

Regression Models and Evaluation Techniques

July 3, 2025

1 Introduction to Regression

1.1 What is Regression?

Regression is the task of predicting continuous numerical values rather than discrete categories. It's one of the fundamental supervised learning tasks in machine learning, focusing on understanding relationships between variables.

1.2 Real-World Examples

- **House price prediction:** Estimating property values based on features
- **Stock price forecasting:** Predicting future market values
- **Sales forecasting:** Estimating future revenue
- **Temperature prediction:** Weather forecasting models
- **Medical dosage:** Determining optimal medication amounts

Today's Goals:

1. Learn 5 powerful regression models
2. Master evaluation techniques for regression
3. Understand when to use each approach

2 Data Preparation

2.1 Setting Up Practice Data

We'll create synthetic regression data to demonstrate each model and evaluation technique.

```
1 import numpy as np
2 import pandas as pd
3 from sklearn.datasets import make_regression
4 from sklearn.model_selection import train_test_split
5 from sklearn.preprocessing import StandardScaler
6 import matplotlib.pyplot as plt
7
8 # Create synthetic regression dataset
9 X, y = make_regression(n_samples=1000, n_features=10, noise=0.1, random_state
    =42)
10
11 # Split data (70% train, 30% test)
12 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.3,
    random_state=42)
```

```

13
14 # Scale features for models that need it
15 scaler = StandardScaler()
16 X_train_scaled = scaler.fit_transform(X_train)
17 X_test_scaled = scaler.transform(X_test)
18
19 print(f"Training set size: {X_train.shape}")
20 print(f"Test set size: {X_test.shape}")

```

Listing 1: Data Setup for Regression

3 Regression Models

3.1 Model 1: Linear Regression

How it works: Finds the best-fitting straight line (or hyperplane) through the data by minimizing the sum of squared residuals.

Mathematical foundation: $y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \dots + \beta_p x_p + \epsilon$

Best for: Linear relationships, interpretable coefficients

Limitations: Assumes linear relationships, sensitive to outliers

Example use: House price prediction based on size, location

```

1 from sklearn.linear_model import LinearRegression
2 from sklearn.metrics import mean_squared_error, r2_score
3
4 # Create and train model
5 linear_reg = LinearRegression()
6 linear_reg.fit(X_train, y_train)
7
8 # Make predictions
9 y_pred = linear_reg.predict(X_test)
10
11 # Evaluate performance
12 mse = mean_squared_error(y_test, y_pred)
13 r2 = r2_score(y_test, y_pred)
14
15 print(f"Linear Regression Results:")
16 print(f"MSE: {mse:.2f}")
17 print(f"R Score: {r2:.3f}")

```

Listing 2: Linear Regression Implementation

3.2 Model 2: Ridge Regression

How it works: Linear regression with L2 regularization that adds a penalty term to prevent overfitting by shrinking coefficients.

Mathematical foundation: Minimizes: $\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^p \beta_j^2$

Best for: High-dimensional data, multicollinearity issues

Limitations: Doesn't perform feature selection

Example use: Gene expression analysis, text analysis

```

1 from sklearn.linear_model import Ridge
2 from sklearn.model_selection import GridSearchCV
3
4 # Create model with regularization
5 ridge_reg = Ridge()
6
7 # Tune hyperparameter alpha
8 param_grid = {'alpha': [0.1, 1.0, 10.0, 100.0]}
9 ridge_cv = GridSearchCV(ridge_reg, param_grid, cv=5, scoring='
    neg_mean_squared_error')
10 ridge_cv.fit(X_train_scaled, y_train)
11
12 # Best model predictions
13 y_pred_ridge = ridge_cv.predict(X_test_scaled)
14
15 # Evaluate performance
16 mse_ridge = mean_squared_error(y_test, y_pred_ridge)
17 r2_ridge = r2_score(y_test, y_pred_ridge)
18
19 print(f"Ridge Regression Results:")
20 print(f"Best alpha: {ridge_cv.best_params_['alpha']}")
21 print(f"MSE: {mse_ridge:.2f}")
22 print(f"R    Score: {r2_ridge:.3f}")

```

Listing 3: Ridge Regression Implementation

3.3 Model 3: Lasso Regression

How it works: Linear regression with L1 regularization that can shrink coefficients to exactly zero, performing automatic feature selection.

