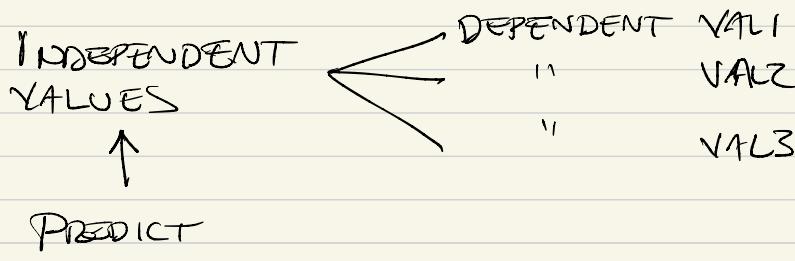


MODEL DEVELOPMENT

- SIMPLE AND MULTIPLE LINEAR REGRESSION
- MODEL EVALUATION USING VISUALIZATION POLYNOMIAL REGRESSION AND PIPELINES
- R-SQUARED AND MSE IN-SAMPLE EVALUATION
- PREDICTION AND DECISION MAKING

MODEL ESTIMATOR



LINEAR

1 INDEPENDENT VARIABLE

MULTIPLE LINEAR

2 OR MORE INDEPENDENT VARIABLES

SINGLE LINEAR REGRESSION

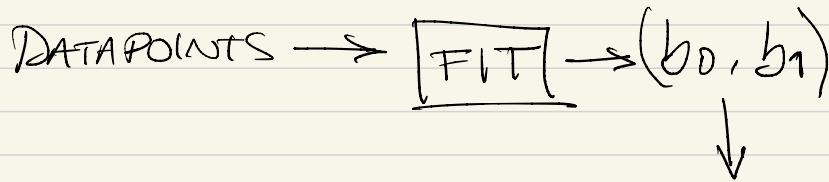
$$y = b_0 + b_1 x$$

↑ ↗
INTERCEPT SLOPE

A LOT OF MATHEMATICS TO DETERMINE THESE VALUES

DATA IN NUMPY ARRAYS
 X AND Y

PROCESS



from sklearn.linear_model
import LinearRegression
lm = LinearRegression()

$$\hat{y} = b_0 + b_1 x$$

$$X = df[["indVar"]]$$
$$Y = df[["depVar"]]$$

lm.intercept
lm.coef

lm.fit(X, Y) # Find b_0 and b_1
 \hat{y} = lm.predict(X)

MULTIPLE LINEAR REGRESSION

RELATIONSHIP BETWEEN
ONE CONTINUOUS TARGET AND
TWO OR MORE PREDICTOR VARIABLES

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2$$

lm. intercept

lm. coef # array ([2,3])

MODEL EVALUATION USING VISUALIZATION

1. The relationship between 2 variables
2. The strength of the correlation
3. The direction of the correlation (+ or -)

REGRESSION PLOT

- Scatter plot

- Fitted linear regression line (\hat{y})

import seaborn as sns

sns.regplot(x = "indepVar", y = "depVar", data = df)

plt.ylim(0.)

- DISTRIBUTED EVENLY AROUND THE X AXIS WITH
SIMILAR VARIANCE SUGGESTS A LINEAR PLOT IS
APPROPRIATE

- IF THE RESIDUALS ARE NOT RANDOMLY SEPARATED, THIS SUGGESTS THE LINEAR ASSUMPTION IS INCORRECT.
- IF THE VARIANCE INCREASES WITH X, THE MODEL IS INCORRECT.

RESIDUAL PLOT

import seaborn as sns

```
sns.residplot(df["indVar"], df["depVar"])
    X           Y
```

DISTRIBUTION PLOT

VISUALIZING MODELS WITH MORE THAN 1 INDEPENDENT VARIABLE

import seaborn as sns

```
ax1 = sns.displot(df["depVar"], hist=False,
                   color="r", label="Actual Value")
```

```
sns.displot(yhat, hist=False, color="b",
             label="Fitted Values", ax=ax1)
```

