LinearRegression().fit(  
2   X = pd.get_dummies(books[["volume", "cover"]]),  
3   y = books.weight  
4 )  
5  
6 lm.intercept_
```

```
105.93920788192202
```

```
1 lm.coef_
```

```
array([ 0.71795374, 92.02363569, -92.02363569])
```


Quick comparison with R

```
1 d = read.csv('data/daag_books.csv')
2 d['cover_hb'] = ifelse(d$cover == "hb", 1, 0)
3 d['cover_pb'] = ifelse(d$cover == "pb", 1, 0)
4 lm = lm(weight~volume+cover_hb+cover_pb, data=d)
5 summary(lm)
```

Call:

```
lm(formula = weight ~ volume + cover_hb + cover_pb, data = d)
```

Residuals:

Min	1Q	Median	3Q	Max
-110.10	-32.32	-16.10	28.93	210.95

Coefficients: (1 not defined because of singularities)

	Estimate	Std. Error	t value	Pr(> t)	
(Intercept)	13.91557	59.45408	0.234	0.818887	
volume	0.71795	0.06153	11.669	6.6e-08	***
cover_hb	184.04727	40.49420	4.545	0.000672	***
cover_pb	NA	NA	NA	NA	

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 78.2 on 12 degrees of freedom

Multiple R-squared: 0.9275, Adjusted R-squared: 0.9154

F-statistic: 76.73 on 2 and 12 DF, p-value: 1.455e-07

Sta 663 - Spring 2023

Avoiding co-linearity

```
1 lm = LinearRegression(  
2     fit_intercept = False  
3 ).fit(  
4     X = pd.get_dummies(books[["volume", "cover"]])  
5     y = books.weight  
6 )  
7  
8 lm.intercept_
```

0.0

```
1 lm.coef_
```

```
array([ 0.71795374, 197.96284357, 13.91557219])
```

```
1 lm.feature_names_in_
```

```
array(['volume', 'cover_hb', 'cover_pb'], dtype=object)
```

```
1 lm = LinearRegression(  
2 ).fit(  
3     X = pd.get_dummies(  
4         books[["volume", "cover"]],  
5         drop_first=True  
6     ),  
7     y = books.weight  
8 )  
9  
10 lm.intercept_
```

197.96284357271753

```
1 lm.coef_
```

```
array([ 0.71795374, -184.04727138])
```

```
1 lm.feature_names_in_
```

```
array(['volume', 'cover_pb'], dtype=object)
```

Preprocessors

Preprocessors

These are a set of transformer classes present in the `sklearn.preprocessing` submodule that are designed to help with the preparation of raw feature data into quantities more suitable for downstream modeling tools.

Like the modeling classes, they have an object oriented design that shares a common interface (methods and attributes) for bringing in data, transforming it, and returning it.

OneHotEncoder

For dummy coding we can use the `OneHotEncoder` preprocessor, the default is to use one hot encoding but standard dummy coding can be achieved via the `drop` parameter.

```
1 from sklearn.preprocessing import OneHotEncoder
```

```
1 enc = OneHotEncoder(sparse_output=False)
2 enc.fit(X = books[["cover"]])
```

▼ OneHotEncoder

OneHotEncoder(sparse_output=False)

```
1 enc.transform(X = books[["cover"]])
```

```
array([[1., 0.],
       [1., 0.],
       [1., 0.],
       [1., 0.],
       [1., 0.],
       [1., 0.],
       [1., 0.],
       [1., 0.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.],
       [0., 1.]])
```

```
1 enc = OneHotEncoder(sparse_output=False, drop="first")
2 enc.fit_transform(X = books[["cover"]])
```

```
array([[0.],
       [0.],
       [0.],
       [0.],
       [0.],
       [0.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.],
       [1.]])
```

Other useful bits

```
1 enc.get_feature_names_out()
```

```
array(['cover_hb', 'cover_pb'], dtype=object)
```

```
1 f = enc.transform(X = books[["cover"]])  
2 f
```

```
array([[1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [1., 0.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.],  
       [0., 1.]])
```

```
1 enc.inverse_transform(f)
```

```
array([[ 'hb' ],  
       [ 'hb' ],  
       [ 'hb' ],  
       [ 'hb' ],  
       [ 'hb' ],  
       [ 'hb' ],  
       [ 'hb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ],  
       [ 'pb' ]], dtype=object)
```

A cautionary note

Unlike `pd.get_dummies()` it is not safe to use `OneHotEncoder` with both numerical and categorical features, as the former will also be transformed.

