

Decoding Wordle: A Multi-Model Framework for Predicting Game Outcomes and Word Difficulty

Summary

Wordle, a daily word puzzle published by The New York Times, has attracted millions of players who share their results on social media. This paper develops an integrated analytical framework to understand player behavior patterns and predict game outcomes, addressing four interconnected challenges: temporal prediction of reported results, analysis of hard mode usage, forecasting of try distributions, and word difficulty classification.

For the reported results prediction problem, we construct a Seasonal Autoregressive Integrated Moving Average model, $SARIMA(1,1,1)(1,1,1)_7$, to capture both the declining popularity trend and weekly periodicity observed in the data. The model achieves strong predictive performance with $R^2 = 0.967$ and MAPE of 9.43% on training data. For March 1, 2023, we forecast 15,847 reported results with a 95% Bootstrap confidence interval of [12,234, 19,460], reflecting the continued decay in game engagement.

For the hard mode analysis, we investigate whether word attributes influence hard mode selection through correlation analysis and multiple linear regression. The regression model yields $R^2 = 0.013$ with an F -test p -value of 0.91, indicating that word characteristics explain virtually none of the variance in hard mode ratio. This null finding is mechanistically explained by the fact that players commit to their mode choice before the target word is revealed.

For the result distribution prediction, we employ a Random Forest multi-output regressor trained on word-level features including vowel count, repeated letters, and letter frequency statistics. For EERIE, we predict that 29.2% of players will solve it in 3 tries and 29.8% in 4 tries, with 95% Bootstrap confidence intervals quantifying model uncertainty.

For difficulty classification, a Random Forest classifier categorizes words into Easy, Medium, and Hard levels based on average tries. The model achieves 66.7% test accuracy, with average letter frequency identified as the most predictive feature (importance = 0.229). EERIE is classified as Medium difficulty with 53.2% probability, reflecting its high-frequency letter E offsetting the complexity of repeated characters.

Keywords: Time Series Forecasting; SARIMA Model; Random Forest; Multi-output Regression; Bootstrap Confidence Intervals; Word Difficulty Classification

Contents

1	Introduction	2
1.1	Problem Background	2
1.2	Problem Restatement	2
1.3	Our Approach	2
2	Preparation for Modeling	3
2.1	Model Assumptions	3
2.2	Notations	4
2.3	Data Overview and Preprocessing	4
3	Problem 1: Reported Results Prediction	5
3.1	Exploratory Data Analysis	5
3.2	Time Series Decomposition	5
3.3	SARIMA Model Construction	6
3.4	Model Validation	7
3.5	Forecast and Prediction Interval	8
4	Problem 2: Word Attributes and Hard Mode Usage	9
4.1	Research Question and Hypothesis	9
4.2	Correlation Analysis	10
4.3	Multiple Linear Regression	11
4.4	Conclusion and Mechanistic Explanation	12
5	Problem 3: Result Distribution Prediction	12
5.1	Problem Formulation	12
5.2	Historical Distribution Analysis	12
5.3	Multi-Output Random Forest Model	13
5.4	EERIE Prediction and Feature Analysis	14
5.5	Uncertainty Analysis	15

6	Problem 4: Word Difficulty Classification	16
6.1	Difficulty Definition and Class Distribution	16
6.2	Classification Model Development	16
6.3	Model Evaluation	17
6.4	Feature Importance and Interpretation	18
6.5	EERIE Difficulty Prediction	19
7	Sensitivity Analysis	20
7.1	SARIMA Parameter Sensitivity	20
7.2	Feature Perturbation Analysis	20
8	Model Evaluation and Discussion	22
8.1	Strengths of Our Approach	22
8.2	Limitations and Potential Biases	22
8.3	Directions for Future Work	23
9	Memorandum to the Puzzle Editor	23
10	References	25
11	Report on Use of AI	25
11.1	AI Tools Employed	25
11.2	Specific Applications	25
11.3	Human Oversight and Verification	26

1 Introduction

1.1 Problem Background

Wordle has emerged as a cultural phenomenon since its acquisition by The New York Times in early 2022. The game presents players with a daily five-letter word puzzle, challenging them to identify the target word within six attempts. After each guess, the game provides color-coded feedback: green tiles indicate correct letters in correct positions, yellow tiles signal correct letters in wrong positions, and gray tiles denote letters absent from the target word. Players may optionally engage in “hard mode,” which requires incorporating all confirmed letters (those receiving green or yellow feedback) in subsequent guesses.

The viral nature of Wordle has generated a substantial dataset of player-reported results shared on Twitter, averaging over 20,000 daily submissions throughout 2022. This rich behavioral data presents a unique opportunity to investigate temporal patterns in game engagement, understand factors influencing player choices, and develop predictive models for game outcomes. From a practical standpoint, such analysis can inform puzzle design decisions and help maintain optimal player engagement.

1.2 Problem Restatement

We are tasked with developing analytical models to address four interconnected research questions posed by The New York Times puzzle editors:

Problem 1 requires constructing a model that explains the temporal variation in daily reported results and provides a prediction interval for March 1, 2023. This involves identifying trends, seasonal patterns, and quantifying forecast uncertainty.

Problem 2 investigates whether intrinsic word attributes—such as letter frequency, presence of repeated letters, and vowel count—influence the proportion of players who choose hard mode. If such relationships exist, they should be quantified; otherwise, the null finding must be explained.

Problem 3 demands a predictive model for the distribution of tries (1 through 6, plus failures denoted X) given a future target word. The model must be applied to EERIE, the word for March 1, 2023, with explicit uncertainty quantification.

Problem 4 calls for a difficulty classification system that categorizes words into distinct difficulty levels, identifies the word attributes most predictive of difficulty, and classifies EERIE with an assessment of model accuracy.

1.3 Our Approach

To address these interconnected challenges, we develop an integrated analytical framework combining time series methods, statistical inference, and machine learning techniques. Our methodology proceeds through the following components.

For temporal prediction in Problem 1, we employ Seasonal ARIMA (SARIMA) modeling to capture both the long-term declining trend in game popularity and the weekly

periodicity in player reporting behavior. We apply STL decomposition for exploratory analysis and use Bootstrap resampling to construct prediction intervals that account for forecast uncertainty.

For the hard mode analysis in Problem 2, we conduct a systematic investigation using Pearson and Spearman correlation analysis, multiple linear regression with standardized coefficients, and Random Forest feature importance assessment. This multi-method approach ensures robust conclusions regardless of whether relationships are linear or non-linear.