Mathematical foundation: Minimizes: $\sum_{i=1}^n (y_i - \hat{y}_i)^2 + \alpha \sum_{j=1}^p |\beta_j|$

Best for: Feature selection, sparse models

Limitations: Can be unstable with correlated features

Example use: Biomarker discovery, economics modeling

```

1 from sklearn.linear_model import Lasso
2
3 # Create model with L1 regularization
4 lasso_reg = Lasso()
5
6 # Tune hyperparameter alpha
7 param_grid = {'alpha': [0.01, 0.1, 1.0, 10.0]}
8 lasso_cv = GridSearchCV(lasso_reg, param_grid, cv=5, scoring='
    neg_mean_squared_error')
9 lasso_cv.fit(X_train_scaled, y_train)
10
11 # Best model predictions
12 y_pred_lasso = lasso_cv.predict(X_test_scaled)
13
14 # Evaluate performance
15 mse_lasso = mean_squared_error(y_test, y_pred_lasso)
16 r2_lasso = r2_score(y_test, y_pred_lasso)
17
18 print(f"Lasso Regression Results:")
19 print(f"Best alpha: {lasso_cv.best_params_['alpha']}")
20 print(f"MSE: {mse_lasso:.2f}")
21 print(f"R    Score: {r2_lasso:.3f}")

```

```
22 print(f"Features selected: {np.sum(lasso_cv.best_estimator_.coef_ != 0)}")
```

Listing 4: Lasso Regression Implementation

3.4 Model 4: Random Forest Regression

How it works: Ensemble of decision trees that averages predictions from multiple trees, each trained on different subsets of data and features.

Best for: Non-linear relationships, robust performance

Limitations: Less interpretable, can overfit with small datasets

Example use: Real estate valuation, demand forecasting

```
1 from sklearn.ensemble import RandomForestRegressor
2
3 # Create random forest model
4 rf_reg = RandomForestRegressor(n_estimators=100, random_state=42)
5 rf_reg.fit(X_train, y_train)
6
7 # Make predictions
8 y_pred_rf = rf_reg.predict(X_test)
9
10 # Evaluate performance
11 mse_rf = mean_squared_error(y_test, y_pred_rf)
12 r2_rf = r2_score(y_test, y_pred_rf)
13
14 print(f"Random Forest Regression Results:")
15 print(f"MSE: {mse_rf:.2f}")
16 print(f"R Score: {r2_rf:.3f}")
17
18 # Feature importance
19 feature_importance = rf_reg.feature_importances_
20 print(f"Top 3 most important features: {np.argsort(feature_importance)[-3:]}")
```

Listing 5: Random Forest Regression Implementation

3.5 Model 5: Support Vector Regression (SVR)

How it works: Finds a function that deviates from actual values by at most ϵ while being as flat as possible.

Best for: High-dimensional data, non-linear relationships (with kernels)

Limitations: Sensitive to feature scaling, computationally expensive

Example use: Financial modeling, time series prediction

```
1 from sklearn.svm import SVR
2
3 # Create SVR model with RBF kernel
4 svr_reg = SVR(kernel='rbf', C=100, gamma=0.1, epsilon=0.1)
5 svr_reg.fit(X_train_scaled, y_train)
6
7 # Make predictions
8 y_pred_svr = svr_reg.predict(X_test_scaled)
9
10 # Evaluate performance
```

```

11 mse_svr = mean_squared_error(y_test, y_pred_svr)
12 r2_svr = r2_score(y_test, y_pred_svr)
13
14 print(f"Support Vector Regression Results:")
15 print(f"MSE: {mse_svr:.2f}")
16 print(f"R Score: {r2_svr:.3f}")

```

Listing 6: Support Vector Regression Implementation

4 Evaluation Techniques

4.1 Common Regression Metrics

4.1.1 Mean Squared Error (MSE)

Formula: $MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2$

Interpretation: Average of squared differences between actual and predicted values. Lower is better.

