# 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

This is one of several other "scikits" (e.g. scikit-image) which are scientific toolboxes built on top of scipy. For a

# 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.
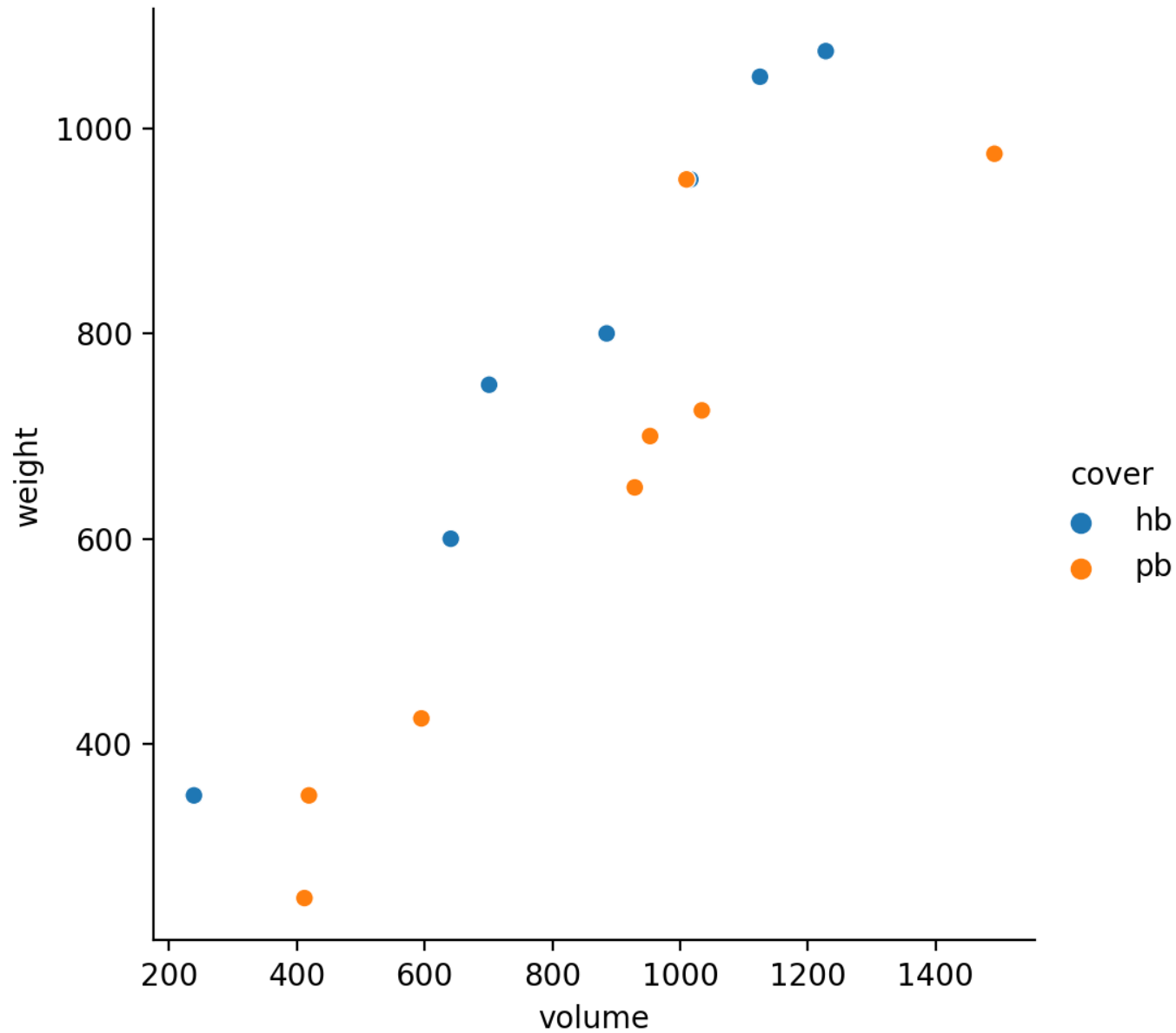
The included columns are:

- `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    |

These data come from the `allbacks` data set from the `DAAG` package in R

```
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  )
```

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

```
Error: ValueError: Expected 2D array, got 1D array i
array=[ 885 1016 1125  239  701  641 1228  412  953
 1034].
Reshape your data either using array.reshape(-1, 1)
```

```
Error: ValueError: Expected 2D array, got 1D array i
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  )
```

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

▼ LinearRegression

LinearRegression()

