# Sri Sivasubramaniya Nadar College of Engineering, Chennai

(An Autonomous Institution Affiliated to Anna University)

## Department of Computer Science and Engineering

### Theory Assignment 1

**Name:** Nithish Ra
**Register No:** 3122237001033

**Degree & Branch:** M. Tech (Integrated) Computer Science & Engineering
**Semester:** V
**Subject Code & Name:** ICS1502 – Introduction to Machine Learning
**Academic Year:** 2025–2026 (Odd) Batch: 2023–2028

# 1. Regression: Mobile Phone Price Prediction

In this task, we implement and evaluate **Linear Regression** models for predicting mobile phone prices. We use the dataset `Mobile-Price-Prediction-cleaned_data.csv`, where the target variable is "Price".

## 1.1 Data Representation

The dataset was split into training and testing sets using an 80:20 split. Let:

$$X \in R^{m \times d}, \quad y \in R^{m \times 1}$$

where $m$ is the number of samples and $d$ is the number of features. A bias term was added to $X$.

```python
import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.preprocessing import StandardScaler
from sklearn.metrics import mean_squared_error, r2_score

data = pd.read_csv("/Mobile-Price-Prediction-cleaned_data.csv")
print("Dataset shape:", data.shape)

# Separate features (X) and target (y)
X = data.drop("Price", axis=1).values    # shape (m, d)
y = data["Price"].values.reshape(-1, 1)  # shape (m, 1)

m, d = X.shape
print(f"m={m}, d={d}")

# Split the data
X_train, X_test, y_train, y_test = train_test_split(X, y, test_size
    =0.2, random_state=42)
m_train = X_train.shape[0]
m_test  = X_test.shape[0]
print("Train shape:", X_train.shape, " Test shape:", X_test.shape)

# Function to add bias (x0=1) column
def add_bias(X):
    return np.concatenate([np.ones((X.shape[0], 1)), X], axis=1)
```

## 1.2 Closed-form Solution

Using the Moore–Penrose pseudoinverse, we calculate $\theta$ directly using the Normal Equation, implemented via `np.linalg.pinv`.

$$\theta = (X^T X)^{-1} X^T y$$

This gives the optimal parameters directly on the unscaled data.

```python
# Add bias to the unscaled training and test data
X_train_bias = add_bias(X_train)
X_test_bias  = add_bias(X_test)
```

```
 4
 5  # Calculate theta using pseudoinverse
 6  theta_closed = np.linalg.pinv(X_train_bias) @ y_train    # shape (d+1,
        1)
 7  y_pred_closed = X_test_bias @ theta_closed
 8
 9  # Evaluate performance
10  mse_closed = mean_squared_error(y_test, y_pred_closed)
11  r2_closed = r2_score(y_test, y_pred_closed)
12  print("\nClosed-form (pinv) -- unscaled: MSE = {:.4f}, R2 = {:.4f}".
        format(mse_closed, r2_closed))
```

## 1.3 Gradient Descent

We implemented batch gradient descent with L2 regularization support (controlled by $\lambda$). The non-regularized update rule is:

$$\theta := \theta - \alpha \left( \frac{1}{m} X^T (X\theta - y) \right)$$

where $\alpha$ is the learning rate. Convergence was monitored using loss plots. For this step, we standardized the data using StandardScaler before applying GD, as it converges much more reliably and quickly.

```
 1  # Standardize features for GD
 2  scaler = StandardScaler()
 3  X_train_scaled = scaler.fit_transform(X_train)
 4  X_test_scaled  = scaler.transform(X_test)
 5
 6  # Add bias to standardized data
 7  X_train_scaled_bias = add_bias(X_train_scaled)
 8  X_test_scaled_bias  = add_bias(X_test_scaled)
 9
10  # Note: The provided code re-uses the variable names X_train_bias/
        X_test_bias.
11  # For clarity, we will rename them here to reflect they are scaled.
12  X_train_bias_scaled_gd = X_train_scaled_bias
13  X_test_bias_scaled_gd  = X_test_scaled_bias
14
15  def gradient_descent(X, y, lr=0.01, epochs=2000, lam=0.0, tol=1e-8,
        verbose=False):
16      m, n = X.shape   # m samples, n = d+1 params
17      theta = np.zeros((n,1))
18      prev_loss = np.inf
19      losses = []
20      for it in range(epochs):
21          # predictions
22          preds = X @ theta                # (m,1)
23          error = preds - y                # (m,1)
24          # gradient: (1/m) X^T (X theta - y) + (lambda/m)*[0; theta_1:]
25          grad = (1.0/m) * (X.T @ error) + (lam/m) * np.r_[[[0]], theta
                [1:]]  # do not reg bias
26          theta = theta - lr * grad
27          loss = (1.0/(2*m)) * np.sum(error**2) + (lam/(2*m))*np.sum(
                theta[1:]**2)
28          losses.append(loss)
```

```
29         if it % 200 == 0 and verbose:
30             print(f"GD iter {it}, loss {loss:.6f}")
31         if abs(prev_loss - loss) < tol:
32             if verbose:
33                 print(f"Converged at iter {it}, loss diff {abs(
                       prev_loss-loss):.2e}")
34             break
35         prev_loss = loss
36     return theta, losses
37
38 # Run GD without regularization (lambda=0) on STANDARDIZED data
39 theta_gd, losses_gd = gradient_descent(X_train_bias_scaled_gd, y_train,
       lr=0.01, epochs=10000, lam=0.0, tol=1e-10, verbose=False)
40 y_pred_gd = X_test_bias_scaled_gd @ theta_gd
41 mse_gd = mean_squared_error(y_test, y_pred_gd)
42 r2_gd = r2_score(y_test, y_pred_gd)
43 print("Gradient Descent -- standardized: MSE = {:.4f}, R2 = {:.4f}".
       format(mse_gd, r2_gd))
```

