

Data Science I: Homework 1 Python Code

Neal Kuperman

January 15, 2026

Contents

1.	Helper Functions	2
2.	ISLP 3.10	6
3.	ISLP 3.14	11
4.	ISLP 3.15	18
5.	ESL 3.17	23

Helper Functions

1.1 Imports

```
1 import ISLP
2 import matplotlib.pyplot as plt
3 import pandas as pd
4 import numpy as np
5 import statsmodels.api as smf
6 from IPython.display import display
7 from lin_reg_plots import LinearRegDiagnostic
8 import pickle
9 from ucimlrepo import fetch_ucirepo
10 from sklearn.linear_model import LinearRegression, Ridge, Lasso
11 from sklearn.model_selection import train_test_split
12 from sklearn.preprocessing import StandardScaler
13 from sklearn.metrics import mean_squared_error, r2_score
14 from sklearn.model_selection import GridSearchCV
15 from sklearn.feature_selection import SequentialFeatureSelector
```

1.2 Summarize Function

```
1 def summarize(model, vars = [], verbose = True):
2     if isinstance(model.params, np.ndarray):
3         params = model.params
4         if len(vars) != (len(params) - 1):
5             if verbose:
6                 print("Warning: The number of variables does not match
7                     the number of parameters")
8         vars = ["intercept"] + [f"xi{i}" for i in range(len(params)
9             - 1)]
10    else:
11        params = model.params.values
12        vars = model.params.index.to_list()
13
14    tvalues = np.round(model.tvalues, 4)
15    pvalues = np.round(model.pvalues, 4)
16    std_err = np.round(model.bse, 4)
17
18    param_summary = pd.DataFrame(index=vars)
19    param_summary["coef"] = params
20    param_summary["t value"] = tvalues
21    param_summary["p value"] = pvalues
22    param_summary["std err "] = std_err
23
24    r_squared = np.round(model.rsquared, 4)
25    F_stat = np.round(model.fvalue, 4)
26    # model_summary = pd.DataFrame(columns=["R-squared", "F-statistic
27    #"])
28    model_summary = pd.DataFrame({"R-squared": [r_squared], "F-
29        statistic": [F_stat]}, index = ["value"])
30
31    return param_summary, model_summary
```

1.3 Leverage Functions

```
1 def calc_leverage(x):
2     """
3         Observations with high leverage have an unusual value for xi.
4         High leverage observations tend to have a sizable impact on the
5         estimated regression line.
6
7     SLR:  hi = 1/n + (xi - x_bar)^2 / sum((xi - x_bar)^2)
8     MLR:  Hii = diag(H), where H = X^T (X^T X)^{-1} X
9     """
10    H = x @ np.linalg.inv(x.T @ x) @ x.T
11    leverage_points = np.diag(H)
12    return leverage_points
13
14 def find_high_leverage_points(x):
15     """
16         Observations with high leverage have an unusual value for xi.
17         hii > 2p/n
18             p = number of predictors
19             n = number of observations
20     """
21     leverage_points = calc_leverage(x)
22     high_leverage_points = leverage_points > 2 * x.shape[1]/x.shape[0]
23     index = np.where(high_leverage_points)[0]
24     return index, leverage_points[index]
```

1.4 Cook's Distance

```
1  def calc_cooks_distance(x, y):
2      """
3          Cook's distance is a measure of the influence of an observation
4          on the estimated regression coefficients.
5
6          D_i = (e_i^2 / (p * MSE)) * (h_ii / (1 - h_ii)^2)
7
8          where:
9              e_i = residual (y_i - y_hat_i)
10             h_ii = leverage (diagonal of hat matrix)
11             p = number of predictors (including intercept)
12             MSE = mean squared error = SS_res / (n - p)
13
14      n, p = x.shape
15
16      # Fit model
17      beta = np.linalg.inv(x.T @ x) @ x.T @ y
18      y_hat = x @ beta
19      residuals = y - y_hat
20
21      # MSE
22      SS_res = np.sum(residuals**2)
23      MSE = SS_res / (n - p)
24
25      # Leverage
26      h = calc_leverage(x)
27
28      # Cook's distance
29      cooks_d = (residuals**2 / (p * MSE)) * (h / (1 - h)**2)
30
31      return cooks_d
```