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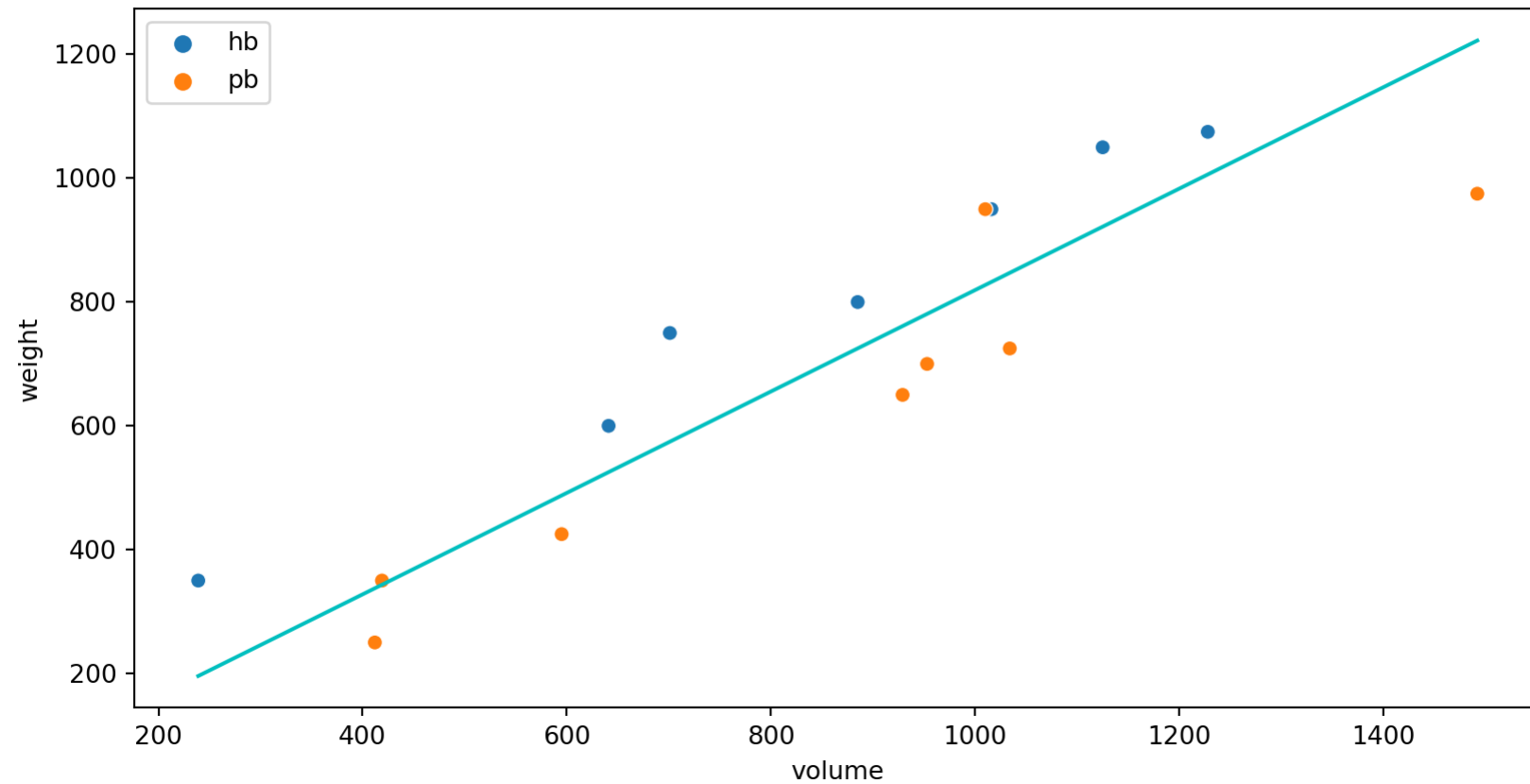
```
1  lm.predict(X = books[["volume"]])
```

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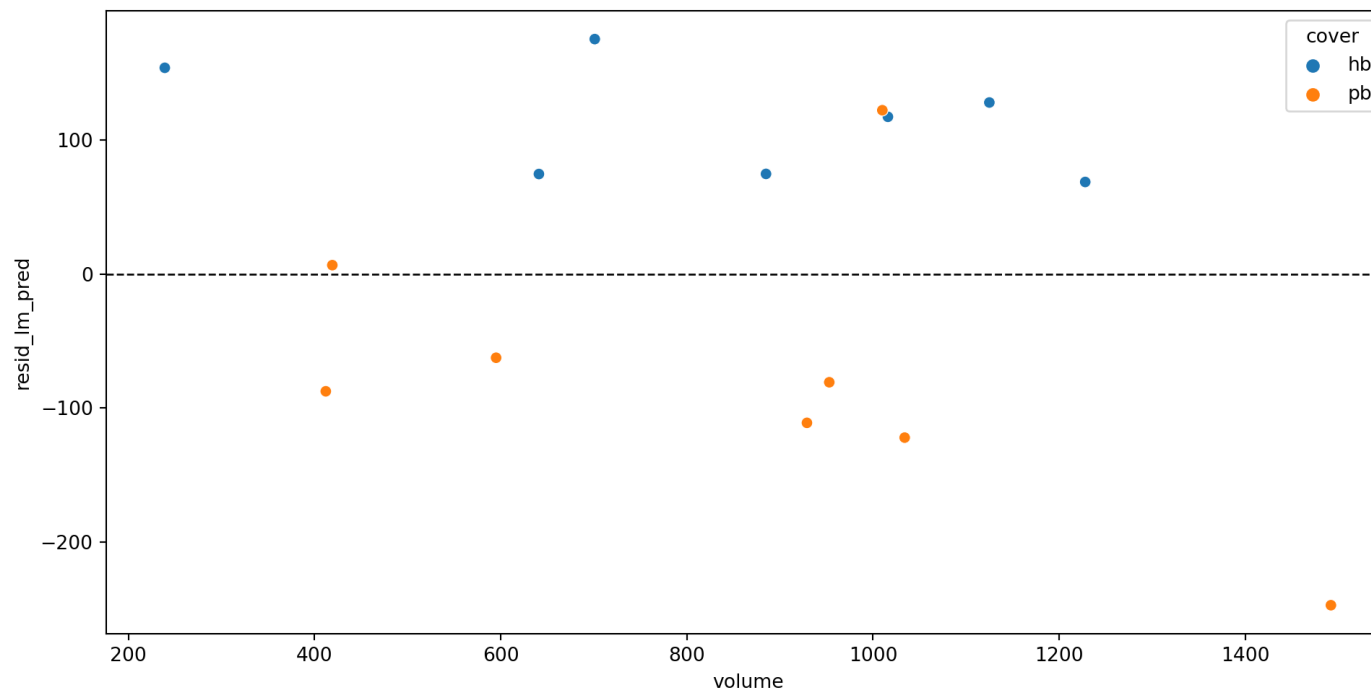
```
1  plt.figure()
2  sns.scatterplot(data=books, x="volume", y="weight", hue="cover")
3  sns.lineplot(data=books, x="volume", y="weight_lm_pred", color="c")
4  plt.show()
```