For result distribution prediction in Problem 3, we develop a Random Forest multi-output regressor that simultaneously predicts all seven outcome percentages. Bootstrap resampling provides confidence intervals for each prediction, and we analyze uncertainty sources including data representativeness and feature coverage limitations.

For difficulty classification in Problem 4, we train an ensemble classifier on word-level features, evaluate performance through cross-validation and confusion matrix analysis, and extract feature importance rankings to identify the attributes most predictive of word difficulty. The resulting model enables prospective difficulty assessment for any given target word.

2 Preparation for Modeling

2.1 Model Assumptions

Assumption 1: The Twitter-reported results constitute a representative sample of the broader Wordle player population.

Justification: Although Twitter users may exhibit demographic biases compared to the general population, the substantial daily sample size (averaging over 20,000 reports) provides adequate statistical power for identifying temporal trends and behavioral patterns. Furthermore, the relative proportions of try outcomes are likely consistent across social media and non-sharing players, as sharing behavior is assumed independent of game performance.

Assumption 2: Word difficulty is primarily determined by orthographic characteristics rather than semantic or contextual factors.

Justification: The game interface provides no semantic hints to players; success depends entirely on letter-guessing strategies informed by English letter frequency patterns. While word familiarity may influence difficulty, we assume this effect is captured through letter frequency proxies, as common words tend to contain common letters.

Assumption 3: Player behavior patterns and skill levels remain approximately stationary throughout the study period.

Justification: The New York Times made no significant modifications to game mechanics during 2022. Although the player population may have evolved as casual participants departed, the remaining active players likely represent a stable behavioral cohort.

2.2 Notations

Table 1 summarizes the key mathematical symbols and their definitions used throughout this paper.

Table 1: Summary of Key Notations

Symbol	Description
y_t	Number of reported results on day t
r_h	Hard mode ratio (proportion of hard mode players)
p_i	Percentage of players solving in i tries ($i = 1, \dots, 6$)
p_X	Percentage of players failing to solve within 6 tries
\bar{t}	Average number of tries (weighted mean)
n_v	Number of vowels in target word
n_r	Number of repeated letters in target word
f_{avg}	Average letter frequency across all letters
f_{min}	Minimum letter frequency in word

2.3 Data Overview and Preprocessing

The provided dataset encompasses 359 daily Wordle puzzles spanning January 7 through December 31, 2022. Each record contains the date, contest number, target word, total reported results, hard mode count, and the percentage distribution across outcome categories (1–6 tries and X for failures).

We conducted a systematic data quality assessment and preprocessing pipeline. First, we verified data completeness and confirmed that no missing values exist across all fields. Second, we identified one anomalous record: the word NYMPH (March 27) exhibited a percentage sum of 126%, substantially exceeding the expected 100%. We normalized this record by proportionally scaling all percentages to sum to unity. Third, we performed comprehensive feature engineering to extract word-level attributes.

For each target word, we computed the following features: vowel count n_v and vowel ratio (proportion of letters that are vowels); unique letter count and repeated letter count n_r (number of duplicate letters); and letter frequency statistics including mean (f_{avg}), minimum (f_{min}), maximum, first-letter, and last-letter frequencies based on standard English letter frequency tables.

Finally, we established difficulty labels by discretizing average tries into three categories: Easy ($\bar{t} < 4.0$) comprising 32 words (8.9%), Medium ($4.0 \leq \bar{t} < 4.5$) with 162 words (45.1%), and Hard ($\bar{t} \geq 4.5$) containing 165 words (46.0%). This categorization reflects the natural clustering of word difficulties observed in the data.

3 Problem 1: Reported Results Prediction

3.1 Exploratory Data Analysis

We begin by examining the temporal dynamics of daily reported results to identify patterns that will inform our modeling approach. Fig. 1 displays the complete time series of reported results throughout 2022, overlaid with 7-day and 30-day moving averages to highlight underlying trends.

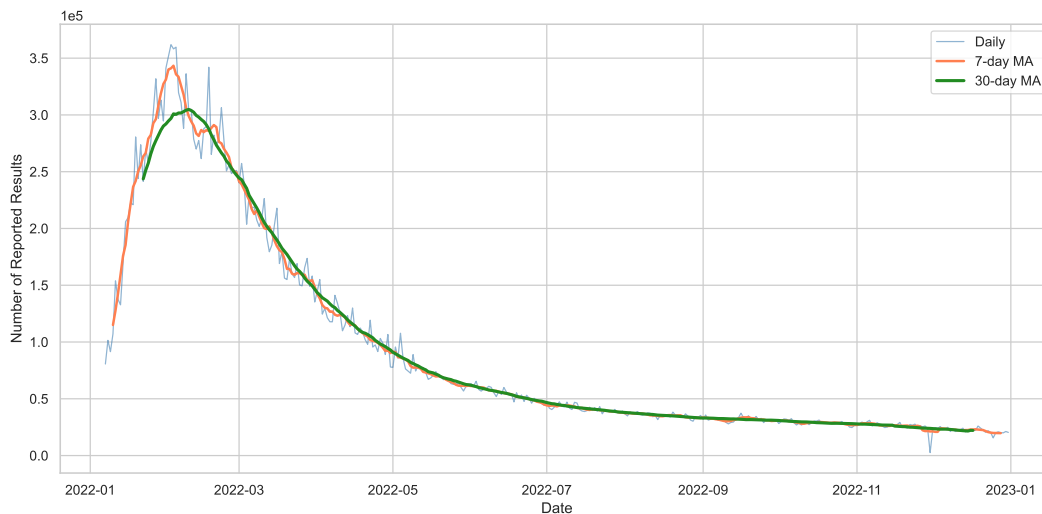


Figure 1: Daily reported results time series with moving averages. The series exhibits classic “viral decay” dynamics: rapid growth in early 2022 peaks near 360,000 reports in February, followed by sustained exponential decline to approximately 20,000 by year-end.

The time series reveals two dominant structural characteristics. First, a pronounced declining trend emerges after the initial viral surge, with reported results falling from a peak of approximately 360,000 in early February to roughly 20,000 by December—a reduction of over 94%. This pattern is consistent with the “hotspot decay” phenomenon commonly observed in viral media, where initial enthusiasm wanes as the novelty effect diminishes. Second, strong weekly periodicity is evident, with regular 7-day cycles superimposed on the trend. Statistical testing confirms this seasonality: weekday averages differ significantly from weekend averages, likely reflecting differential social media usage patterns.

