

# In-Class Exercises: Supervised Learning

Work through these during the lecture. Ask if you get stuck.

Data Science with Python – BSc Course

---

## Exercise 1: OLS by Hand

8 min | Pen & Paper

Five training points:

Point	$x$	$y$
1	1	2.0
2	2	3.5
3	3	5.0
4	4	6.2
5	5	7.8

Fit an ordinary least squares (OLS) line  $\hat{y} = \hat{\beta}_0 + \hat{\beta}_1 x$  using the closed-form formulas:

$$\hat{\beta}_1 = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad \hat{\beta}_0 = \bar{y} - \hat{\beta}_1 \bar{x}$$

(a) Compute  $\bar{x}$  and  $\bar{y}$ .

(b) Fill in the working table:

$x_i$	$y_i$	$x_i - \bar{x}$	$y_i - \bar{y}$	$(x_i - \bar{x})(y_i - \bar{y})$
1	2.0			
2	3.5			
3	5.0			
4	6.2			
5	7.8			

(c) Compute the numerator  $\sum (x_i - \bar{x})(y_i - \bar{y})$  and the denominator  $\sum (x_i - \bar{x})^2$ , then obtain  $\hat{\beta}_1$  and  $\hat{\beta}_0$ .

(d) Predict  $\hat{y}$  at  $x = 6$ . What is the residual if the true  $y_6 = 9.2$ ?

## Exercise 2: Read the Coefficient Path

5 min | Conceptual

Refer to the ridge regression coefficient path shown in **lecture slide 22** (Coefficient Shrinkage: Ridge vs. Lasso). The plot shows four coefficients  $\beta_1, \beta_2, \beta_3, \beta_4$  on the vertical axis and the regularization strength  $\lambda$  on the horizontal axis (increasing to the right).

(a) As  $\lambda$  increases, in which direction do the coefficients move — toward zero or away from zero? Why?

(b) Which feature has the *largest absolute coefficient* at small  $\lambda$ ? In practical terms, what does this tell us about the feature?

(c) At what  $\lambda$  do all four coefficients become essentially indistinguishable from zero? What is the resulting model at that point?

(d) In a **Lasso** (L1) path plot, some coefficients would become *exactly* zero at finite  $\lambda$ , while in a Ridge (L2) path they only approach zero asymptotically. Why does this difference matter in practice?

### Exercise 3: Compute Precision, Recall, F1

7 min | Pen & Paper

A fraud-detection classifier produced the following confusion matrix on a test set of 1,000 transactions. The positive class is “fraudulent.”

	Predicted Fraud	Predicted Legit
Actual Fraud	TP = 80	FN = 15
Actual Legit	FP = 20	TN = 885

Recall the definitions:

$$\text{Precision} = \frac{TP}{TP + FP}, \quad \text{Recall} = \frac{TP}{TP + FN},$$
$$\text{F1} = 2 \cdot \frac{\text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}, \quad \text{Accuracy} = \frac{TP + TN}{TP + FP + FN + TN}.$$

(a) Compute precision, recall, F1, and accuracy. Show your working.

(b) A baseline that labels *every* transaction as legitimate would achieve what accuracy? What would its precision and recall on the fraud class be?

(c) In fraud detection, investigators pay a cost for each false positive (wasted investigation time), and the bank loses money on each false negative (fraud goes through). If false negatives are  $10\times$  more costly than false positives, would you prefer a model with high precision or high recall?

## Exercise 4: Which Model?

5 min | Discussion

Match each scenario to the most appropriate supervised learning method. Write your reasoning — more than one method can be defensible.

#	Scenario	Best Method?
A	Predict a continuous stock return (real-valued target) given a small set of clean, roughly linear factors.	
B	Classify whether a loan applicant will default. The bank's risk committee requires that each decision be <i>explainable</i> to regulators.	
C	Classify credit-card transactions as fraud vs. legitimate. The training data has only 1% fraud, and the bank wants top performance.	
D	Predict quarterly revenue using 50 candidate features, many of which are strongly correlated with each other.	

**Methods to choose from:** Linear Regression, Logistic Regression, Decision Tree, Random Forest, Lasso, SMOTE + XGBoost.

**Bonus:** For scenario C, why is plain logistic regression a poor choice? What does SMOTE fix, and why add XGBoost on top?