# Residuals?

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

```python
books["resid_lm_pred"] = books["weight"] - books["weight_lm_pred"]

plt.figure(layout="constrained")
ax = sns.scatterplot(data=books, x="volume", y="resid_lm_pred", hue="cover")
ax.axhline(c="k", ls="--", lw=1)
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
```

# 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=obje
```

```
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"]])
```

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

▼              OneHotEncoder

OneHotEncoder(sparse_output=False)

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

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

```
array([[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.]])
```

# 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.],
       [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']], 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]
```
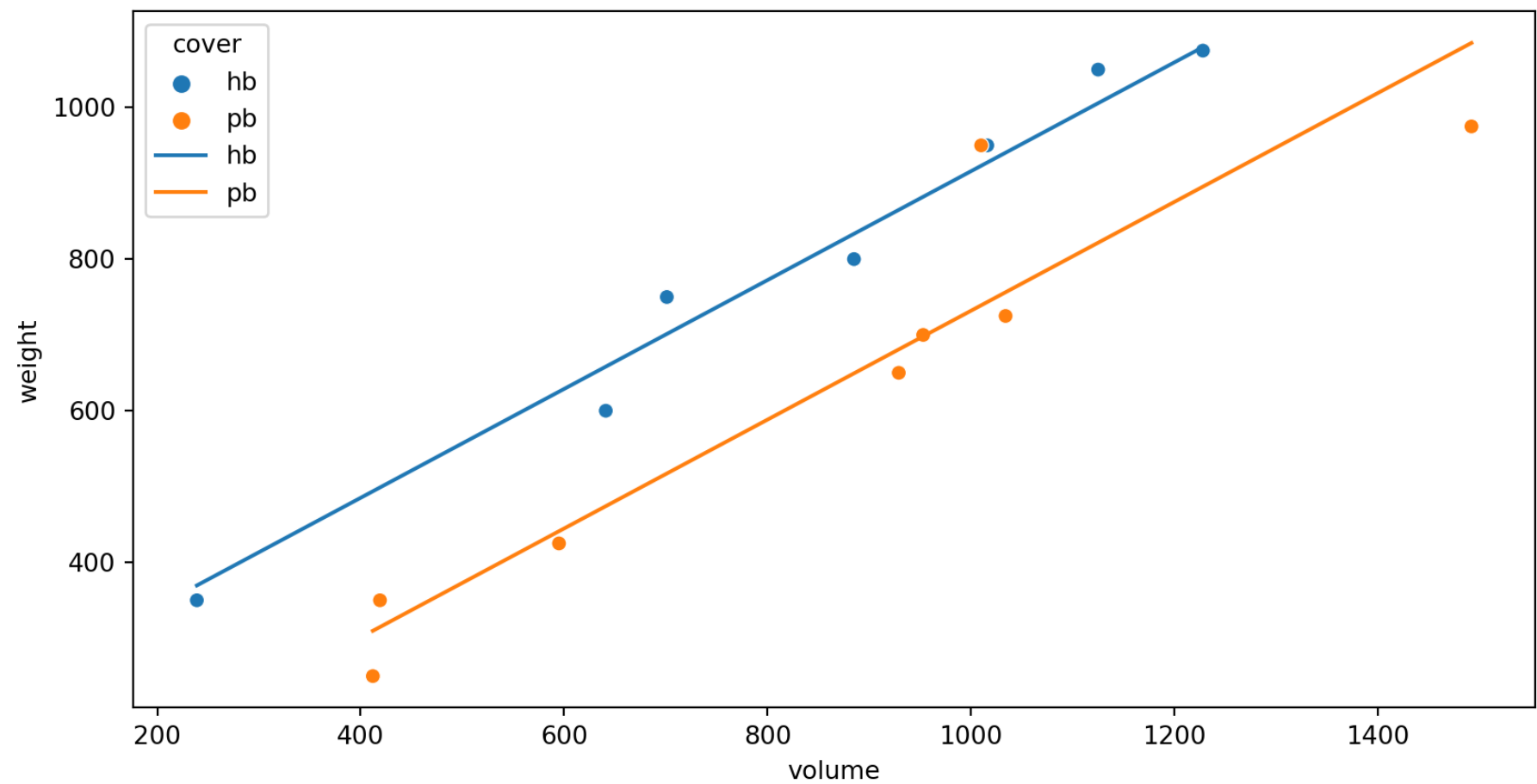
# Putting it together

```
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_
```