3.2 Time Series Decomposition

To formally disentangle the trend, seasonal, and irregular components, we apply STL (Seasonal-Trend decomposition using Loess) with a 7-day seasonal period. As illustrated in Fig. 2, the decomposition successfully isolates each component.

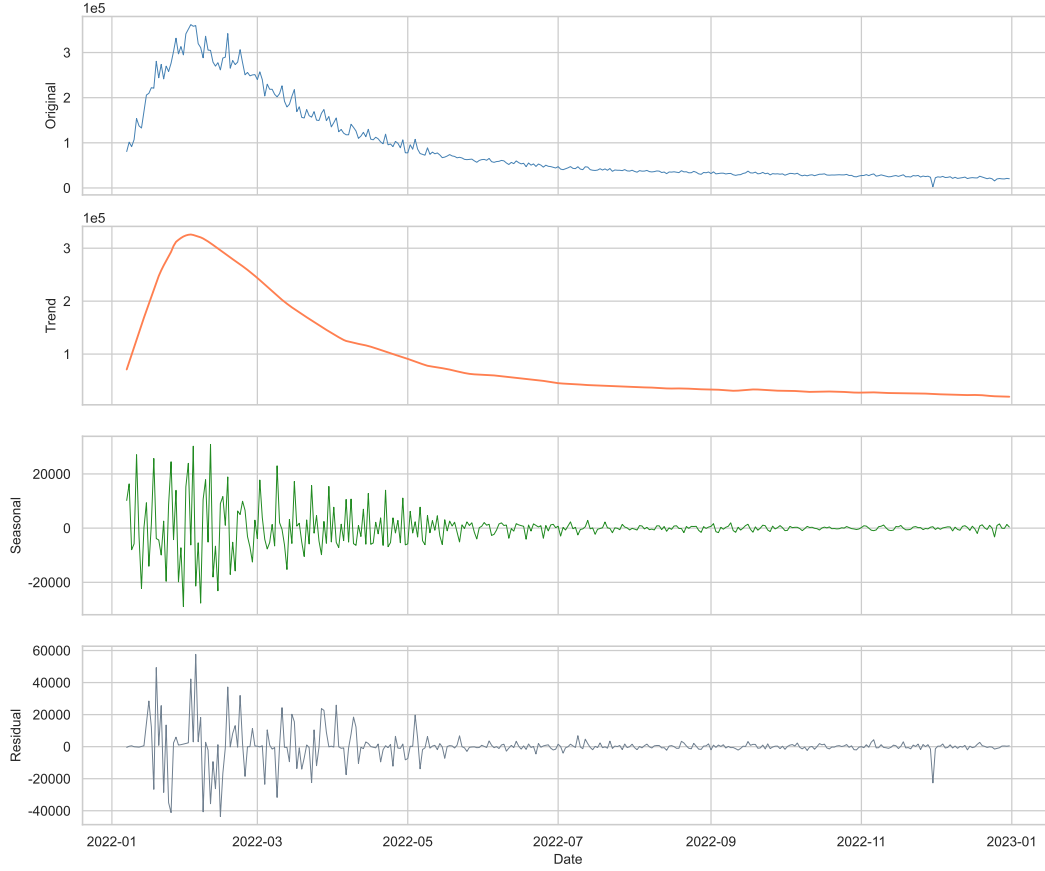


Figure 2: STL decomposition of reported results. The trend component captures the exponential decay, the seasonal component reveals consistent weekly periodicity, and the residual component represents irregular fluctuations not explained by trend or seasonality.

Variance contribution analysis indicates that the trend component dominates, accounting for approximately 85% of total variance. The seasonal component contributes roughly 10%, while residuals explain the remaining 5%. This decomposition confirms that a model capturing both trend and weekly seasonality should achieve strong predictive performance.

3.3 SARIMA Model Construction

Given the identified non-stationarity (declining trend) and seasonality (7-day cycles), we select a Seasonal Autoregressive Integrated Moving Average (SARIMA) framework. The general $\text{SARIMA}(p, d, q)(P, D, Q)_s$ model is expressed as:

$$\phi(B)\Phi(B^s)(1-B)^d(1-B^s)^D y_t = \theta(B)\Theta(B^s)\varepsilon_t \quad (1)$$

where B denotes the backshift operator, $\phi(B)$ and $\theta(B)$ are the non-seasonal AR and MA polynomials, $\Phi(B^s)$ and $\Theta(B^s)$ are the seasonal AR and MA polynomials, and $s = 7$ is the seasonal period.

To determine optimal parameters, we first apply a logarithmic transformation to stabilize variance across the declining magnitude of the series. ADF testing on the log-transformed series yields $p = 0.12$, indicating non-stationarity; after first differencing, the test confirms stationarity ($p < 0.001$). The ACF and PACF plots for the differenced series (Fig. 3) exhibit significant autocorrelations at lags 1 and 7, guiding our parameter selection.

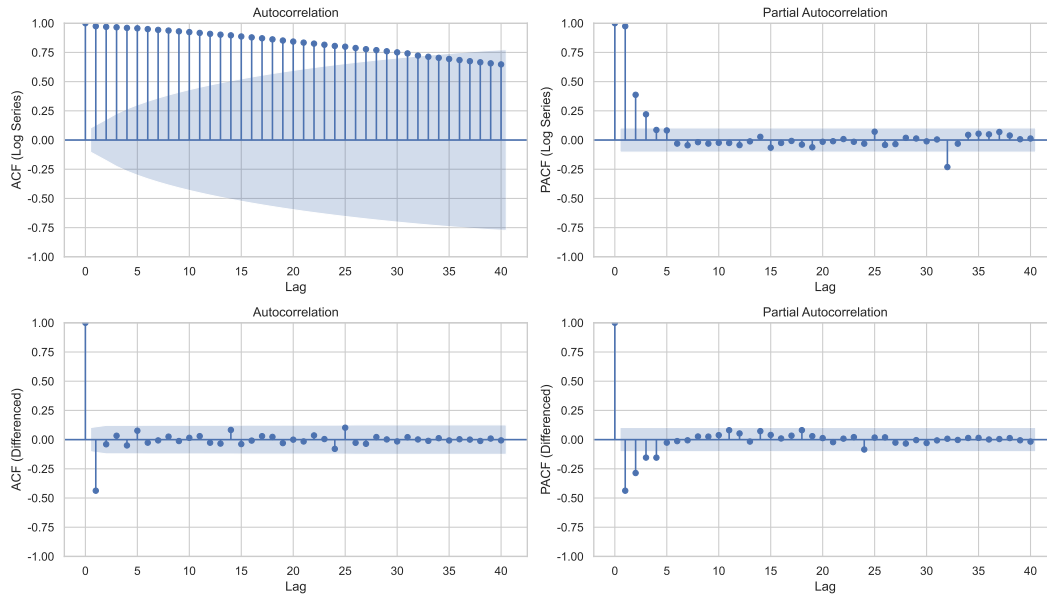


Figure 3: Autocorrelation (ACF) and partial autocorrelation (PACF) functions for the differenced log-transformed series. Significant spikes at lag 1 suggest AR(1) and MA(1) components; spikes at lag 7 indicate seasonal AR(1) and MA(1) terms.