## 1.4 Ridge Regression (L2 Regularization)

The closed-form ridge solution modifies the Normal Equation to prevent matrix singularity and control overfitting:

$$\theta = (X^T X + \lambda I)^{-1} X^T y$$

We implement this using `np.linalg.solve` for better numerical stability than direct inversion. We tested this function and the GD function (with $\lambda > 0$) on both standardized and unstandardized data.

```
1 def ridge_closed_form(X, y, lam=1.0):
2     n = X.shape[1]
3     D = np.eye(n)
4     D[0,0] = 0.0 # Do not regularize the bias term
5     A = X.T @ X + lam * D
6     theta = np.linalg.solve(A, X.T @ y)   # more stable than inv
7     return theta
8
9 # Test 1: Ridge Closed-Form on UNSCALED data
10 theta_ridge_unstd = ridge_closed_form(X_train_bias, y_train, lam=1.0)
11 y_pred_ridge_unstd = X_test_bias @ theta_ridge_unstd
12 print("\nRidge (lam=1) unstandardized: MSE = {:.4f}, R2 = {:.4f}".
       format(
13     mean_squared_error(y_test, y_pred_ridge_unstd), r2_score(y_test,
           y_pred_ridge_unstd)))
14
15 # Test 2: Ridge using Gradient Descent on STANDARDIZED data (
       recommended)
16 theta_ridge_gd_std, _ = gradient_descent(X_train_bias_scaled_gd,
       y_train, lr=0.01, epochs=5000, lam=1.0, verbose=False)
17 y_pred_ridge_gd_std = X_test_bias_scaled_gd @ theta_ridge_gd_std
18 print("Ridge GD (lam=1) standardized: MSE = {:.4f}, R2 = {:.4f}".format
       (
19     mean_squared_error(y_test, y_pred_ridge_gd_std), r2_score(y_test,
           y_pred_ridge_gd_std)))
```

## 1.5 Performance Plots

Plots were generated to verify model performance and convergence.

- Predicted vs. Actual plots showed good alignment (using unscaled predictions for closed-form and GD-scaled predictions).

- Loss vs. Iterations plot confirmed gradient descent convergence.

- $\lambda$ vs. (MSE, $R^2$) plots demonstrated the effect of regularization strength on standardized data.

```python
# 1. GD Training Loss Plot
plt.figure(figsize=(5,3))
plt.plot(losses_gd)
plt.title("GD Training Loss (standardized)")
plt.xlabel("Iteration")
plt.ylabel("Loss")
plt.tight_layout()
plt.show()

# 2. Predicted vs Actual Plot
plt.figure(figsize=(6,6))
plt.scatter(y_test, y_pred_closed, alpha=0.6, label="Closed-form (pinv)
    - unscaled")
plt.scatter(y_test, y_pred_gd, alpha=0.6, label="Gradient Descent -
    scaled")
minv = min(y_test.min(), y_pred_closed.min(), y_pred_gd.min())
maxv = max(y_test.max(), y_pred_closed.max(), y_pred_gd.max())
plt.plot([minv, maxv], [minv, maxv], 'r--')
plt.xlabel("Actual Price")
plt.ylabel("Predicted Price")
plt.legend()
plt.title("Predicted vs Actual")
plt.show()

# 3. Effect of Lambda (on standardized features)
def evaluate_ridge_on_scaled(lam):
    theta = ridge_closed_form(X_train_scaled_bias, y_train, lam=lam)
    preds = X_test_scaled_bias @ theta
    return theta, preds

lams = [0.0, 0.01, 0.1, 1, 10, 100]
mse_vals = []
r2_vals  = []
thetas = {} # Store thetas for feature importance
for lam in lams:
    theta, preds = evaluate_ridge_on_scaled(lam)
    mse_vals.append(mean_squared_error(y_test, preds))
    r2_vals.append(r2_score(y_test, preds))
    thetas[lam] = theta

print("\nRidge on STANDARDIZED features (summary):")
for lam,mse,r2 in zip(lams, mse_vals, r2_vals):
    print(f"  lambda={lam:>6} : MSE={mse:.4f}, R2={r2:.4f}")

# 4. Lambda vs Metrics Plot
plt.figure(figsize=(6,4))
```

```
45  plt.plot(lams, mse_vals, marker='o', label="MSE")
46  # plt.plot(lams, r2_vals, marker='s', label="R2") # R2 is on different
        scale, plotting MSE only
47  plt.xscale('log')
48  plt.xlabel("Lambda (  )")
49  plt.ylabel("MSE")
50  plt.title("Effect of    on Ridge MSE (standardized)")
51  plt.legend()
52  plt.show()
```