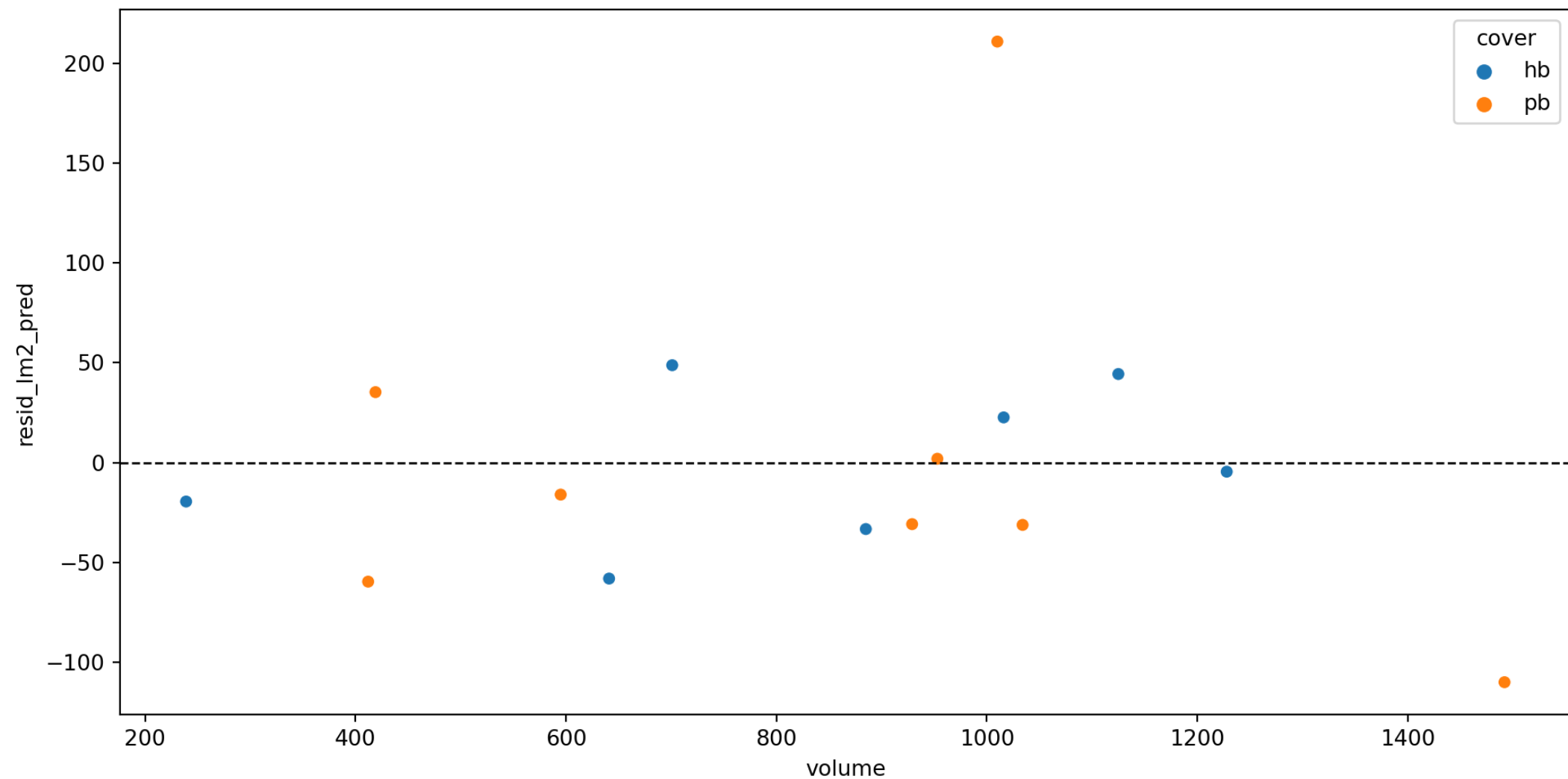array([  0.71795374, 197.96284357,  13.91557219])

```
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      |

We'll see a more elegant way of doing this in the near future

# 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

See API Docs for a list of available metrics

# 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.

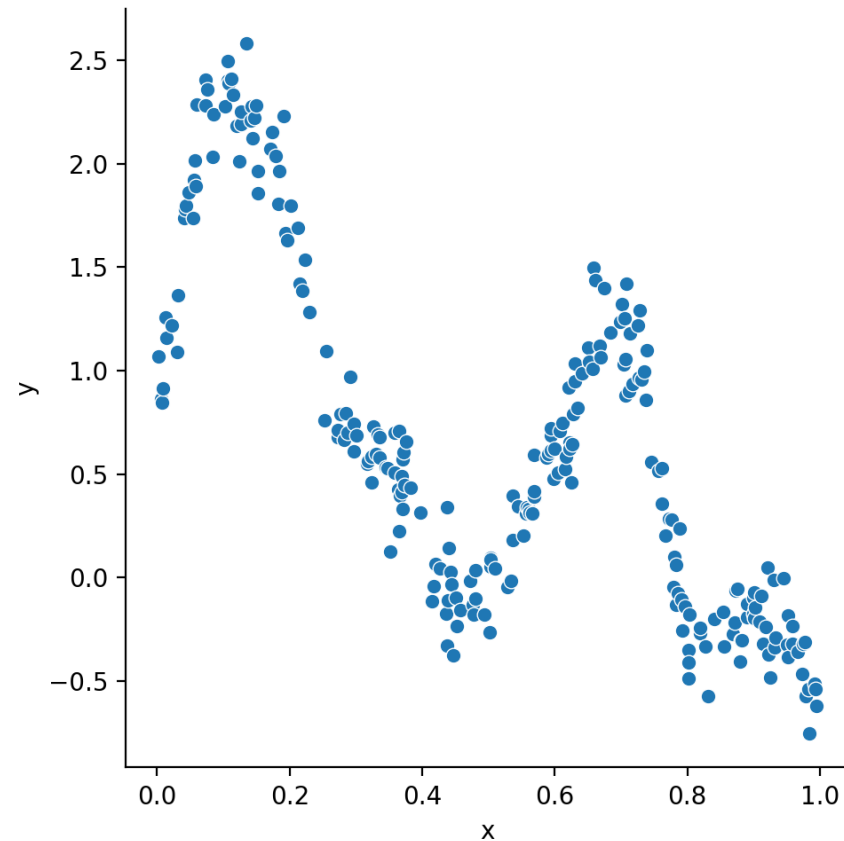The data can be read into pandas with,