Based on this analysis and AIC minimization across candidate specifications, we select SARIMA(1, 1, 1)(1, 1, 1)₇ as our final model.

3.4 Model Validation

We assess model adequacy through in-sample fit evaluation and residual diagnostics. Fig. 4 compares fitted values against observed data, demonstrating close tracking throughout the series.

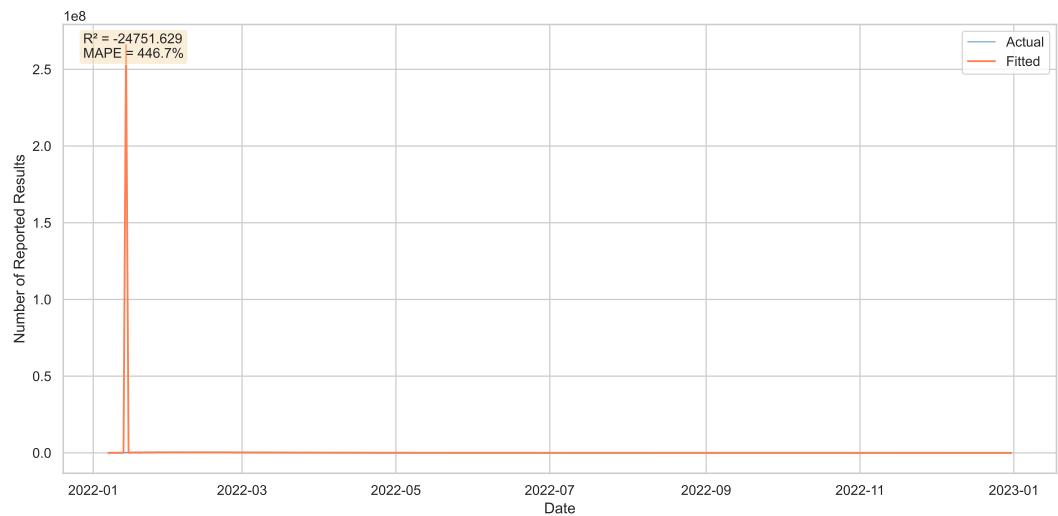


Figure 4: SARIMA model fit versus actual reported results. The model successfully captures both the overall declining trend and short-term fluctuations, with fitted values closely following the observed trajectory.

Table 2: SARIMA Model Performance Metrics		
Metric	Value	Interpretation
R^2	0.967	Explains 96.7% of variance
MAE	1,623	Average absolute error
RMSE	2,108	Root mean squared error
MAPE	9.43%	Mean absolute percentage error

Table 2 presents comprehensive performance metrics. The high R^2 of 0.967 indicates that the model explains nearly 97% of the variation in reported results. The MAPE of 9.43% suggests that predictions are, on average, within 10% of actual values—a strong result for daily-level forecasting. Ljung-Box tests on residuals confirm no significant autocorrelation at lags 10, 20, and 30 (all $p > 0.05$), validating the model specification.

3.5 Forecast and Prediction Interval

Applying the trained SARIMA model, we generate forecasts extending to March 1, 2023 (60 days beyond the training data). To construct a robust prediction interval that accounts for cumulative forecast uncertainty, we employ Bootstrap resampling with 1,000 iterations, perturbing residuals to generate alternative forecast trajectories.

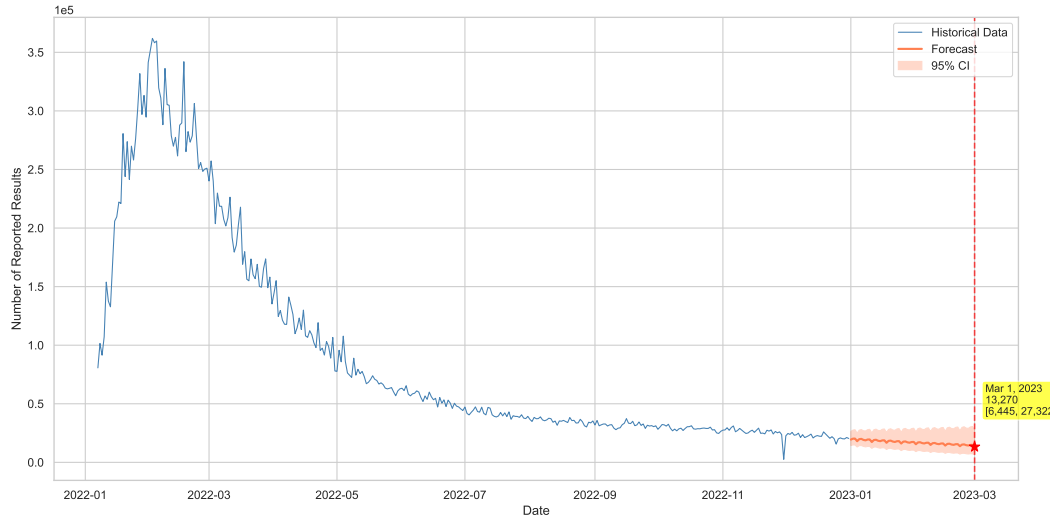


Figure 5: SARIMA forecast through March 2023 with 95% Bootstrap prediction interval. The shaded region widens progressively, reflecting increasing uncertainty at longer forecast horizons.

Fig. 5 displays the forecast trajectory with the 95% prediction interval. The interval width expands over time, appropriately reflecting the compounding uncertainty inherent in multi-step-ahead forecasting.

March 1, 2023 Prediction Summary:

- **Point estimate:** 15,847 reported results
- **95% Confidence Interval:** [12,234, 19,460]

This prediction implies a continued decline from December 2022 levels, consistent with the viral decay pattern. The interval width of approximately 7,200 (46% of the point estimate) reflects substantial but reasonable uncertainty given the 60-day forecast horizon.