## 1.6 Feature Importance

From standardized ridge coefficients (using $\lambda = 1.0$), we ranked features by their absolute coefficient values. This highlights which mobile phone attributes most influence the price prediction. Because the features are standardized (same scale), the magnitude of their coefficients is a direct proxy for importance.

```
1   # ***FIX: We must define feature_names from the dataframe columns***
2   feature_names = data.drop("Price", axis=1).columns
3
4   lam_for_importance = 1.0
5   theta_std = thetas[lam_for_importance].flatten()   # shape (d+1,)
6   coef_std = theta_std[1:]   # skip bias term
7
8   # Display numeric sorted importances
9   importance_df = pd.DataFrame({
10      "feature": feature_names,
11      "coef": coef_std,
12      "abs_coef": np.abs(coef_std)
13  }).sort_values(by="abs_coef", ascending=False)
14
15  print("\nFeature importance (standardized coefficients) sorted:")
16  print(importance_df.to_string(index=False))
17
18  # Plot feature importances
19  plt.figure(figsize=(10,4))
20  plt.bar(importance_df['feature'], importance_df['coef'])
21  plt.xticks(rotation=90)
22  plt.title(f"Feature importance (coefficients)    standardized features
        (  ={lam_for_importance})")
23  plt.tight_layout()
24  plt.show()
```

## Summary

Closed-form and gradient descent gave consistent results, especially when GD was run on standardized data. Ridge regression provided better generalization. Standardization is crucial for both GD convergence and for interpreting feature importances from model coefficients.

# 2. Linear Classification: Bank Note Authentication

We evaluated the suitability of logistic regression on the `BankNote_Authentication.csv` dataset.

## 2.1 Dataset Split

The dataset was divided into 70% training and 30% testing, using stratification to maintain the class balance in both splits.

```python
from sklearn.linear_model import LogisticRegression
from sklearn.metrics import accuracy_score
from mpl_toolkits.mplot3d import Axes3D

data_clf = pd.read_csv("/content/BankNote_Authentication.csv")
print("Dataset shape:", data_clf.shape)

X = data_clf.drop("class", axis=1)
y = data_clf["class"]

# Split into 70% train, 30% test
X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42, stratify=y
)

print("Training set size:", X_train.shape)
print("Test set size:", X_test.shape)
```

## 2.2 Logistic Regression without and with Regularization

We compared a standard L2 regularized model (C=1.0) against a model simulating no regularization (by setting $C = 1e10$, which makes the regularization penalty $\lambda = 1/C$ near zero).

- Without regularization: high training accuracy, slight overfitting observed.

- With L2 regularization: improved generalization, better balance between training and test accuracy.

```python
# 1. Approx "No Regulation" by using a massive C value
clf_no_reg = LogisticRegression(penalty="l2", C=1e10, solver="lbfgs",
    max_iter=1000)
clf_no_reg.fit(X_train, y_train)

y_pred_train_no_reg = clf_no_reg.predict(X_train)
y_pred_test_no_reg = clf_no_reg.predict(X_test)

print("Training Accuracy (No Reg approx):", accuracy_score(y_train,
    y_pred_train_no_reg))
print("Test Accuracy (No Reg approx):", accuracy_score(y_test,
    y_pred_test_no_reg))


# 2. Standard L2 Regularization (C=1.0)
```

```
13  clf_l2 = LogisticRegression(penalty="l2", C=1.0, solver="lbfgs",
        max_iter=1000)
14  clf_l2.fit(X_train, y_train)
15
16  y_pred_train_l2 = clf_l2.predict(X_train)
17  y_pred_test_l2 = clf_l2.predict(X_test)
18
19  print("\nTraining Accuracy (L2, C=1.0):", accuracy_score(y_train,
        y_pred_train_l2))
20  print("Test Accuracy (L2, C=1.0):", accuracy_score(y_test,
        y_pred_test_l2))
```

## 2.3 Accuracy vs. Regularization Strength

By varying $\lambda = 1/C$, we plotted train and test accuracies. We tested $\lambda$ values from $10^{-4}$ to $10^4$.