```
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")
```
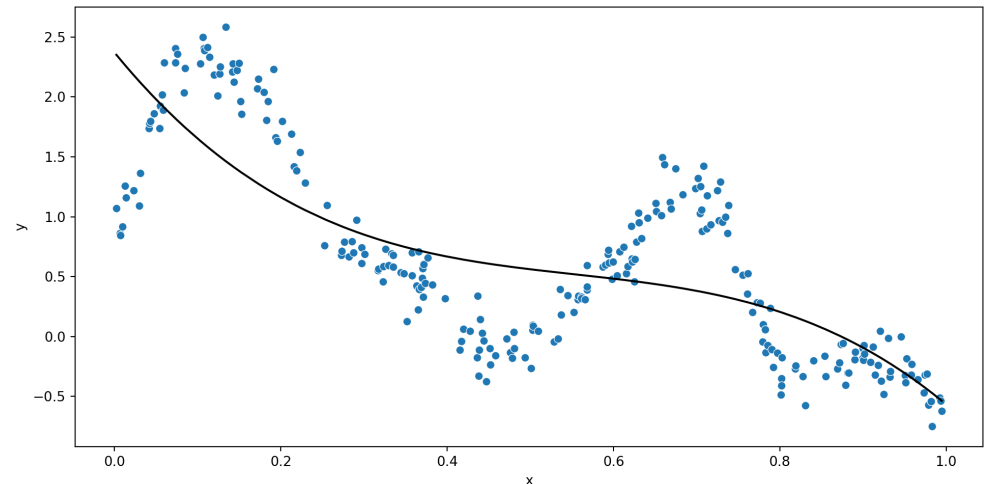
# By hand

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 then desirable generally - particularly if we want to do anything fancy (e.g. cross validation)

```python
 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, -8.392

```python
 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="
 6  plt.show()
```

# PolynomialFeatures

This is another transformer class from `sklearn.preprocessing` that simplifies the process of constructing polynormial 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.,    (
       [  2.,    3.,    4.,    6.,    9.,    8.,   12.,   18
       [  4.,    5.,   16.,   20.,   25.,   64.,   80.,  100
```