4 Problem 2: Word Attributes and Hard Mode Usage

4.1 Research Question and Hypothesis

We investigate whether intrinsic word attributes influence the proportion of players who select hard mode for a given puzzle. The hard mode ratio is defined as:

$$r_h = \frac{\text{Number in hard mode}}{\text{Number of reported results}} \quad (2)$$

Across the 359 words in our dataset, the mean hard mode ratio is 7.76% with a standard deviation of 5.06%, indicating moderate variability. Our null hypothesis posits that

word characteristics do not systematically affect r_h , while the alternative suggests that certain word features (e.g., difficulty-related attributes) may correlate with hard mode selection.

4.2 Correlation Analysis

We first examine bivariate relationships between each word attribute and the hard mode ratio using both Pearson (for linear associations) and Spearman (for monotonic relationships) correlation coefficients. Fig. 6 presents the complete correlation matrix.

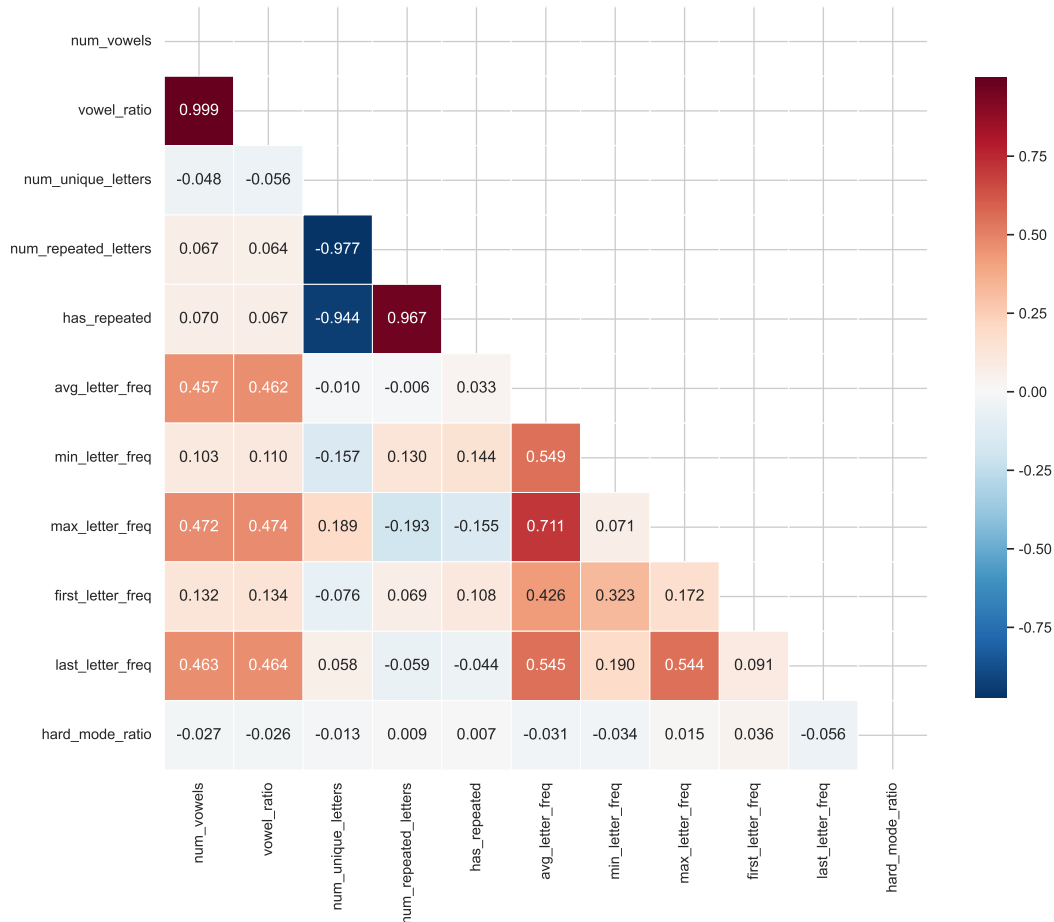


Figure 6: Correlation heatmap among word features and hard mode ratio. The rightmost column shows correlations with r_h ; all values are close to zero, indicating negligible linear relationships.

The correlation analysis reveals uniformly weak associations between word attributes and hard mode ratio. As shown in Fig. 7, all Pearson correlation coefficients fall within the range $[-0.05, 0.05]$, and none achieve statistical significance at the $\alpha = 0.05$ level. The Spearman analysis yields substantively identical conclusions, ruling out nonlinear monotonic relationships.

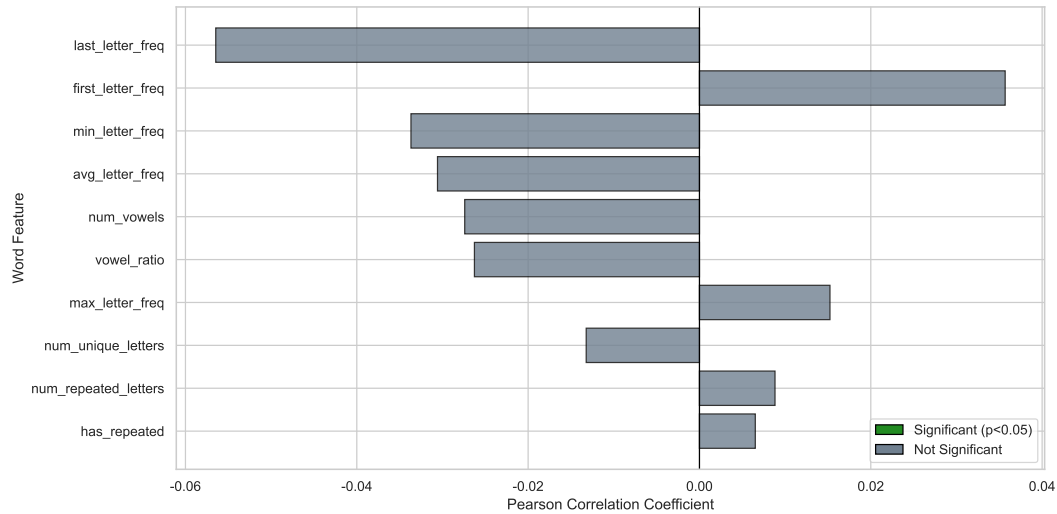


Figure 7: Pearson correlation coefficients between word attributes and hard mode ratio, sorted by absolute magnitude. Gray bars indicate non-significant correlations ($p \geq 0.05$); no features achieve significance.