```
1 enc = OneHotEncoder(sparse_output=False)
2 X = enc.fit_transform(X = books[["volume", "cover"]])
3 pd.DataFrame(data=X, columns = enc.get_feature_names_out())
```

	volume_239	volume_412	volume_419	...	volume_1492	cover_hb	cover_pb
0	0.0	0.0	0.0	...	0.0	1.0	0.0
1	0.0	0.0	0.0	...	0.0	1.0	0.0
2	0.0	0.0	0.0	...	0.0	1.0	0.0
3	1.0	0.0	0.0	...	0.0	1.0	0.0
4	0.0	0.0	0.0	...	0.0	1.0	0.0
5	0.0	0.0	0.0	...	0.0	1.0	0.0
6	0.0	0.0	0.0	...	0.0	1.0	0.0
7	0.0	1.0	0.0	...	0.0	0.0	1.0
8	0.0	0.0	0.0	...	0.0	0.0	1.0
9	0.0	0.0	0.0	...	0.0	0.0	1.0
10	0.0	0.0	0.0	...	1.0	0.0	1.0
11	0.0	0.0	1.0	...	0.0	0.0	1.0
12	0.0	0.0	0.0	...	0.0	0.0	1.0
13	0.0	0.0	0.0	...	0.0	0.0	1.0
14	0.0	0.0	0.0	...	0.0	0.0	1.0

[15 rows x 17 columns]

Putting it together

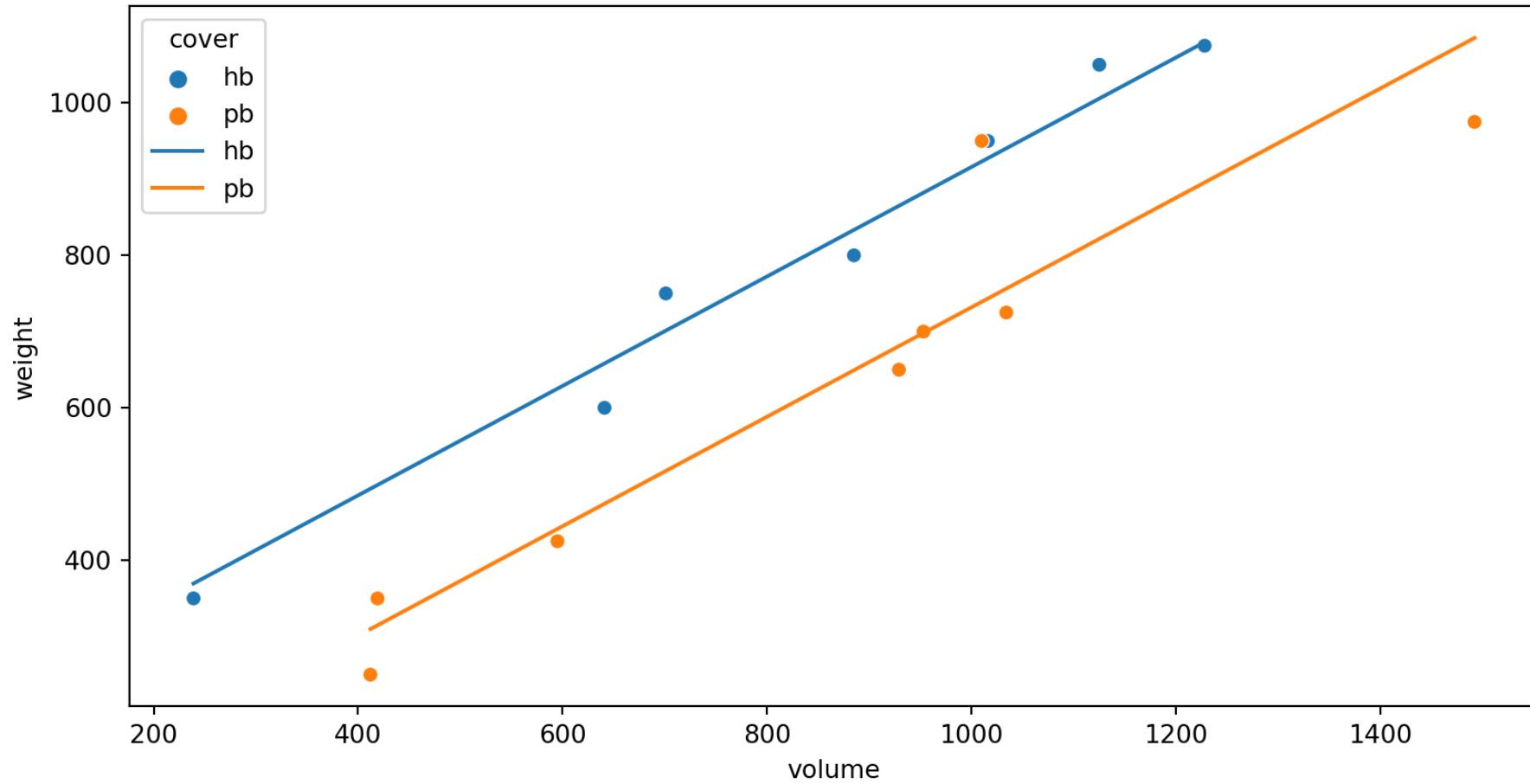
```
1 cover = OneHotEncoder(  
2     sparse_output=False  
3 ).fit_transform(  
4     books[["cover"]]  
5 )  
6 X = np.c_[books.volume, cover]  
7  
8 lm2 = LinearRegression(  
9     fit_intercept=False  
10 ).fit(  
11     X = X,  
12     y = books.weight  
13 )  
14  
15 lm2.coef_
```