```
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', 'x
       '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=F
```
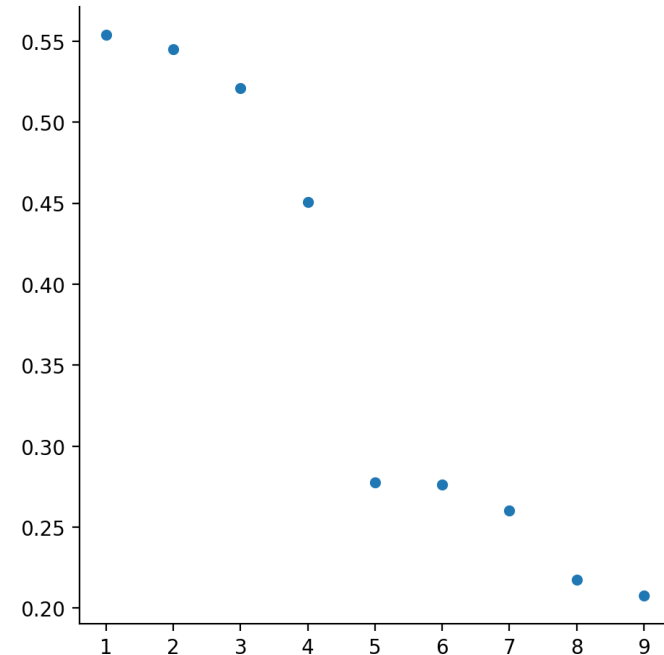
```
1  poly_model(X = df[["x"]], y = df.y, degree = 2)
```
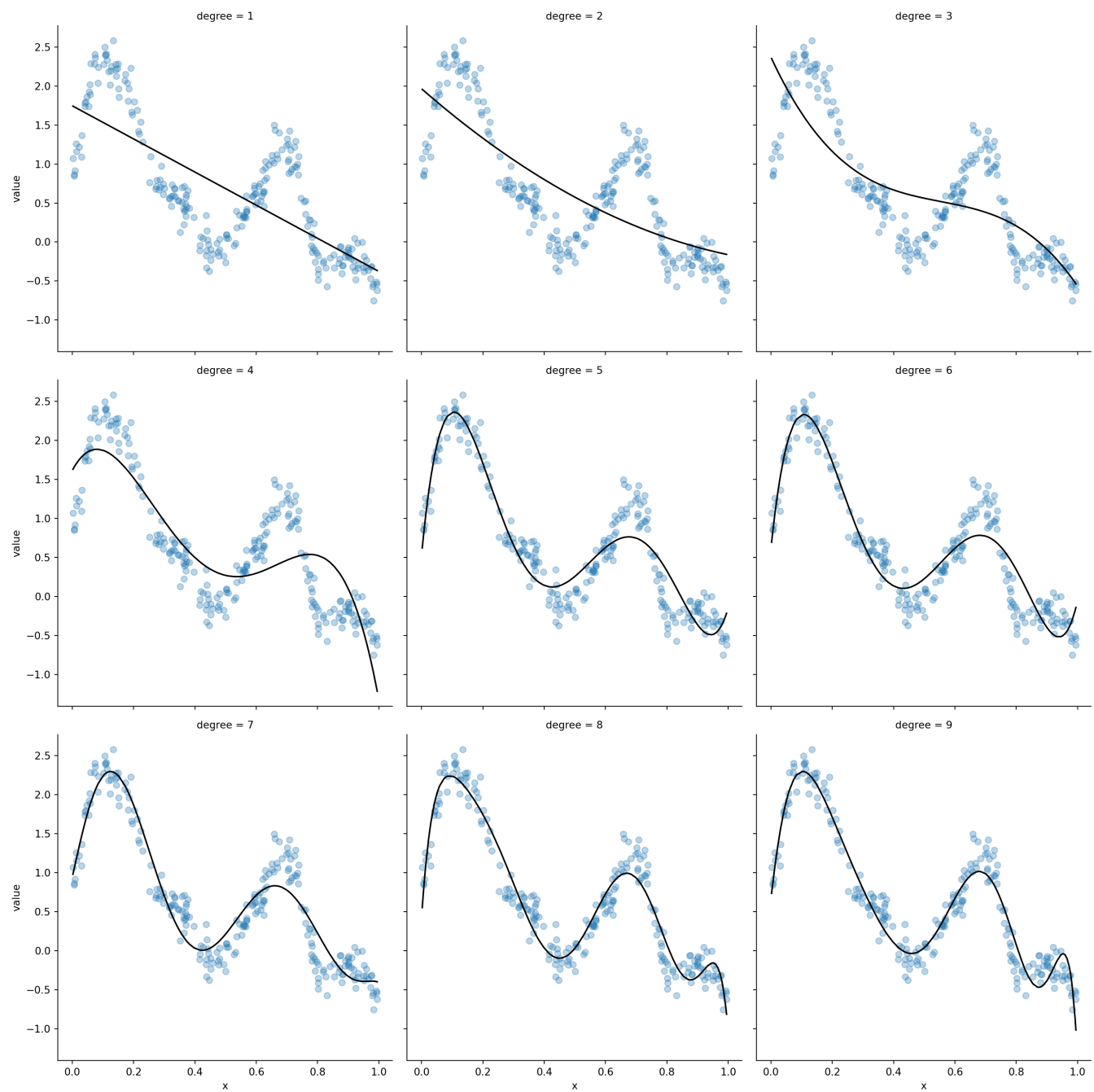
0.5449418707295371

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

0.5208157900621085

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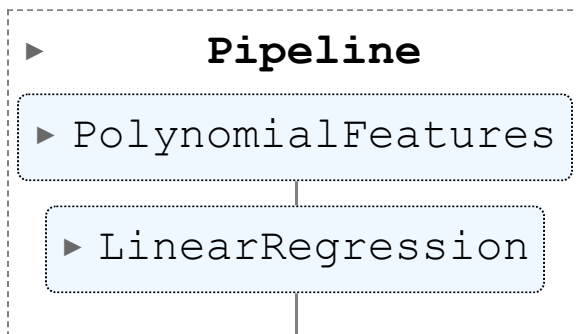
# Pipelines

You may have noticed that `PolynomialFeatures` takes a model matrix as input and returns a new model matrix as output which is then used as the input for `LinearRegression`. This is not an accident, and by structuring the library in this way sklearn is designed to enable the connection of these steps together, into what sklearn calls a *pipeline*.

