

Data Science Canvas		Project:	Online News Popularity				
		Team:	Rajat Chaudhary, Srinivas Shavukapu Kattegummula, Sonu Goyal, Yuvaraj G				
Problem Statement				Execution & Evaluation		Data Collection & Preparation	
<b>Business Case &amp; Value Added</b> Media companies and content creators struggle to predict which articles will gain traction online. Accurate predictions help optimize content strategy, improve engagement and increase ad revenue by focusing on high-impact articles.	<b>Model Selection</b>  <b>Regression models</b> Linear regression, logistic regression, Random Forest, XGBoost, KNN  <b>Classification models</b> XGBoostClassifier, LightGBM, RandomForest, Adaboost, Naive Bayes, Stacked Ensemble	<b>Model Requirements</b> Models must handle multivariate, tabular data with both numerical and categorical features. Should be robust to skewed distributions and class imbalance.	<b>Skills</b> Python skills for data cleaning, feature engineering, and EDA; ML expertise for model building, tuning, and evaluation. Deployment and interpretability using Streamlit/Flask, SHAP/LIME, plus Git for collaboration.	<b>Model Evaluation</b> Use MAE, MSE, RMSE for regression; accuracy, precision, F1-score, macro-averaged metrics for classification. Monitor model performance, especially on minority classes (Trending, Viral).	<b>Data Storytelling</b> Present results with clear visuals (charts, tables). Focus on actionable insights for editors (e.g., which features drive virality). Communicate limitations and next steps.	<b>Data Selection &amp; Cleansing</b> Use all relevant features except identifiers (e.g., URL). Remove or correct outliers (e.g., invalid ratios, negative keyword shares). One-hot encode categorical, standardize numerical.	<b>Data Collection</b> Dataset is pre-collected from UCI; no additional data required. Ensure data meets requirements: sufficient size, relevant features, no missing values.
				<b>Data Integration</b> All data is from a single source (UCI dataset); no integration needed.	<b>Explorative Data Analysis</b> Check for skewness, outliers, and feature distributions. Use univariate, bivariate, and multivariate analysis to understand relationships. Identify that popularity is driven by complex, non-linear interactions.		