4.1.2 Root Mean Squared Error (RMSE)

Formula: $RMSE = \sqrt{MSE}$

Interpretation: Same units as target variable, easier to interpret than MSE.

4.1.3 Mean Absolute Error (MAE)

Formula: $MAE = \frac{1}{n} \sum_{i=1}^n |y_i - \hat{y}_i|$

Interpretation: Average absolute difference, less sensitive to outliers than MSE.

4.1.4 R-squared (R^2)

Formula: $R^2 = 1 - \frac{SS_{res}}{SS_{tot}}$

Interpretation: Proportion of variance explained by the model. Range: 0 to 1 (higher is better).

```

1 from sklearn.metrics import mean_absolute_error, mean_squared_error, r2_score
2 import numpy as np
3
4 def evaluate_regression_model(y_true, y_pred, model_name):
5     """
6     Comprehensive evaluation of regression model performance
7     """
8     mse = mean_squared_error(y_true, y_pred)
9     rmse = np.sqrt(mse)
10    mae = mean_absolute_error(y_true, y_pred)
11    r2 = r2_score(y_true, y_pred)
12
13    print(f"\n{model_name} Performance:")
14    print(f"{'='*50}")
15    print(f"MSE: {mse:.4f}")
16    print(f"RMSE: {rmse:.4f}")
17    print(f"MAE: {mae:.4f}")
18    print(f"R : {r2:.4f}")
19
20    return {'MSE': mse, 'RMSE': rmse, 'MAE': mae, 'R2': r2}
21
22 # Example usage for all models
23 results = {}

```

```

24 results['Linear'] = evaluate_regression_model(y_test, y_pred, "Linear
    Regression")
25 results['Ridge'] = evaluate_regression_model(y_test, y_pred_ridge, "Ridge
    Regression")
26 results['Lasso'] = evaluate_regression_model(y_test, y_pred_lasso, "Lasso
    Regression")
27 results['Random Forest'] = evaluate_regression_model(y_test, y_pred_rf, "Random
    Forest")
28 results['SVR'] = evaluate_regression_model(y_test, y_pred_svr, "Support Vector
    Regression")

```

Listing 7: Comprehensive Evaluation Function

4.2 Cross-Validation

Cross-validation provides more robust performance estimates by using multiple train/test splits.

```

1 from sklearn.model_selection import cross_val_score
2
3 def cross_validate_model(model, X, y, cv=5, scoring='neg_mean_squared_error'):
4     """
5     Perform cross-validation for regression model
6     """
7     scores = cross_val_score(model, X, y, cv=cv, scoring=scoring)
8
9     print(f"Cross-Validation Results (CV={cv}):")
10    print(f"Mean MSE: {-scores.mean():.4f}")
11    print(f"Std MSE: {scores.std():.4f}")
12    print(f"95% CI: [{-scores.mean() - 1.96*scores.std():.4f}, "
13          f"{-scores.mean() + 1.96*scores.std():.4f}]" )
14
15    return scores
16
17 # Example: Cross-validate Random Forest
18 rf_scores = cross_validate_model(RandomForestRegressor(n_estimators=100,
19    random_state=42),
20    X_train, y_train)