```python
from sklearn.pipeline import make_pipeline

p = make_pipeline(
    PolynomialFeatures(degree=4),
    LinearRegression()
)
p
```

```
▶         Pipeline
┌──────────────────────────┐
│ ▶ PolynomialFeatures     │
└──────────────────────────┘
   ┌───────────────────────┐
   │ ▶ LinearRegression    │
   └───────────────────────┘
```
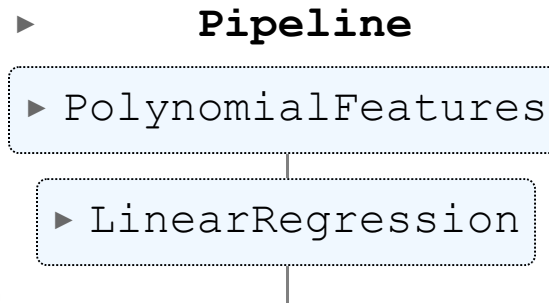
# Using Pipelines

Once constructed, this object can be used just like our previous `LinearRegression` model (i.e. fit to our data and then used for prediction)

```
1  p = p.fit(X = df[["x"]], y = df.y)
2  p
```

```
▶       **Pipeline**
  ▶ PolynomialFeatures
    ▶ LinearRegression
```

```
1  p.predict(X = df[["x"]])
```

```
array([ 1.6295693 ,  1.65734929,  1.6610466 ,  1.6777
        1.70475286,  1.75280126,  1.78471392,  1.7904
        1.82966357,  1.83376043,  1.84494343,  1.8600
        1.86619112,  1.86837909,  1.87065283,  1.8841
        1.88527174,  1.88577463,  1.88544367,  1.8689
        1.86252922,  1.86047349,  1.85377801,  1.8493
        1.82623453,  1.82024199,  1.81799793,  1.7976
        1.77034143,  1.76574288,  1.75371272,  1.7438
        1.73356954,  1.65527727,  1.64812184,  1.6186
        1.5960389 ,  1.56080881,  1.55036459,  1.5400
        1.45096594,  1.43589836,  1.41886389,  1.3942
        1.23072992,  1.21355164,  1.11776117,  1.1152
        1.06449719,  1.04672121,  1.03662739,  1.0140
        0.98081577,  0.96176797,  0.87491417,  0.8711
        0.84171166,  0.82875003,  0.8085086 ,  0.7916
        0.78078036,  0.73538157,  0.7181484 ,  0.7004
        0.67229069,  0.64782899,  0.64050946,  0.6372
        0.62323271,  0.61965166,  0.61705548,  0.6141
        0.60347713,  0.5909255 ,  0.566617  ,  0.5090
        0.44177711,  0.43291379,  0.40957833,  0.3848
        0.38067928,  0.3791518 ,  0.37610476,  0.3693
        0.35806518,  0.3475729 ,  0.3466828 ,  0.3333
        0.3006981 ,  0.29675876,  0.29337641,  0.2933
```

```
1  plt.figure(layout="constrained")
2  sns.scatterplot(data=df, x="x", y="y")
3  sns.lineplot(x=df.x, y=p.predict(X = df[["x"]]),
4  plt.show()
```

# Model coefficients (or other attributes)

The attributes of steps are not directly accessible, but can be accessed via `steps` or `named_steps` attributes,

```
1  p.coef_
```

Error: AttributeError: 'Pipeline' object has no attribute 'coef_'

```
1  p.named_steps["linearregression"].intercept_
```

1.6136636604768615

```
1  p.steps[1][1].coef_
```

array([  0.        ,   7.39051417, -57.67175293, 102.72227443,
       -55.38181361])

```
1  p.steps
```

[('polynomialfeatures', PolynomialFeatures(degree=4)), ('linearregression', LinearRegression())]

```
1  p.steps[0][1].get_feature_names_out()
```

array(['1', 'x', 'x^2', 'x^3', 'x^4'], dtype=object)

Anyone notice a problem?

# What about step parameters?

By accessing each step we can adjust their parameters (via `set_params()`),
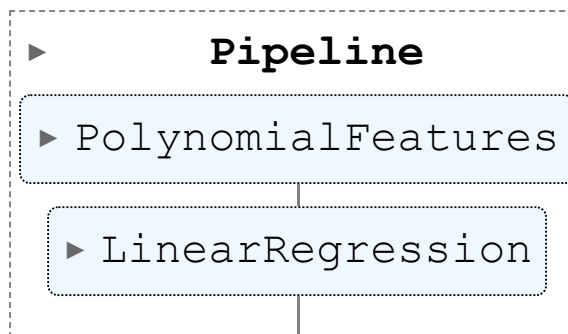
```
1  p.named_steps[
2     "linearregression"
3  ].get_params()
```

`{'copy_X': True, 'fit_intercept': True, 'n_jobs': No`

```
1  p.named_steps[
2     "linearregression"
3  ].set_params(
4     fit_intercept=False
5  )
```

▼          LinearRegression

LinearRegression(fit_intercept=False)