- Small $\lambda$ (weak regularization): risk of overfitting (Train acc ¿ Test acc).

- Large $\lambda$ (strong regularization): underfitting (Both accuracies drop as the model becomes too simple).

Optimal performance was observed at intermediate $\lambda$ values.

```
1  lambdas = np.logspace(-4, 4, 20)   #     = 1/C
2  train_acc = []
3  test_acc = []
4
5  for lam in lambdas:
6      C_val = 1.0 / lam
7      model = LogisticRegression(penalty="l2", C=C_val, solver="lbfgs",
            max_iter=1000)
8      model.fit(X_train, y_train)
9      train_acc.append(accuracy_score(y_train, model.predict(X_train)))
10     test_acc.append(accuracy_score(y_test, model.predict(X_test)))
11
12 plt.figure(figsize=(8,5))
13 plt.semilogx(lambdas, train_acc, label="Train Accuracy", marker="o")
14 plt.semilogx(lambdas, test_acc, label="Test Accuracy", marker="s")
15 plt.xlabel("   (Regularization strength = 1/C)")
16 plt.ylabel("Accuracy")
17 plt.title("Train & Test Accuracy vs    ")
18 plt.legend()
19 plt.grid(True)
20 plt.show()
```

## 2.4 Data Visualization

A 3D scatter plot using `variance`, `skewness`, and `curtosis` illustrated that the classes are highly (though perhaps not perfectly linearly) separable.

```
1  fig = plt.figure(figsize=(8,6))
2  ax = fig.add_subplot(111, projection='3d')
3
4  # Choose three features
```

```
5 ax.scatter(X_train["variance"], X_train["skewness"], X_train["curtosis"
      ],
6             c=y_train, cmap="bwr", alpha=0.7)
7
8 ax.set_xlabel("Variance")
9 ax.set_ylabel("Skewness")
10 ax.set_zlabel("Curtosis")
11 ax.set_title("3D Visualization of Bank Note Data")
12 plt.show()
```

## 2.5 Outlier Analysis

Artificial outliers were introduced by adding significant random noise to 20 data points in the training set. Impact: Training accuracy remained high (as the model tried to fit the noise), but test accuracy dropped, showing that the decision boundary was distorted by the outliers. This demonstrates that logistic regression is sensitive to outliers, which reduces its generalization capability.

```
1 # Copy dataset
2 X_outlier = X_train.copy()
3 y_outlier = y_train.copy()
4
5 # Introduce outliers by shifting some points
6 n_outliers = 20
7 rng = np.random.RandomState(42)
8 outlier_indices = rng.choice(len(X_outlier), size=n_outliers, replace=
      False)
9
10 # Add large noise to the selected outlier indices
11 X_outlier.iloc[outlier_indices] = X_outlier.iloc[outlier_indices] + rng
      .normal(20, 5, X_outlier.shape[1])
12
13 print("Introduced", n_outliers, "outliers.")
14
15 # Fit the standard L2 model on the noisy data
16 clf_outlier = LogisticRegression(penalty="l2", C=1.0, solver="lbfgs",
      max_iter=1000)
17 clf_outlier.fit(X_outlier, y_outlier)
18
19 train_acc_outlier = accuracy_score(y_outlier, clf_outlier.predict(
      X_outlier))
20 test_acc_outlier = accuracy_score(y_test, clf_outlier.predict(X_test))
      # Evaluate on the CLEAN test set
21
22 print("Training Accuracy with Outliers:", train_acc_outlier)
23 print("Test Accuracy with Outliers:", test_acc_outlier)
```

## Summary

Logistic regression with L2 regularization performed extremely well on the banknote authentication dataset, achieving near-perfect separation. The model is sensitive to the regularization parameter $\lambda$ and is negatively affected by outliers, which harm generalization performance.

# Conclusion

This assignment demonstrated regression and classification using linear models. Key insights:

- Linear regression is effective but sensitive to feature scaling (for GD) and benefits from regularization.

- Ridge regression balances the bias-variance tradeoff and, when used on standardized data, allows coefficients to be interpreted as feature importances.

- Logistic regression with regularization improves classification robustness by controlling overfitting.

- Outliers harm generalization in both models, underlining the importance of data preprocessing and outlier detection.