```

Listing 8: Cross-Validation for Regression

4.3 Residual Analysis

Analyzing residuals helps identify model assumptions violations and potential improvements.

```

1 import matplotlib.pyplot as plt
2
3 def plot_residuals(y_true, y_pred, model_name):
4     """
5     Create residual plots for regression analysis
6     """
7     residuals = y_true - y_pred
8
9     fig, axes = plt.subplots(1, 2, figsize=(12, 5))
10
11    # Residuals vs Predicted
12    axes[0].scatter(y_pred, residuals, alpha=0.6)
13    axes[0].axhline(y=0, color='r', linestyle='--')
14    axes[0].set_xlabel('Predicted Values')
15    axes[0].set_ylabel('Residuals')
16    axes[0].set_title(f'{model_name}: Residuals vs Predicted')
17
18    # Q-Q plot for normality

```

```

19 from scipy import stats
20 stats.probplot(residuals, dist="norm", plot=axes[1])
21 axes[1].set_title(f'{model_name}: Q-Q Plot')
22
23 plt.tight_layout()
24 plt.show()
25
26 # Statistical tests
27 print(f"\nResidual Analysis for {model_name}:")
28 print(f"Mean of residuals: {np.mean(residuals):.6f}")
29 print(f"Std of residuals: {np.std(residuals):.4f}")
30
31 # Example usage
32 plot_residuals(y_test, y_pred_rf, "Random Forest")

```

Listing 9: Residual Analysis

5 Model Comparison

Model	Strengths	Weaknesses	Best Use Cases
Linear Regression	Simple, interpretable, fast	Assumes linearity, sensitive to outliers	Linear relationships, baseline models
Ridge Regression	Handles multicollinearity, prevents overfitting	Doesn't select features	High-dimensional data, correlated features
Lasso Regression	Feature selection, sparse models	Unstable with correlated features	Feature selection, interpretable models
Random Forest	Handles non-linearity, robust	Less interpretable, can overfit	Complex relationships, robust predictions
SVR	Flexible with kernels, robust to outliers	Computationally expensive, requires scaling	High-dimensional data, non-linear patterns

Table 1: Comparison of Regression Models

6 Evaluation Metrics Comparison

7 Best Practices

7.1 Model Selection Guidelines

1. **Start simple:** Begin with linear regression as baseline
2. **Consider data size:** Use regularized models for small datasets
3. **Check assumptions:** Verify linearity, normality, homoscedasticity
4. **Handle multicollinearity:** Use Ridge regression when features are correlated
5. **Feature selection:** Use Lasso when you need interpretable models

Metric	Formula	Interpretation	When to Use
MSE	$\frac{1}{n} \sum (y_i - \hat{y}_i)^2$	Penalizes large errors heavily	When large errors are costly
RMSE	\sqrt{MSE}	Same units as target	General purpose, interpretable
MAE	$\frac{1}{n} \sum y_i - \hat{y}_i $	Robust to outliers	When outliers present
R ²	$1 - \frac{SS_{res}}{SS_{tot}}$	Proportion of variance explained	Model comparison, goodness of fit
MAPE	$\frac{100}{n} \sum \frac{ y_i - \hat{y}_i }{y_i}$	Percentage error	Relative performance

Table 2: Regression Evaluation Metrics

6. **Non-linear relationships:** Try Random Forest or SVR with appropriate kernels

7.2 Evaluation Best Practices

1. **Multiple metrics:** Don't rely on a single metric
2. **Cross-validation:** Use k-fold CV for robust estimates
3. **Residual analysis:** Check model assumptions
4. **Domain context:** Consider business/domain-specific requirements
5. **Overfitting check:** Compare training vs. validation performance

8 Conclusion

Regression modeling requires careful consideration of data characteristics, model assumptions, and evaluation metrics. Start with simple linear models to establish baselines, then explore more complex approaches as needed. Always validate your models using appropriate metrics and cross-validation techniques, and don't forget to analyze residuals to ensure model assumptions are met.

The key to successful regression modeling lies in understanding your data, choosing appropriate models for your specific use case, and thoroughly evaluating performance using multiple metrics and validation techniques.