```
array([ 0.71795374, 197.96284357, 13.91557219])
```

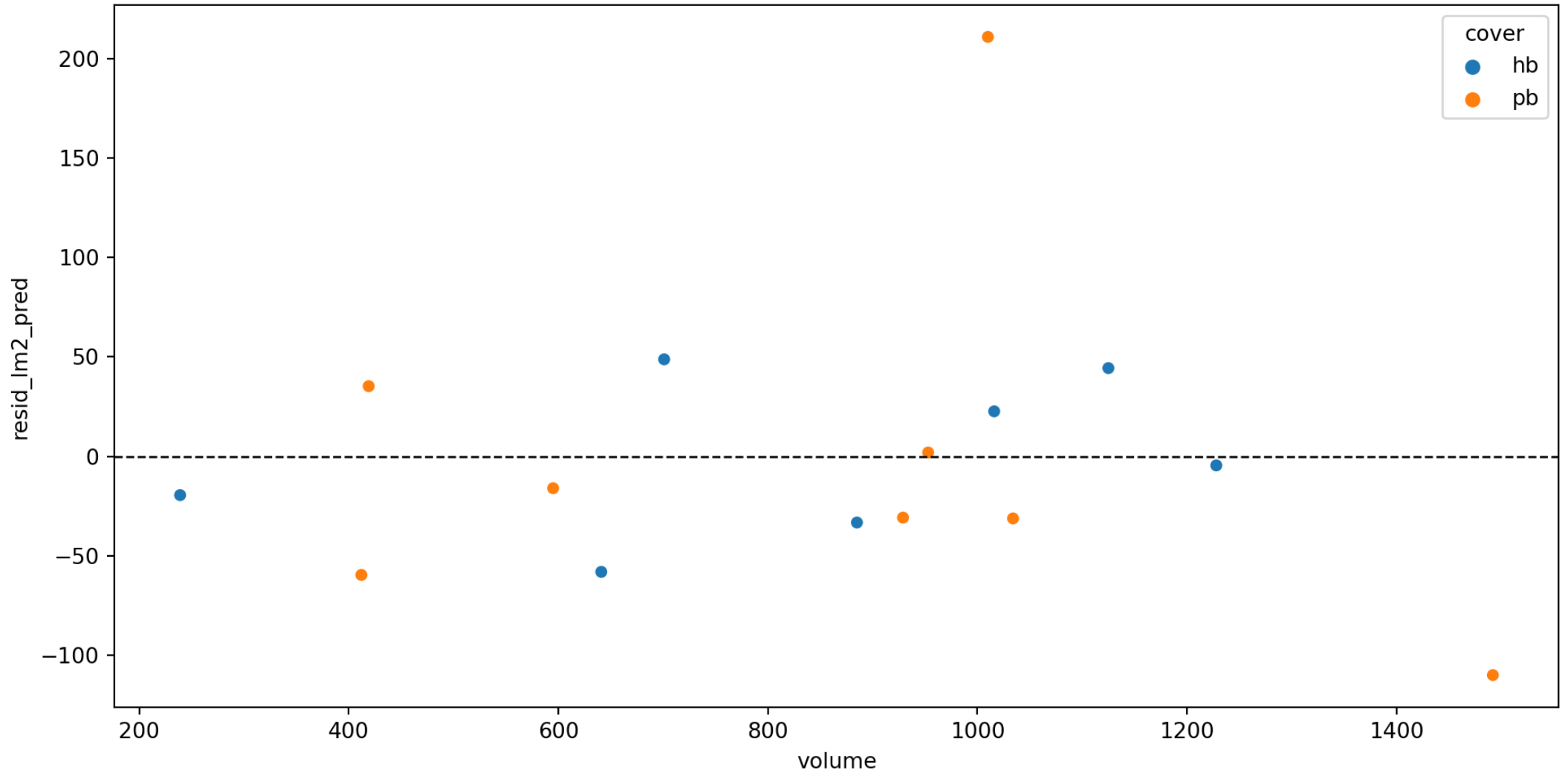
```
1 books["weight_lm2_pred"] = lm2.predict(X=X)  
2 books.drop(["weight_lm_pred", "resid_lm_pred"],
```

	volume	weight	cover	weight_lm2_pred
0	885	800	hb	833.351907
1	1016	950	hb	927.403847
2	1125	1050	hb	1005.660805
3	239	350	hb	369.553788
4	701	750	hb	701.248418
5	641	600	hb	658.171193
6	1228	1075	hb	1079.610041
7	412	250	pb	309.712515
8	953	700	pb	698.125490
9	929	650	pb	680.894600
10	1492	975	pb	1085.102558
11	419	350	pb	314.738191
12	1010	950	pb	739.048853
13	595	425	pb	441.098050
14	1034	725	pb	756.279743

Model fit



Model residuals



Model performance

Scikit-learn comes with a number of builtin functions for measuring model performance in the `sklearn.metrics` submodule - these are generally just functions that take the vectors `y_true` and `y_pred` and return a scalar score.