```
1  p.fit(X = df[["x"]], y = df.y)
```

```
▶          Pipeline

▶ PolynomialFeatures

▶ LinearRegression
```

```
1  p.named_steps["linearregression"].intercept_
```

0.0

```
1  p.named_steps["linearregression"].coef_
```

```
array([  1.61366366,   7.39051417, -57.67175293, 102.
       -55.38181361])
```
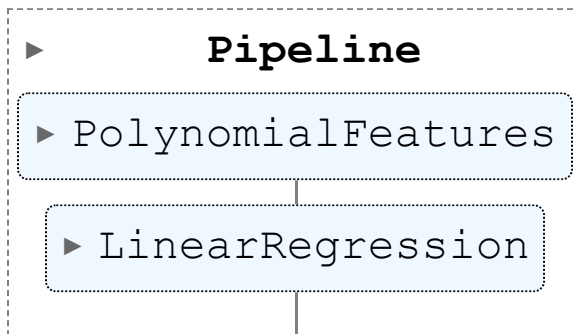
# Pipeline parameter names

These parameters can also be directly accessed at the pipeline level, note how the names are constructed:
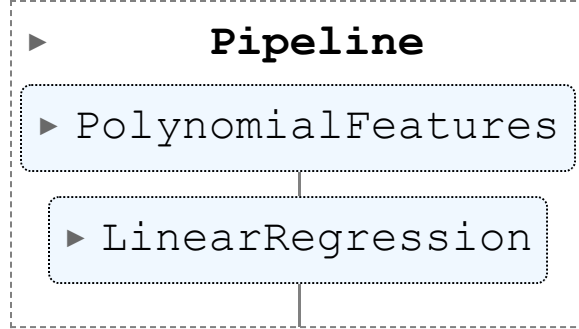
```
1  p.get_params()
```

```
{'memory': None, 'steps': [('polynomialfeatures', PolynomialFeatures(degree=4)), ('linearregression', Linea
```

```
1  p.set_params(
2     linearregression__fit_intercept=True,
3     polynomialfeatures__include_bias=False
4  )
```

```
►          Pipeline
► PolynomialFeatures

  ► LinearRegression
```

```
1  p.fit(X = df[["x"]], y = df.y)
```

```
▶        Pipeline
▶ PolynomialFeatures
    ▶ LinearRegression
```

```
1  p.named_steps["linearregression"].intercept_
```

1.6136636604768375

```
1  p.named_steps["linearregression"].coef_
```

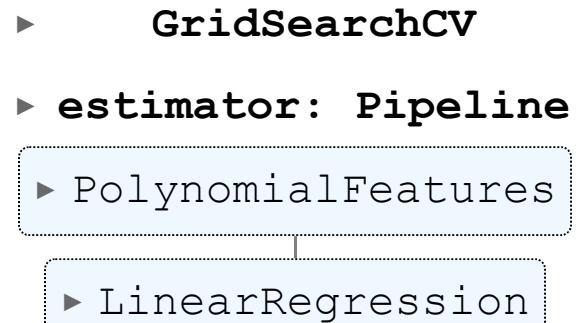array([  7.39051417, -57.67175293, 102.72227443, -55.38181361])

# Tuning parameters

We've already seen a manual approach to tuning models over the degree parameter, scikit-learn also has built in tools to aide with this process. Here we will leverage `GridSearchCV` to tune the degree parameter in our pipeline.

```python
from sklearn.model_selection import GridSearchCV, KFold

p = make_pipeline(
    PolynomialFeatures(include_bias=True),
    LinearRegression(fit_intercept=False)
)


grid_search = GridSearchCV(
  estimator = p,
  param_grid = {"polynomialfeatures__degree": range(1,10)},
  scoring = "neg_root_mean_squared_error",
  cv = KFold(shuffle=True)
)
```

Much more detail on this next time – including the proper way to do cross-validation

# Preview - Performing a grid search

```
1  grid_search.fit(X = df[["x"]], y = df.y)
```

```
▶       GridSearchCV

▶ estimator: Pipeline

 ▶ PolynomialFeatures

   ▶ LinearRegression
```

```
1  grid_search.best_index_
```

8

```
1  grid_search.best_params_
```

{'polynomialfeatures__degree': 9}

```
1  grid_search.best_score_
```

−0.21584378186687297

# cv_results_

```
1  grid_search.cv_results_["mean_test_score"]
```

```
array([-0.55788374, -0.54949818, -0.52704221, -0.45732017, -0.27965894,
       -0.27855276, -0.26232186, -0.22311198, -0.21584378])
```

```
1  grid_search.cv_results_["rank_test_score"]
```

```
array([9, 8, 7, 6, 5, 4, 3, 2, 1], dtype=int32)
```

```
1  grid_search.cv_results_["mean_fit_time"]
```

```
array([0.00098572, 0.00087113, 0.00086641, 0.00087099, 0.0008688 ,
       0.00087166, 0.0008811 , 0.00088105, 0.00092745])
```

```
1  grid_search.cv_results_.keys()
```

```
dict_keys(['mean_fit_time', 'std_fit_time', 'mean_score_time', 'std_score_time', 'param_polynomialfeatures__
```