2.1 Load Data

```
1 carseats = ISLP.load_data('Carseats')
2
3 quant_cols = carseats.select_dtypes(include=['number']).columns
4 cat_cols = carseats.select_dtypes(include=['category']).columns
5
6 # convert categorical columns to dummy variables. Creates a new column
   # for each category level for each categorical column.
7 carseats = pd.get_dummies(carseats, columns=cat_cols, dtype=float)
```

2.2 (a) Multiple Regression: Sales ~ Price + Urban + US

```
1 # Must add an intercept column when using statsmodels.api.OLS
2 X = carseats[["Price", "Urban_Yes", "US_Yes"]]
3 X.insert(0, 'Intercept', 1.0)
4 y = carseats["Sales"]
5
6 model = smf.OLS(y, X).fit()
7 print(model.summary())
```

2.3 (e) OLS with Price and US_Yes Only

```
1 X = carseats[["Price", "US_Yes"]]
2 X.insert(0, 'Intercept', 1.0)
3 y = carseats["Sales"]
4
5 model = smf.OLS(y, X).fit()
6 print(model.summary())
7 # print(model.summary().as_latex())
```

2.4 (g) Confidence Intervals

```
1 model.conf_int(alpha=0.05)
```

2.5 (h) High Leverage Points and Cook's Distance

```
1 index, high_leverage_points = find_high_leverage_points(X.values)
2 lvg_df = pd.DataFrame({"index": index, "high_leverage_points":
3     high_leverage_points})
# print(lvg_df.to_latex(index=False))
4
5 # Get influence measures
6 influence = model.get_influence()
7 cooks_d = influence.cooks_distance[0] # [0] is the values, [1] is p-
8     values
9
10 # Plot Cook's D
11 plt.stem(range(len(cooks_d)), cooks_d, markerfmt=",")
12 plt.xlabel('Observation')
13 plt.ylabel("Cook's Distance")
14 plt.savefig("../images/3_10_h_cooks_d.png")
15 plt.show()
```

3.1 (a) Generate Data

```
1 np.random.seed(5)
2 rng = np.random.default_rng(10)
3 x1 = rng.uniform(0, 1, size=100)
4 x2 = 0.5 * x1 + rng.normal(size=100) / 10
5 y = 2 + 2 * x1 + 0.3 * x2 + rng.normal(size=100)
```

3.2 (b) Correlation Between x1 and x2

```
1 # Manual calculation of correlation
2 x1_mean = np.mean(x1)
3 x2_mean = np.mean(x2)
4
5 numerator = np.sum((x1 - x1_mean) * (x2 - x2_mean))
6 denominator = np.sqrt(np.sum((x1 - x1_mean)**2) * np.sum((x2 - x2_mean)
    **2))
7
8 correlation = numerator / denominator
9 print(f"Correlation between x1 and x2: {correlation:.3f}")
10
11 plt.scatter(x1, x2)
12 plt.xlabel('x1')
13 plt.ylabel('x2')
14 plt.title('Scatterplot of x1 and x2')
15 plt.show()
```

3.3 (c) OLS with x1 and x2

```
1  intercept = np.ones(len(x1))
2  X = np.column_stack((intercept, x1, x2))
3  model = smf.OLS(y, X).fit()
4  print(model.summary())
5  # print(model.summary().as_latex())
```

3.4 (d) OLS with x1 Only

```
1  intercept = np.ones(len(x1))
2  X = np.column_stack((intercept, x1))
3  model = smf.OLS(y, X).fit()
4  print(model.summary())
5  # print(model.summary().as_latex())
```

3.5 (e) OLS with x2 Only

```
1  intercept = np.ones(len(x1))
2  X = np.column_stack((intercept, x2))
3  model = smf.OLS(y, X).fit()
4  print(model.summary())
5  # print(model.summary().as_latex())
```