POLYNOMIAL REGRESSION

MORE THAN 1 DIMENSION

- TRANSFORM DATA INTO POLYNOMIAL
- USE LINEAR REGRESSION TO FIT THE PARAMETER

PR is a SPECIAL CASE OF LR

USEFUL FOR DESCRIBING CURVILINEAR RELATIONSHIPS

2ND ORDER

$$\hat{y} = \underline{b_0} + \underline{b_1}x + \underline{b_2}x^2$$

↑
SQUARING OR SETTING
A HIGHER ORDER

3RD ORDER

$$\hat{y} = \underline{b_0} + \underline{b_1}x + \underline{b_2}x^2 + \underline{b_3}x^3$$

$$f = np.polyfit(x, y, 3)$$
$$p = np.poly1d(f)$$

print(p)

We can also have MULTIDIMENSIONAL
POLYNOMIAL LINEAR REGRESSION

NUMPY'S POLYFIT FUNCTION CANNOT PERFORM THIS TYPE OF REGRESSION

HAVE TO USE PREPROCESSING LIB IN SCI-KIT

from sklearn.preprocessing import PolynomialFeatures

pr = PolynomialFeatures(degree=2, include_bias=False)

X_poly = pr.fit_transform(X[["hp", "wm"]])

BECAUSE OF MULTIPLE DIMENSIONS IT IS CONVENIENT TO SCALE THE FEATURES:

from sklearn.preprocessing import StandardScaler

SCALE = StandardScaler()

SCALE.fit(X_data[["hp", "wm"]])

X_scale = SCALE.transform(X_data[["hp", "wm"]])

THERE ARE MANY STEPS TO GETTING A PREDICTION:

NORMALIZATION → POLYNOMIAL → LINEAR
TRANSFORM REGRESSION

USE A PIPELINE

PIPELINE

1. Import library Pipeline
2. Create a list of tuples
3. Import the list of step 2 in the pipeline constructor to create the pipeline object.
4. Apply the train method to the pipeline object.
5. Produce a prediction in the pipeline object.

```
Input = [("scale", StandardScaler()), ("polynomial",  
        PolynomialFeatures(degree=2)),  
        ("model", LinearRegression())]
```

```
pipe = Pipeline(Input)
```

MEASURES FOR IN-SAMPLE EVALUATION
A WAY TO NUMERICALLY DETERMINE HOW GOOD THE MODEL FITS ON DATASET

- COEFFICIENT OF DETERMINATION R^2
- MEAN SQUARE ERROR MSE

MSE

```
from sklearn.metrics import mean_squared_error  
mean_squared_error(df['price'], Y_predict_simple_fit)
```

↑ ↑
ACTUAL VALUE OF PREDICTED VALUE
TARGET VARIABLE OF TARGET VAR

R^2

THE % OF VARIATION OF THE TARGET VAR(Y) THAT IS EXPLAINED BY THE LINEAR MODEL.

$R^2 \approx 1$ GOOD FIT

$R^2 \approx 0$ BAD FIT

} CLOSER TO 1 INDICATES A BETTER FIT

$X = df[["highway-mpg"]]$

$y = df['price']$

$lm.fit(X, y)$

$lm.score(X, y)$

$R^2 < 1$ OVERRFITTING

EXAMPLE

$R^2 = 0.49$

MEANS APPROX 49% OF THE VARIATION OF PRICE IS EXPLAINED BY THIS SIMPLE LINEAR MODEL.

PREDICTION AND DECISION MAKING

How can we determine if our model is correct?

1. DO THE PREDICTED VALUES MAKE SENSE?
2. VISUALIZATION
3. NUMERICAL MEASURES FOR EVALUATION
4. COMPARING MODELS

Ex.

lm.fit(x, y)

lm.predict(value) # how the predicted value looks like? in between a reasonable range?

To generate a sequence of values for testing the model:

import numpy as np

new_input = np.arange(1, 101, 1).reshape(-1, 1)

↑ ↑ ↗
STARTING ENDING +1 STEP SIZE

yhat = lm.predict(new_input)

Using a REGRESSION PLOT SHOULD BE TRIED

RESIDUAL PLOT

DISTRIBUTION PLOT

SHOULD BE

- SOME AUTHORS SUGGEST $R^2 \geq 0.1$
- IS A LOWER MSE ALWAYS IMPLYING A BETTER FIT? NOT NECESSARILY

$$\text{MSE}_{\text{MLR}} < \text{MSE}_{\text{SLR}}$$

$$\text{MSE}_{\text{PR}} < \text{MSE}_{\text{LR}}$$

A SIMILAR INVERSE
RELATIONSHIP HOLDS
FOR R^2