```
1 from sklearn.metrics import mean_squared_error, r2_score
```

```
1 r2_score(books.weight, books.weight_lm_pred)
```

0.7800969547785038

```
1 mean_squared_error(  
2     books.weight, books.weight_lm_pred  
3 )
```

14833.68208377448

```
1 mean_squared_error(  
2     books.weight, books.weight_lm_pred,  
3     squared=False  
4 )
```

121.79360444528473

```
1 r2_score(books.weight, books.weight_lm2_pred)
```

0.9274775756821679

```
1 mean_squared_error(  
2     books.weight, books.weight_lm2_pred  
3 )
```

4892.040422595093

```
1 mean_squared_error(  
2     books.weight, books.weight_lm2_pred,  
3     squared=False  
4 )
```

69.94312276839727

Exercise 1

Create and fit a model for the `books` data that includes an interaction effect between `volume` and `cover`.

You will need to do this manually with `pd.getdummies()` and some additional data munging.

The data can be read into pandas with,

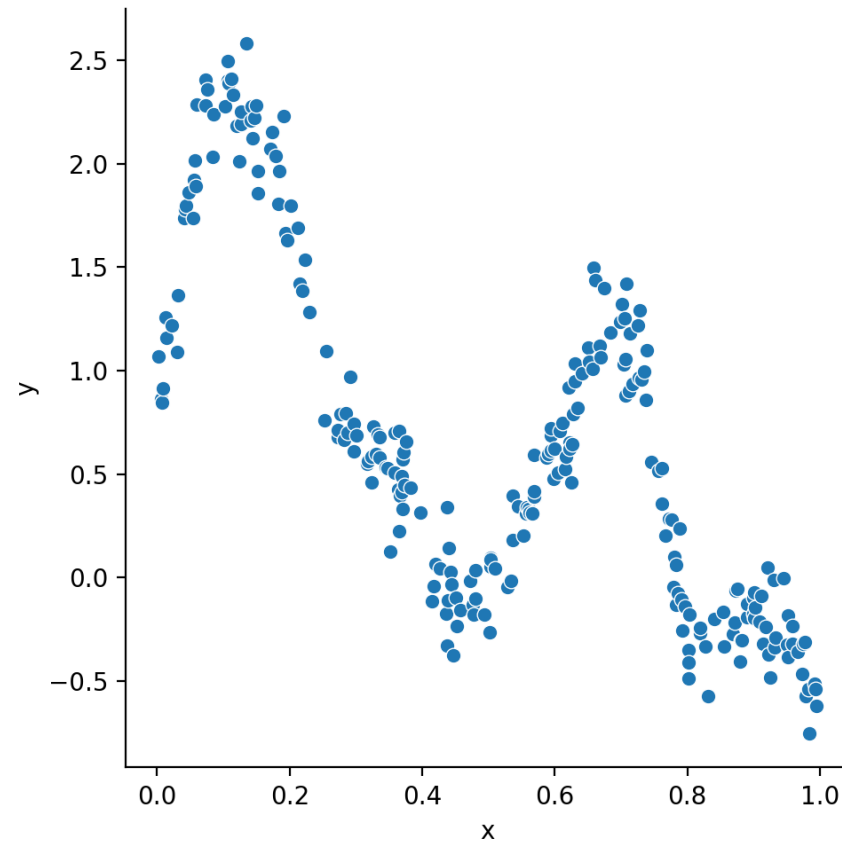
```
1 books = pd.read_csv(  
2     "https://sta663-sp23.github.io/slides/data/daag_books.csv"  
3 )
```

Other transformers

Polynomial regression