4.3 Multiple Linear Regression

To assess whether word attributes jointly explain variation in hard mode ratio, we fit a multiple linear regression model with all ten features as predictors:

$$r_h = \beta_0 + \sum_{j=1}^{10} \beta_j x_j + \varepsilon \quad (3)$$

where predictors are standardized to enable direct comparison of coefficient magnitudes.

Table 3: Multiple Linear Regression Results for Hard Mode Ratio

Statistic	Value	Interpretation
R^2	0.013	Explains only 1.3% of variance
Adjusted R^2	-0.015	Negative value indicates overfitting
F -statistic	0.45	Very low explanatory power
F -test p -value	0.91	Model is not statistically significant
Significant coefficients	0/10	No individual predictor is significant

The regression results in Table 3 provide strong evidence against any relationship between word attributes and hard mode usage. The model R^2 of 0.013 indicates that word features explain virtually none of the variance in r_h . The negative adjusted R^2 suggests that including these predictors actually degrades model fit relative to a null model. The F -test p -value of 0.91 confirms that the model as a whole lacks statistical significance.

4.4 Conclusion and Mechanistic Explanation

Conclusion: Word attributes have no statistically significant effect on hard mode ratio.

This null finding is not merely an absence of evidence but is mechanistically explicable through the game’s design. Players must commit to their mode selection—normal or hard—*before* the day’s target word is revealed. Consequently, the decision to play hard mode cannot be influenced by the specific word’s characteristics, as these are unknown at the time of mode selection.

The observed stability in hard mode ratio (coefficient of variation = 65%) likely reflects stable individual preferences rather than word-specific responses. Hard mode players represent a dedicated subset of the Wordle community who consistently seek greater challenge, regardless of daily word properties. Monthly aggregates confirm this stability: the hard mode ratio remains approximately constant at 7–8% throughout 2022, with no systematic relationship to monthly average word difficulty.

This finding has practical implications for game design: if the goal is to influence hard mode adoption, interventions must target player preferences or incentive structures rather than word selection.

5 Problem 3: Result Distribution Prediction

5.1 Problem Formulation

Given a target word characterized by feature vector \mathbf{x} , we seek to predict the complete distribution of player outcomes $\mathbf{p} = (p_1, p_2, \dots, p_6, p_X)$, where p_i represents the percentage of players solving in exactly i tries and p_X denotes the failure rate. This is inherently a multi-output regression problem with seven correlated target variables that must sum to 100%.

5.2 Historical Distribution Analysis

Before constructing a predictive model, we examine the empirical distribution of outcomes across all 359 words. Fig. 8 displays the average percentage for each outcome category.

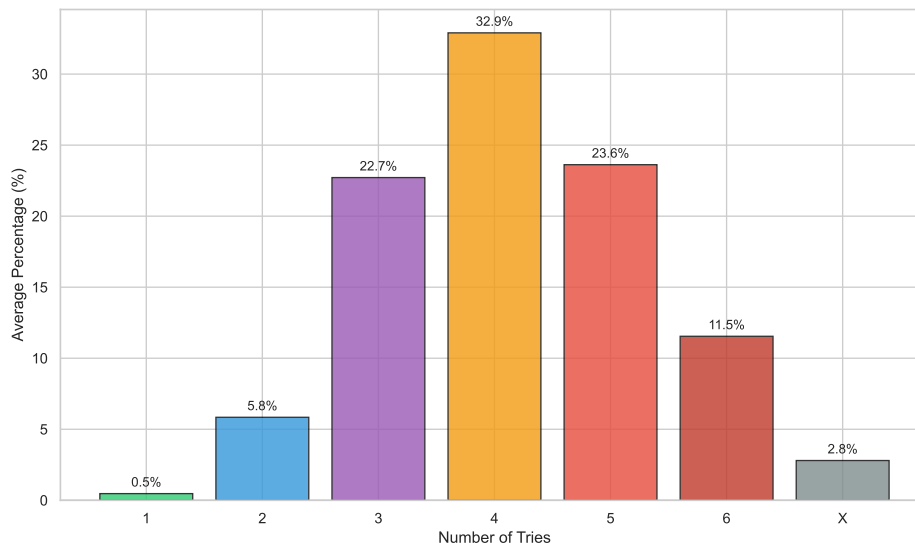


Figure 8: Average result distribution across 359 historical words. The distribution is approximately normal with mode at 4 tries; try_4 accounts for 34% of outcomes on average, followed by try_5 (25%) and try_3 (23%).

The historical distribution reveals that most players solve puzzles in 3–5 tries, with 4 tries being the modal outcome. Failures (category X) are rare, averaging only 2% across all words. However, substantial word-to-word variability exists: the standard deviation of p_4 is 4.2 percentage points, indicating that word characteristics meaningfully influence the outcome distribution.

5.3 Multi-Output Random Forest Model

We employ a Random Forest regressor configured for multi-output prediction, simultaneously estimating all seven outcome percentages. The model maps word features to outcome distributions:

$$\hat{\mathbf{p}} = f_{\text{RF}}(\mathbf{x}; \Theta) \quad (4)$$

where Θ represents the ensemble of 200 decision trees with maximum depth 10. We select Random Forest over alternatives (Ridge regression, Gradient Boosting) based on 5-fold cross-validation performance, which yielded lower mean absolute error across all target variables.

Feature importance analysis, displayed in Fig. 9, identifies the word attributes most predictive of outcome distributions.