3.6 (g) New Observation Analysis

```
1  x1_new = np.concatenate([x1, [0.1]])
2  x2_new = np.concatenate([x2, [0.8]])
3  y_new = np.concatenate([y, [6]])
4
5  intercept = np.ones(len(x1_new))
6  X_full = np.column_stack((intercept, x1_new, x2_new))
7  X_x1 = np.column_stack((intercept, x1_new))
8  X_x2 = np.column_stack((intercept, x2_new))
9
10 model_full = smf.OLS(y_new, X_full).fit()
11 model_x1 = smf.OLS(y_new, X_x1).fit()
12 model_x2 = smf.OLS(y_new, X_x2).fit()
13
14 print("Full Model")
15 print("*"*60)
16 param_summary_full, model_summary_full = summarize(model_full, verbose=False)
17 display(param_summary_full)
18 display(model_summary_full)
19
20 print("\nx1 Model")
21 print("*"*60)
22 param_summary_x1, model_summary_x1 = summarize(model_x1, verbose=False)
23 display(param_summary_x1)
24 display(model_summary_x1)
25
26 print("\nx2 Model")
27 print("*"*60)
28 param_summary_x2, model_summary_x2 = summarize(model_x2, verbose=False)
29 display(param_summary_x2)
30 display(model_summary_x2)
31
32 # Check for high leverage and influential points
33 index, high_leverage_points = find_high_leverage_points(X_full)
34 print(index)
35 print(high_leverage_points)
36
37 influential_stats = model_full.get_influence()
38 cooks_D = influential_stats.cooks_distance[0]
39 influential_pts_index = np.where(cooks_D > 1)[0]
40 influential_pts_vals = cooks_D[influential_pts_index]
41
42 print(influential_pts_index)
43 print(influential_pts_vals)
44
45 cls = LinearRegDiagnostic(model_full)
46 vif, fig, ax = cls()
47 print(vif)
```


4.1 (a) Simple Linear Regression for Each Predictor

```

1   bos = ISLP.load_data('Boston')
2
3   # Get predictor columns (excluding 'crim')
4   predictors = [col for col in bos.columns if col != 'crim']
5   n_predictors = len(predictors)
6
7   # Create subplot grid
8   n_cols = 3
9   n_rows = (n_predictors + n_cols - 1) // n_cols
10  fig_1, axes_1 = plt.subplots(n_rows, n_cols, figsize=(16, 3*n_rows))
11  axes_1 = axes_1.flatten()
12
13  fig_2, axes_2 = plt.subplots(n_rows, n_cols, figsize=(16, 3*n_rows))
14  axes_2 = axes_2.flatten()
15
16  models = {}
17  r_squared = []
18
19  for i, col_name in enumerate(predictors):
20      if col_name == 'crim':
21          continue
22      X = bos[col_name].values
23      intercept = np.ones(len(X))
24      X = np.column_stack((intercept, X))
25      model = smf.OLS(bos['crim'], X).fit()
26      models[col_name] = model
27      predictions = model.predict(X)
28      residual = bos['crim'] - predictions
29      r_squared.append(model.rsquared)
30
31      axes_1[i].scatter(bos[col_name], bos['crim'], alpha=0.5)
32      axes_1[i].plot(bos[col_name].sort_values(), predictions[bos[
33          col_name].argsort()], color='red')
34      axes_1[i].set_xlabel(col_name, fontsize=14)
35      axes_1[i].set_ylabel('crim', fontsize=14)
36      axes_1[i].set_title(f'crim vs {col_name}', fontsize=14)
37
38      axes_2[i].scatter(bos[col_name], residual, alpha=0.5)
            axes_2[i].set_xlabel(col_name, fontsize=14)

```

```
39     axes_2[i].set_ylabel('residual', fontsize=14)
40     axes_2[i].set_title(f'residual - {col_name}', fontsize=14)
41
42     fig_1.tight_layout(pad=1.5, h_pad=2, w_pad=1)
43     fig_1.savefig('../images/3_15_a_scatter.png', dpi=300, bbox_inches='tight')
44
45     fig_2.tight_layout(pad=1.5, h_pad=2, w_pad=1)
46     fig_2.savefig('../images/3_15_a_residuals.png', dpi=300, bbox_inches='tight')
47
48     plt.show()
49
50     r_squared_df = pd.DataFrame({'predictor': predictors, 'r_squared':
51         r_squared})
52     r_squared_df = r_squared_df.sort_values(by='r_squared', ascending=False
53     )
54     print(r_squared_df)
```