We will now look at another flavor of regression model, that involves preprocessing and a hyperparameter - namely polynomial regression.

```
1 df = pd.read_csv("data/gp.csv")  
2 sns.relplot(data=df, x="x", y="y")
```



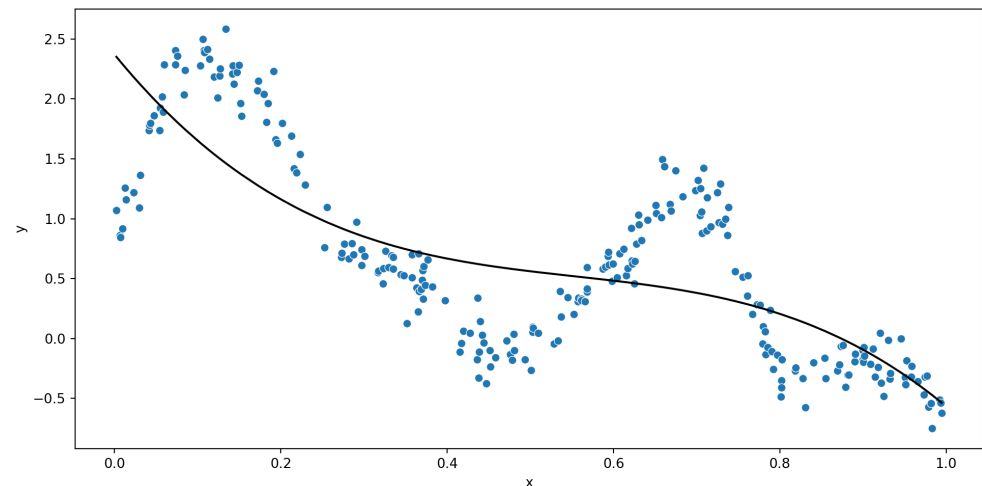
By hand

It is certainly possible to construct the necessary model matrix by hand (or even use a function to automate the process), but this is less than desirable generally - particularly if we want to do anything fancy (e.g. cross validation)

```
1 X = np.c_[
2     np.ones(df.shape[0]),
3     df.x,
4     df.x**2,
5     df.x**3
6 ]
7
8 plm = LinearRegression(
9     fit_intercept = False
10 ).fit(
11     X=X, y=df.y
12 )
13
14 plm.coef_
```

```
array([ 2.36985684, -8.49429068, 13.95066369, -8.3921
```

```
1 df["y_pred"] = plm.predict(X=X)
2
3 plt.figure(layout="constrained")
4 sns.scatterplot(data=df, x="x", y="y")
5 sns.lineplot(data=df, x="x", y="y_pred", color="red")
6 plt.show()
```



PolynomialFeatures

This is another transformer class from `sklearn.preprocessing` that simplifies the process of constructing polynomial features for your model matrix. Usage is similar to that of `OneHotEncoder`.