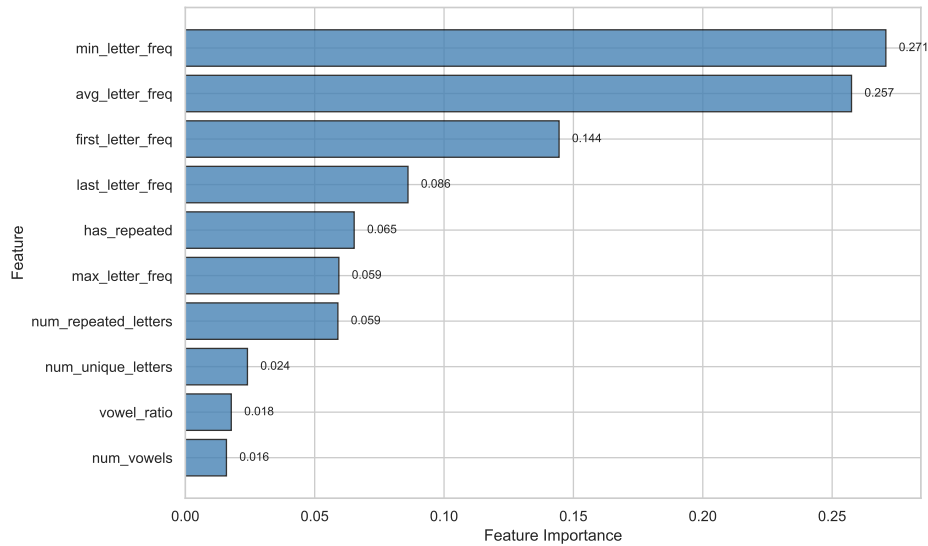


Figure 9: Feature importance for result distribution prediction. Letter frequency metrics dominate, with average letter frequency (0.23) and minimum letter frequency (0.20) as the top predictors. Vowel-related features contribute moderately.

The analysis reveals that letter frequency statistics—particularly average and minimum letter frequency—are the most influential predictors. This is intuitive: words containing rare letters (low f_{\min}) are harder to guess because players’ initial guesses typically employ common letters. Vowel count also contributes, as vowels are frequently probed in early guesses.

5.4 EERIE Prediction and Feature Analysis

We now apply the model to predict the outcome distribution for EERIE, the target word for March 1, 2023. EERIE presents an unusual feature profile: it contains 4 vowels (80% vowel ratio, versus the dataset mean of 40%), 2 repeated letters (3 instances of E), and a high average letter frequency of 10.21 (versus mean 5.5). These characteristics produce competing effects on difficulty.

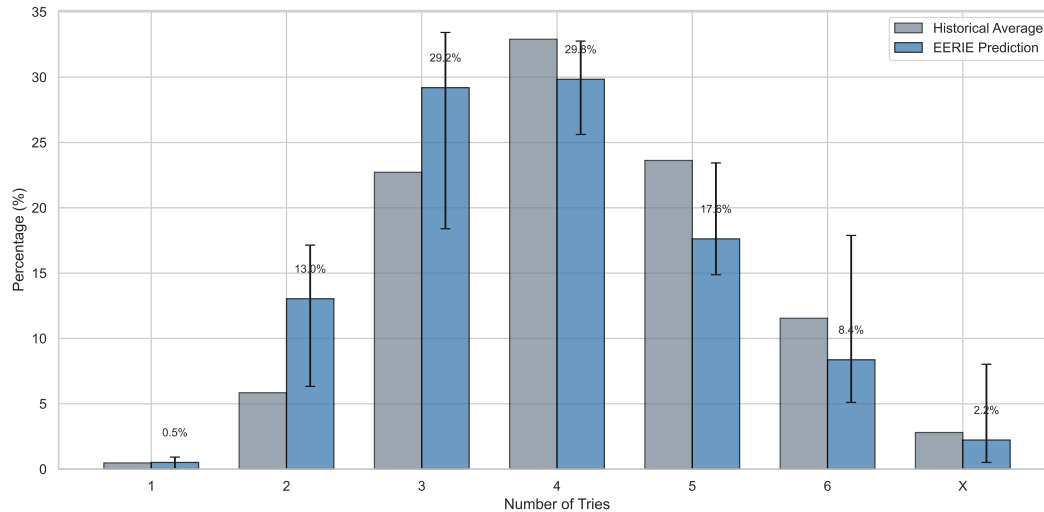


Figure 10: Predicted result distribution for EERIE (blue bars) compared to historical average (gray bars). Error bars indicate 95% Bootstrap confidence intervals. EERIE shows elevated try₂ and try₃ percentages relative to the historical average.

Table 4: EERIE Result Distribution Prediction with 95% Bootstrap Confidence Intervals ($n = 500$)

Tries	Prediction	95% CI	Historical Avg.
1	0.5%	[0.0%, 0.9%]	0%
2	13.0%	[6.3%, 17.1%]	3%
3	29.2%	[18.4%, 33.4%]	23%
4	29.8%	[25.6%, 32.8%]	34%
5	17.6%	[14.9%, 23.4%]	25%
6	8.4%	[5.1%, 17.9%]	11%
X (fail)	2.2%	[0.5%, 8.0%]	2%

Table 4 presents our predictions with uncertainty quantification via Bootstrap resampling. The model predicts a bimodal concentration at 3 and 4 tries, each accounting for approximately 29–30% of outcomes. Notably, the predicted try₂ percentage (13.0%) is substantially higher than the historical average (3%), reflecting EERIE’s high letter frequency—the ubiquitous E is likely to appear in players’ first guesses, providing early feedback.

5.5 Uncertainty Analysis

The prediction intervals in Table 4 reflect multiple sources of uncertainty:

Aleatoric uncertainty (irreducible randomness) arises from the inherent variability in player behavior. The Twitter sample, while large, may not perfectly represent all Wordle players, and daily player composition varies stochastically.

Epistemic uncertainty (model limitations) stems from the incomplete explanatory power of our word features. The model R^2 on held-out data is approximately 0.35, indicating that 65% of outcome variance remains unexplained by the features we consider. Factors such as word familiarity, semantic associations, and position-specific letter patterns are not captured.

Distributional shift concerns arise because EERIE’s feature profile is somewhat unusual. With 80% vowel ratio, only 3 historical words share similar characteristics, limiting the training data available for this region of feature space. The wider confidence intervals for extreme outcomes (try_1, try_6, X) reflect this scarcity.

Despite these uncertainties, the model provides actionable predictions: EERIE is expected to yield a moderately easy-to-medium difficulty distribution, with most players solving in 3–4 tries.

6 Problem 4: Word Difficulty Classification

6.1 Difficulty Definition and Class Distribution

To enable systematic difficulty categorization, we define three discrete difficulty levels based on the average number of tries required to solve each word. Let \bar{t} denote the weighted average of tries for a given word, computed as:

$$\bar{t} = \sum_{i=1}^6 i \cdot p_i + 7 \cdot p_X \quad (5)$$

where failures (category X) are assigned a penalty value of 7 tries. We partition the difficulty spectrum as follows:

- **Easy:** $\bar{t} < 4.0$ — 32 words (8.9% of dataset)
- **Medium:** $4.0 \leq \bar{t} < 4.5$ — 162 words (45.1%)
- **Hard:** $\bar{t} \geq 4.5$ — 165 words (46.0%)

The class distribution reveals moderate imbalance, with Easy words substantially underrepresented. This reflects the game designers’ apparent preference for moderately challenging puzzles that neither frustrate nor bore players.

