

Week 4

## Model Development

1. Simple/Multiple Lin Reg
2. Model Eval Using Visualization
3. Polynomial Regression ; Pipelines
4.  $R^2$  + MSE for In-Sample Eval
5. Prediction + Decision Making

1. Linear regression refers to 1 I.V. to make a prediction.

- Simple lin reg (SLR):  $\hat{y} = b_0 + b_1 x$  (Data pts are usually stored into df or numpy arrays)

- predictor (IV)  $\rightarrow x$
- target (DV)  $\rightarrow y$
- In Python:

- 1 from sklearn.linear\_model import LinearRegression()
- 2 lm = LinearRegression() uses lm as a construct for lin reg function
- 3 x = df['var1']
- 4 y = df['var2']
- 5 lm.fit(x, y)

- Multiple lin reg (MLR):  $\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \dots$

- Predictors (IVs)  $\rightarrow x_1, x_2, \dots$
- Target (DV)  $\rightarrow y$
- In Python:

1, 2 " "

3 z = df[['var1', 'var2', ...]]

4 y = df['var\_y']

5 lm.fit(z, y)  $\leftarrow \hat{y}_{\text{hat}} = \text{lm.predict}(x)$ : prediction

2. Residual Plots

$$\text{Residual} = \hat{y} - y$$



x-value = x-value  
y-value = residual

If the points on a residual plot are spread out around the x-axis, lin reg is appropriate. If there is a slope or curve of residual points, then lin reg isn't accurate.

or var changes w/ x

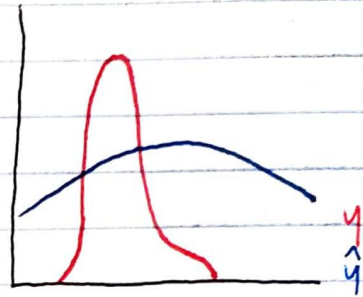
In Python:

- 1 import seaborn as sns
- 2 sns.residplot(df['var<sub>1</sub>'], df['var<sub>2</sub>'])

Distribution Plots count  $\hat{y}$  vs  $y$

In Python:

- 1 ax1 = sns.distplot(df['var<sub>y</sub>'], hist=False, color='r', label='Actual Value')
- 2 sns.distplot(yhat, hist=False, color='b', label='Fitted Values', ax=ax1)



### 3. Polynomial Regression

- Best for desc curvilinear rel's
- Quadratic:  $\hat{y} = b_0 + b_1 x_1 + b_2 (x_2)^2 \dots$

- In Python:

- 1 f = np.polyfit(x, y, deg)
- 2 p = np.polyd(f)
- 3 print(p)

Multidimensional Poly Reg

- 1 from sklearn.preprocessing import PolynomialFeatures
- 2 pr = PolynomialFeatures(degree=#)
- 3 x\_polly = pr.fit\_transform(x['var<sub>1</sub>','var<sub>2</sub>'], include\_bias=False)

Pre-processing: Normalize

- 1 from sklearn.preprocessing import StandardScaler
- 2 SCALE = StandardScaler()
- 3 SCALE.fit(x['var<sub>1</sub>','var<sub>2</sub>'])
- 4 x\_scale = SCALE.transform(x['var<sub>1</sub>','var<sub>2</sub>'])

Pipelines are a way to simplify a code.

Normalization → Polynomial Transform → Lin Reg (prediction)

- 1 from sklearn.preprocessing import PolynomialFeatures
- 2 from sklearn.linear\_model import LinearRegression
- 3 from sklearn.preprocessing import StandardScaler
- 4 from sklearn.pipeline import Pipeline
- 5 Input = [('scale', StandardScaler()), ('polynomial', PolynomialFeatures(degree = #)), ('model', LinearRegression())]
- 6 pipe = Pipeline(Input)
- 7 Pipe.train(x['var\_1', 'var\_2', ...], y)
- 8 yhat = Pipe.predict(x['var\_1', 'var\_2'])

#### 4. In-Sample Eval

- A way to determine how good the model fits on dataset
- Mean squared error (MSE)

- Square the residual for each value

- $MSE = \frac{\sum (\text{residual}^2)}{\# \text{ samples}}$

- In Python:

- 1 from sklearn.metrics import mean\_squared\_error
- 2 mean\_squared\_error(df['var\_1'], y - predict\_simple\_fit)

- $R^2$ : coeff of determination

- Measures how close the data is to reg line (or % of var on target var (y) exp. by model)

- $R^2 = (1 - \frac{MSE \text{ of reg line}}{MSE \text{ of avg of data}})$

- See last line of 4.1

#### 5. Prediction + Decision Making

1. Do results make sense?
2. Always visualize
3. Numerical measures for eval
4. Comparing models

MSE for MLR models will be smaller than MSE for SLR models because  $\uparrow \text{vars} \rightarrow \uparrow \text{accuracy}$