4.2 (b) Multiple Regression with All Predictors

```
1 X = bos[predictors]
2 X.insert(0, 'Intercept', 1.0)
3 y = bos['crim']
4 model = smf.OLS(y, X).fit()
5 print(model.summary())
6
7 param_summary, model_summary = summarize(model, verbose=False)
8 display(param_summary)
9 display(model_summary)
```

4.3 (d) Non-linear Associations (Polynomial Regression)

```
1 n_cols = 3
2 n_rows = (n_predictors + n_cols - 1) // n_cols
3 fig_1, axes_1 = plt.subplots(n_rows, n_cols, figsize=(16, 3*n_rows))
4 axes_1 = axes_1.flatten()
5
6 fig_2, axes_2 = plt.subplots(n_rows, n_cols, figsize=(16, 3*n_rows))
7 axes_2 = axes_2.flatten()
8
9 models = {}
10 r_squared = []
11
12 for i, col_name in enumerate(predictors):
13     if col_name == 'crim':
14         continue
15     X = bos[col_name].values
16     X2 = X**2
17     X3 = X**3
18     intercept = np.ones(len(X))
19     X = np.column_stack((intercept, X, X2, X3))
20     _df = pd.DataFrame(X)
21     _df.columns = ['intercept', col_name, f'{col_name}^2', f'{col_name}^3']
22     model = smf.OLS(bos['crim'], _df).fit()
23     models[col_name] = model
24     predictions = model.predict(X)
25     residual = bos['crim'] - predictions
26     r_squared.append(model.rsquared)
27
28     print(f"Model for {col_name}:")
29     print('='*60)
30     param_summary, model_summary = summarize(model, verbose=False)
31     display(param_summary)
32     display(model_summary)
33     print("\n")
34
35     axes_1[i].scatter(bos[col_name], bos['crim'], alpha=0.5)
36     axes_1[i].plot(bos[col_name].sort_values(), predictions[bos[
37         col_name].argsort()], color='red')
38     axes_1[i].set_xlabel(col_name, fontsize=14)
39     axes_1[i].set_ylabel('crim', fontsize=14)
40     axes_1[i].set_title(f'crim vs {col_name}', fontsize=14)
41
42     axes_2[i].scatter(bos[col_name], residual, alpha=0.5)
43     axes_2[i].set_xlabel(col_name, fontsize=14)
44     axes_2[i].set_ylabel('residual', fontsize=14)
45     axes_2[i].set_title(f'residual - {col_name}', fontsize=14)
46
fig_1.tight_layout(pad=1.5, h_pad=2, w_pad=1)
```

```
47     fig_1.savefig('..../images/3_15_d_scatter.png', dpi=300, bbox_inches='
48         tight')
49
50     fig_2.tight_layout(pad=1.5, h_pad=2, w_pad=1)
51     fig_2.savefig('..../images/3_15_d_residuals.png', dpi=300, bbox_inches='
52         tight')
53
54     plt.tight_layout()
55     plt.show()
56
57     r_squared_df = pd.DataFrame({'predictor': predictors, 'r_squared':
58         r_squared})
59     r_squared_df = r_squared_df.sort_values(by='r_squared', ascending=False
60         )
61     print(r_squared_df)
62     # r_squared_df.to_latex(index = False)
```

5.1 Load and Prepare Data

```
1  spambase = fetch_uci_repo(id=94)
2  scaler = StandardScaler()
3
4  X = spambase.data.features
5  y = spambase.data.targets
6
7  X_train, X_test, y_train, y_test = train_test_split(
8      X, y, test_size=0.3, random_state=42
9  )
10
11 X_train_scaled = scaler.fit_transform(X_train)
12 X_test_scaled = scaler.transform(X_test)
```