## Exercise 5: Build a Classifier Pipeline

10 min | Python (Optional)

Complete the blanks in this `sklearn` pipeline that predicts loan default from three borrower features. The pipeline should: (1) standardize features, (2) generate polynomial interactions, (3) fit a logistic regression.

```
1 from sklearn.pipeline import Pipeline
2 from sklearn.preprocessing import _____, _____ # (1)(2)
3 from sklearn.linear_model import _____ # (3)
4 from sklearn.model_selection import train_test_split
5 import numpy as np
6
7 # Synthetic loan data: 500 applicants, 3 features (income,
8 # debt_ratio, credit_score); target y is 0 = no default, 1 = default.
9 np.random.seed(42)
10 X = np.random.randn(500, 3)
11 y = (X[:, 0] + X[:, 1] - X[:, 2] + 0.3 * np.random.randn(500) > 0).
12     astype(int)
13
14 X_train, X_test, y_train, y_test = train_test_split(
15     X, y, test_size=0.2, random_state=42
16 )
17 # Build the pipeline
18 pipe = Pipeline([
19     ('scaler', _____(_____)), # (4) standard scaler
20     ('poly', _____(_____)), # (5) degree-2 polynomial
21     ('clf', _____(_____)), # (6) logistic reg., C=1.0
22 ])
23
24 # Fit on training data
25 pipe._____(X_train, y_train) # (7) which method?
26
27 # Evaluate
28 train_acc = pipe.score(X_train, y_train)
29 test_acc = pipe.score(X_test, y_test)
30 print(f"Train accuracy: {train_acc:.3f}")
31 print(f"Test accuracy : {test_acc:.3f}")
```

(a) Fill in blanks (1)–(7).

(b) Why is the order (scaler → polynomial features → classifier) important? What would happen if you applied polynomial features *before* standardization?

(c) The regularization strength in `LogisticRegression` is parameterized as `C` (inverse of  $\lambda$ ). Does `C = 0.01` or `C = 100` correspond to stronger regularization?

## Exercise 6: Random Forest Feature Importance 10 min | Python (Bonus)

Use a Random Forest to identify which features drive loan defaults, then visualize the importances.

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.ensemble import RandomForestClassifier
5
6 np.random.seed(42)
7
8 # Synthetic loan data: 1000 applicants, 6 features
9 n = 1000
10 loans = pd.DataFrame({
11     'income': np.random.normal(60_000, 20_000, n),
12     'debt_ratio': np.random.uniform(0.0, 1.0, n),
13     'credit_score': np.random.normal(700, 80, n),
14     'age': np.random.normal(40, 12, n),
15     'employ_yrs': np.random.exponential(5, n),
16     'dependents': np.random.poisson(1.5, n),
17 })
18
19 # Target: default depends mostly on debt_ratio and credit_score
20 logits = (2.5 * loans['debt_ratio']
21           - 0.01 * (loans['credit_score'] - 700)
22           + 0.3 * np.random.randn(n))
23 loans['default'] = (logits > 0.5).astype(int)
24
25 X = loans.drop(columns='default')
26 y = loans['default']
27
28 # Step 1: Fit a Random Forest
29 rf = RandomForestClassifier(n_estimators=200, random_state=42)
30 rf.fit(X, y)
31
32 # Step 2: Extract feature importances
33 importances = pd.Series(rf.feature_importances_, index=X.columns)
34 importances = importances.sort_values(ascending=True)
35 print(importances.round(3))
36
37 # Step 3: Bar chart
38 importances.plot(kind='barh', color='#1e3a5f')
39 plt.xlabel('Feature importance')
40 plt.title('Random Forest feature importance')
41 plt.tight_layout()
42 plt.savefig('feature_importance.pdf', dpi=300)
```

(a) Run the code. Which two features does the Random Forest rank highest? Does this match the way the target was generated?

(b) Why do features like `dependents` or `age` end up with low importance even though the

Random Forest split on them many times during training?

(c) Feature importance in Random Forests is known to *inflate* the importance of high-cardinality numeric features (e.g., income, credit score) relative to low-cardinality ones (e.g., dependents). What alternative method (introduced in L26 or L32) would give a fairer ranking?