

Lesson 21 Summary: Linear Regression

Data Science with Python – Key Concepts

Data Science Program

Linear Regression

OLS Formula

$y = Xw + b$
Minimize SSE

Assumptions

Linearity
Homoscedasticity

Interpretation

Coefficient = slope
p-value < 0.05

Model Fitting

Train/Test Split
Cross-validation

Diagnostics

Residual plots
Normality check

Finance Use

Beta estimation
Factor models

sklearn: LinearRegression().fit(X, y) | statsmodels: OLS for p-values

Linear regression is the foundation of predictive modeling

Ordinary Least Squares minimizes squared errors:

- **Objective:** Minimize $\sum (y_i - \hat{y}_i)^2$
- **Solution:** Closed-form: $w = (X^T X)^{-1} X^T y$
- **Intercept:** Controls baseline prediction level

Key insight:

OLS finds the best linear combination of features to predict the target.

OLS is the most common estimation method for linear regression

Linear regression requires:

- **Linearity:** True relationship is linear
- **Independence:** Observations are independent
- **Homoscedasticity:** Constant error variance
- **Normality:** Errors are normally distributed

Violations of assumptions affect coefficient reliability

Understanding regression outputs:

- **Coefficient:** Change in y per unit change in x
- **p-value:** Probability coefficient is zero
- **Confidence interval:** Range of plausible values

Rule of thumb:

p-value \leq 0.05 typically considered significant.

Use statsmodels for detailed statistical inference

Two main approaches:

- **sklearn**: Fast fitting, no p-values
- **statsmodels**: Full statistical output

sklearn pattern:

```
model = LinearRegression()
model.fit(X_train, y_train)
predictions = model.predict(X_test)
```

Choose sklearn for ML, statsmodels for inference

Prevent overfitting:

- **Training set:** 70-80% of data for fitting
- **Test set:** 20-30% for evaluation
- **Random state:** Ensures reproducibility

Never evaluate on training data!

Test set performance shows real-world capability

Check model assumptions:

- **Residual plot:** Should show no pattern
- **Q-Q plot:** Checks normality of errors
- **Scale-location:** Tests homoscedasticity

Always visualize residuals before trusting results

Estimating systematic risk:

- **CAPM:** $r_i - r_f = \alpha + \beta(r_m - r_f) + \epsilon$
- **Beta:** Sensitivity to market movements
- **Alpha:** Excess return (manager skill)

Beta estimation is a direct application of OLS

Avoid these mistakes:

- **Multicollinearity:** Correlated features
- **Outliers:** Distort coefficient estimates
- **Overfitting:** Too many features

Simple models often outperform complex ones

Essential Commands:

Task	Code
Split data	<code>train_test_split(X, y)</code>
Fit model	<code>LinearRegression().fit(X, y)</code>
Predict	<code>model.predict(X_new)</code>
Coefficients	<code>model.coef_, model.intercept_</code>
R-squared	<code>model.score(X, y)</code>

Master these basics before moving to regularization