5.2 OLS, Ridge, and Lasso Regression

```
1  models = {
2      'LS': {
3          'model': LinearRegression(),
4          'param_grid': {}
5      },
6      'ridge': {
7          'model': Ridge(),
8          'param_grid': {
9              'alpha': np.concatenate([np.arange(0.005, 10, 0.05), np.
10                             arange(10, 2000, 10)]),
11              "max_iter": [1000]
12          }
13      },
14      'lasso': {
15          'model': Lasso(),
16          'param_grid': {
17              'alpha': np.concatenate([np.arange(0.01, .1, 0.001), np.
18                             arange(1, 2000, 10)]),
19              'max_iter': [1000]
20          }
21      }
22
23      results = {
24          "name": [], "model": [], "intercept": [],
25          "coefficients": [], "test_error": [], "test_R2": []
26      }
27
28      for name, model in models.items():
29          print()
30          print(name)
31          print("*"*60)
32          opt_param = {}
33
34          if model["param_grid"]:
35              grid_search = GridSearchCV(
36                  model["model"],
37                  param_grid=model["param_grid"],
38                  cv=5,
39                  scoring='neg_mean_squared_error'
40              )
41              grid_search.fit(X_train_scaled, y_train)
42              opt_param = {
43                  "alpha": grid_search.best_params_['alpha'],
44                  "max_iter": grid_search.best_params_['max_iter']
45              }
46
47          model = model["model"].set_params(**opt_param)
```

```
47     model.fit(X_train_scaled, y_train)
48
49     test_pred = model.predict(X_test_scaled)
50     test_error = mean_squared_error(y_test, test_pred)
51     test_R2 = r2_score(y_test, test_pred)
52
53     print(test_R2)
54
55     results["name"].append(name)
56     results["model"].append(model)
57     results["intercept"].append(model.intercept_[0])
58     results["coefficients"].append(model.coef_.flatten())
59     results["test_error"].append(test_error)
60     results["test_R2"].append(test_R2)
61
62     # Create coefficients DataFrame
63     coefs = np.column_stack([arr for arr in results["coefficients"]])
64     coefs = np.vstack([results["intercept"], coefs])
65
66     cols = X.columns.to_list()
67     cols = ['intercept'] + cols
68     model_names = results["name"]
69
70     df = pd.DataFrame(coefs, index=cols, columns=model_names)
```

5.3 Best Subset Selection

```
1 def best_subset_selection(X, y, X_train, X_test, y_train, y_test, rerun
2     =False):
3     best_subset_models = {}
4     results = {}
5     linear_models = {
6         "model": [],
7         "coefs": [],
8         "intercept": [],
9         "score": []
10    }
11
12    if rerun:
13        for i in range(1, len(X.columns)):
14            print(i)
15            sfs = SequentialFeatureSelector(
16                LinearRegression(),
17                n_features_to_select=i,
18                direction='forward',
19            )
20            sfs.fit(X_train, y_train)
21
22            # Select features from TRAIN data, fit on TRAIN
23            selected_train = X_train[:, sfs.get_support()]
24            selected_test = X_test[:, sfs.get_support()]
25
26            selected_cols = X.columns[sfs.get_support()]
27            df_train = pd.DataFrame(selected_train, columns=
28                selected_cols)
29            df_test = pd.DataFrame(selected_test, columns=selected_cols
30                )
31
32            lin_mod = LinearRegression()
33            lin_mod.fit(df_train, y_train)
34
35            best_subset_models[i] = sfs
36            results[i] = sfs.get_support()
37            linear_models["model"].append(lin_mod)
38            linear_models["coefs"].append(lin_mod.coef_)
39            linear_models["intercept"].append(lin_mod.intercept_)
40            linear_models["score"].append(lin_mod.score(df_test, y_test
41                ))
42
43            results = pd.DataFrame(results)
44            results.index.name = 'Feature'
45
46            save_dict = {
47                "linear_models": linear_models,
48                "results": results,
```

```

45         "best_subset_models": best_subset_models
46     }
47
48     with open("HW_1_ESL_3_17.pkl", "wb") as f:
49         pickle.dump(save_dict, f)
50
51 else:
52     with open('HW_1_ESL_3_17.pkl', 'rb') as f:
53         save_dict = pickle.load(f)
54     results = save_dict["results"]
55     linear_models = save_dict["linear_models"]
56     best_subset_models = save_dict["best_subset_models"]
57
58     plt.scatter(range(1, len(linear_models["score"])+1), linear_models[
59         "score"])
60     plt.savefig("../images/ESL_3_17_best_subset_score.png")
61     plt.show()
62     return results, linear_models, best_subset_models
63
64 # Run best subset selection
65 best_subset_results, best_subset_linear_models, best_subset_models = \
66     best_subset_selection(X, y, X_train, X_test, y_train, y_test, rerun
67     =False)
68
69 best_subset_idx = best_subset_linear_models["score"].index(
70     max(best_subset_linear_models["score"]))
71 best_subset_model = best_subset_models[best_subset_idx]
72
73 col_mask = best_subset_model.get_support()
74 cols = X.columns[col_mask]
75
76 # Select features from TRAIN data, fit on TRAIN
77 selected_train = X_train_scaled[:, col_mask]
78 selected_test = X_test_scaled[:, col_mask]
79
80 # Add intercept
81 selected_train_with_intercept = smf.add_constant(selected_train)
82 selected_test_with_intercept = smf.add_constant(selected_test)
83
84 # Fit OLS model
85 model = smf.OLS(y_train, selected_train_with_intercept).fit()
86
87 # Create best subset df
88 all_indices = df.index.to_list()
89 best_subset_df = pd.DataFrame({"best_subset": model.params[1:].values},
90     index=cols)
91
92 best_subset_df_all_indices = pd.DataFrame(
93     {"best_subset": [0.0]*len(all_indices)},

