

Cheat Sheet: Linear and Logistic Regression

Comparing different regression types

Model Name	Description	Code Syntax
SIMPLE LINEAR REGRESSION	Purpose: TO PREDICT A DEPENDENT VARIABLE BASED ON ONE INDEPENDENT VARIABLE. Pros: EASY TO IMPLEMENT, INTERPRET, AND EFFICIENT FOR SMALL DATASETS. Cons: NOT SUITABLE FOR COMPLEX RELATIONSHIPS; PRONE TO UNDERFITTING. Modeling equation: $Y = B_0 + B_1X$	<pre>from sklearn.linear_model import LinearRegression model = LinearRegression() model.fit(X, y)</pre>
POLYNOMIAL REGRESSION	Purpose: TO CAPTURE NONLINEAR RELATIONSHIPS BETWEEN VARIABLES. Pros: BETTER AT FITTING NONLINEAR DATA COMPARED TO LINEAR REGRESSION. Cons: PRONE TO OVERFITTING WITH HIGH-DEGREE POLYNOMIALS. Modeling equation: $Y = B_0 + B_1X + B_2X^2 + \dots$	<pre>from sklearn.preprocessing import PolynomialFeatures from sklearn.linear_model import LinearRegression poly = PolynomialFeatures(degree=2) X_poly = poly.fit_transform(X) model = LinearRegression().fit(X_poly, y)</pre>
MULTIPLE LINEAR REGRESSION	Purpose: TO PREDICT A DEPENDENT VARIABLE BASED ON MULTIPLE INDEPENDENT VARIABLES. Pros: ACCOUNTS FOR MULTIPLE FACTORS INFLUENCING THE OUTCOME. Cons: ASSUMES A LINEAR RELATIONSHIP BETWEEN PREDICTORS AND TARGET. Modeling equation: $Y = B_0 + B_1X_1 + B_2X_2 + \dots$	<pre>from sklearn.linear_model import LinearRegression model = LinearRegression() model.fit(X, y)</pre>
LOGISTIC REGRESSION	Purpose: TO PREDICT PROBABILITIES OF CATEGORICAL OUTCOMES. Pros: EFFICIENT FOR BINARY CLASSIFICATION PROBLEMS. Cons: ASSUMES A LINEAR RELATIONSHIP BETWEEN INDEPENDENT VARIABLES AND LOG-ODDS. Modeling equation: $\text{LOG}(P/(1-P)) = B_0 + B_1X_1 + \dots$	<pre>from sklearn.linear_model import LogisticRegression model = LogisticRegression() model.fit(X, y)</pre>

Associated functions commonly used

Function/Method Name	Brief Description	Code Syntax
TRAIN_TEST_SPLIT	SPLITS THE DATASET INTO TRAINING AND TESTING SUBSETS TO EVALUATE THE MODEL'S PERFORMANCE.	<pre>from sklearn.model_selection import train_test_split X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=42)</pre>
STANDARDSCALER	STANDARDIZES FEATURES BY REMOVING THE MEAN AND SCALING TO UNIT VARIANCE.	<pre>from sklearn.preprocessing import StandardScaler scaler = StandardScaler() X_scaled = scaler.fit_transform(X)</pre>
LOG_LOSS	CALCULATES THE LOGARITHMIC LOSS, A PERFORMANCE METRIC FOR CLASSIFICATION MODELS.	<pre>from sklearn.metrics import log_loss loss = log_loss(y_true, y_pred_proba)</pre>
MEAN_ABSOLUTE_ERROR	CALCULATES THE MEAN ABSOLUTE ERROR BETWEEN ACTUAL AND PREDICTED VALUES.	<pre>from sklearn.metrics import mean_absolute_error mae = mean_absolute_error(y_true, y_pred)</pre>
MEAN_SQUARED_ERROR	COMPUTES THE MEAN SQUARED ERROR BETWEEN ACTUAL AND PREDICTED VALUES.	<pre>from sklearn.metrics import mean_squared_error mse = mean_squared_error(y_true, y_pred)</pre>
ROOT_MEAN_SQUARED_ERROR	CALCULATES THE ROOT MEAN SQUARED ERROR (RMSE), A COMMONLY USED METRIC FOR REGRESSION TASKS.	<pre>from sklearn.metrics import mean_squared_error import numpy as np rmse = np.sqrt(mean_squared_error(y_true, y_pred))</pre>
R2_SCORE	COMPUTES THE R-SQUARED VALUE, INDICATING HOW WELL THE MODEL EXPLAINS THE VARIABILITY OF THE TARGET VARIABLE.	<pre>from sklearn.metrics import r2_score r2 = r2_score(y_true, y_pred)</pre>

Author(s)

JEFF GROSSMAN
ABHISHEK GAGNEJA



Skills Network