```
1 from sklearn.preprocessing import PolynomialFeatures
2 X = np.array(range(6)).reshape(-1,1)
```

```
1 pf = PolynomialFeatures(degree=3)
2 pf.fit(X)
```

▼ PolynomialFeatures

PolynomialFeatures(degree=3)

```
1 pf.transform(X)
```

```
array([[ 1.,  0.,  0.,  0.],
       [ 1.,  1.,  1.,  1.],
       [ 1.,  2.,  4.,  8.],
       [ 1.,  3.,  9., 27.],
       [ 1.,  4., 16., 64.],
       [ 1.,  5., 25.,125.]])
```

```
1 pf.get_feature_names_out()
```

```
array(['1', 'x0', 'x0^2', 'x0^3'], dtype=object)
```

```
1 pf = PolynomialFeatures(
2     degree=2, include_bias=False
3 )
4 pf.fit_transform(X)
```

```
array([[ 0.,  0.],
       [ 1.,  1.],
       [ 2.,  4.],
       [ 3.,  9.],
       [ 4., 16.],
       [ 5., 25.]])
```

```
1 pf.get_feature_names_out()
```

```
array(['x0', 'x0^2'], dtype=object)
```

Interactions

If the feature matrix X has more than one column then `PolynomialFeatures` transformer will include interaction terms with total degree up to `degree`.

```
1 X.reshape(-1, 2)
```

```
array([[0, 1],
       [2, 3],
       [4, 5]])
```

```
1 pf = PolynomialFeatures(
2   degree=3, include_bias=False
3 )
4 pf.fit_transform(
5   X.reshape(-1, 2)
6 )
```

```
array([[ 0.,  1.,  0.,  0.,  1.,  0.,  0.,  0.,  0.,  0.],
       [ 2.,  3.,  4.,  6.,  9.,  8., 12., 18., 27., 27.],
       [ 4.,  5., 16., 20., 25., 64., 80., 100., 125., 125.]])
```

```
1 pf.get_feature_names_out()
```

```
array(['x0', 'x1', 'x0^2', 'x0 x1', 'x1^2', 'x0^3',
       'x1^3'], dtype=object)
```

```
1 X.reshape(-1, 3)
```

```
array([[0, 1, 2],
       [3, 4, 5]])
```

```
1 pf = PolynomialFeatures(
2   degree=2, include_bias=False
3 )
4 pf.fit_transform(
5   X.reshape(-1, 3)
6 )
```

```
array([[ 0.,  1.,  2.,  0.,  0.,  0.,  1.,  2.,  4.],
       [ 3.,  4.,  5.,  9., 12., 15., 16., 20., 25.]])
```

```
1 pf.get_feature_names_out()
```

```
array(['x0', 'x1', 'x2', 'x0^2', 'x0 x1', 'x0 x2', 'x1^2',
       'x2^2'], dtype=object)
```

Modeling with PolynomialFeatures

```
1 def poly_model(X, y, degree):
2     X = PolynomialFeatures(
3         degree=degree, include_bias=False
4     ).fit_transform(
5         X=X
6     )
7     y_pred = LinearRegression(
8     ).fit(
9         X=X, y=y
10    ).predict(
11        X
12    )
13    return mean_squared_error(y, y_pred, squared=False)
```

```
1 poly_model(X = df[["x"]], y = df.y, degree = 2)
```

0.5449418707295371

```
1 poly_model(X = df[["x"]], y = df.y, degree = 3)
```

0.5208157900621085

```
1 degrees = range(1,10)
2 rmsees = [
3     poly_model(X=df[["x"]], y=df.y, degree=d)
4     for d in degrees
5 ]
6 sns.relplot(x=degrees, y=rmsees)
```