```

```

92     index=all_indices
93 )
94
95 shared_idx = best_subset_df.index.intersection(
96     best_subset_df_all_indices.index)
97 best_subset_df_all_indices.loc[shared_idx, "best_subset"] = \
98     best_subset_df.loc[shared_idx, "best_subset"]
99 best_subset_df_all_indices.loc["intercept", "best_subset"] = model.
100    params.iloc[0]
101
102 # Add best subset model to df of coefficients
103 df["best_subset"] = best_subset_df_all_indices["best_subset"]
104
105 # Add best subset model info to results dictionary
106 test_pred = model.predict(selected_test_with_intercept)
107 test_error = mean_squared_error(y_test, test_pred)
108 test_R2 = r2_score(y_test, test_pred)
109
110 results["name"].append("best subset")
111 results["model"].append(model)
112 results["intercept"].append(model.params.iloc[0])
113 results["coefficients"].append(model.params[1:].values)
114 results["test_error"].append(test_error)
115 results["test_R2"].append(test_R2)

```

5.4 Coefficient Comparison Plot

```
1 # Split features into 10 groups for readability
2 n_features = len(df)
3 chunk_size = (n_features + 9) // 10
4
5 fig, axes = plt.subplots(5, 2, figsize=(14, 20))
6 axes = axes.flatten()
7
8 for i, ax in enumerate(axes):
9     start = i * chunk_size
10    end = min((i + 1) * chunk_size, n_features)
11
12    df.iloc[start:end].plot(kind='bar', ax=ax, width=0.8, legend=False)
13
14    ax.set_ylabel('Coefficient Value')
15    ax.axhline(y=0, color='black', linestyle='--', linewidth=0.5)
16    ax.tick_params(axis='x', rotation=45)
17
18    for label in ax.get_xticklabels():
19        label.set_ha('right')
20
21 axes[0].set_title('Coefficients by Feature')
22
23 handles, labels = axes[0].get_legend_handles_labels()
24 fig.legend(handles, labels, loc='lower center', ncol=len(labels),
25            bbox_to_anchor=(0.5, -0.02))
26
27 plt.tight_layout()
28 fig.subplots_adjust(bottom=0.08)
29
30 plt.savefig("../images/ESL_3_17_coefficients.png", bbox_inches='tight')
31 plt.show()
32
33 # Print results
34 print(df)
35 print(results["test_R2"])
36 print(results["test_error"])
```