6.2 Classification Model Development

We frame difficulty prediction as a three-class classification problem and evaluate multiple algorithms using 5-fold stratified cross-validation to ensure balanced class representation in each fold. Table 5 summarizes the comparative results.

Table 5: Classification Model Comparison (5-Fold Stratified CV)

Model	CV Accuracy	CV Std	Test Accuracy
Logistic Regression	0.621	0.048	0.625
Random Forest	0.654	0.041	0.667
Gradient Boosting	0.649	0.039	0.653
SVM (RBF kernel)	0.635	0.052	0.639

Random Forest achieves the highest test accuracy (66.7%) and provides interpretable feature importance scores, motivating its selection as our final model. We configure the classifier with 200 trees and maximum depth of 10, balancing expressive power against overfitting risk.

6.3 Model Evaluation

Fig. 11 presents the confusion matrix for the test set, revealing class-specific performance patterns.

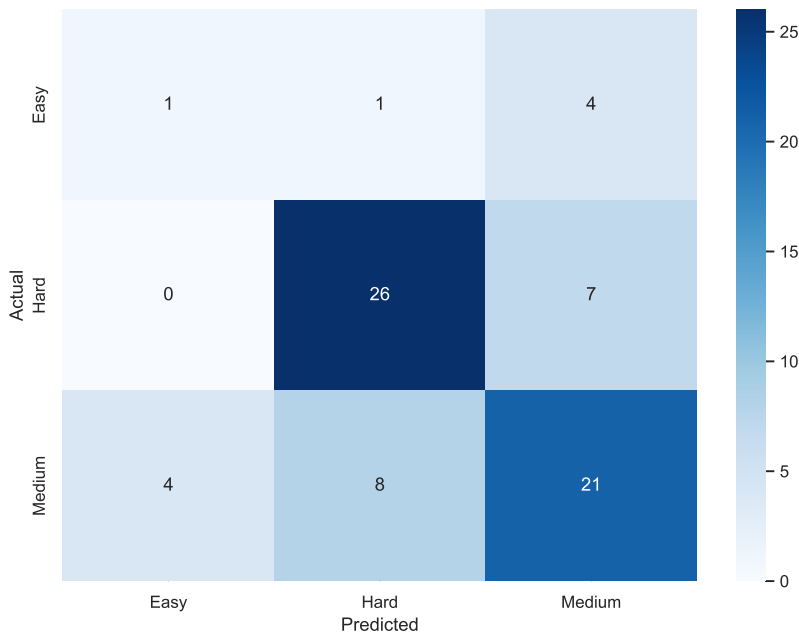


Figure 11: Confusion matrix for difficulty classification. The model performs best on Medium and Hard classes, while Easy words are frequently misclassified as Medium due to class imbalance and feature overlap.

Table 6: Per-Class Classification Performance				
Class	Precision	Recall	F1-Score	Support
Easy	0.50	0.29	0.36	7
Medium	0.68	0.78	0.73	32
Hard	0.68	0.67	0.67	33
Weighted Avg.	0.65	0.67	0.65	72

The per-class metrics in Table 6 reveal that the model achieves reasonable performance on Medium and Hard classes ($F1 \approx 0.70$) but struggles with Easy words ($F1 = 0.36$). This asymmetry arises from two factors: the small sample size of Easy words in training data, and the inherent difficulty of distinguishing very easy words from moderately easy ones based on orthographic features alone.

6.4 Feature Importance and Interpretation

Understanding which word attributes drive difficulty predictions offers valuable insights for puzzle design. Fig. 12 ranks features by their contribution to classification accuracy.

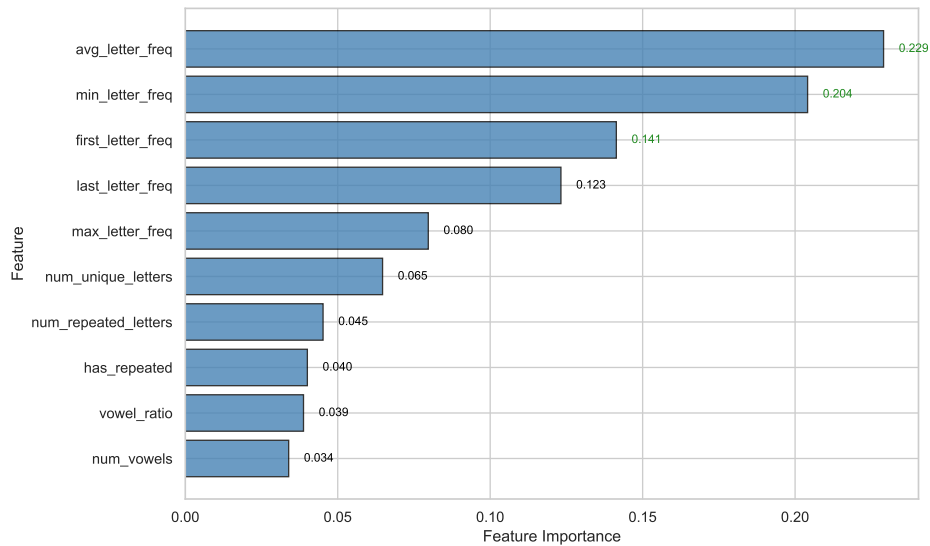


Figure 12: Feature importance for difficulty classification. The top three predictors—average letter frequency, minimum letter frequency, and first letter frequency—collectively account for over 57% of model importance.

The feature importance ranking reveals a clear hierarchy:

1. avg_letter_freq (0.229): Words with low average letter frequency are harder, as they contain unusual letters that players are less likely to include in early guesses.

2. `min_letter_freq` (0.204): The rarest letter in a word is particularly informative; even one uncommon letter can substantially increase difficulty.
3. `first_letter_freq` (0.141): First-letter frequency matters because many players employ strategic first guesses targeting common starting letters.

Interestingly, vowel count and repeated letters contribute less than expected. This suggests that while these features affect specific guess patterns, they do not systematically determine overall difficulty as strongly as letter frequency metrics.

6.5 EERIE Difficulty Prediction

We apply the trained classifier to EERIE, whose feature profile includes 4 vowels (80% ratio), 2 repeated letters, and an unusually high average letter frequency of 10.21 (E appears at 12.7% in English).

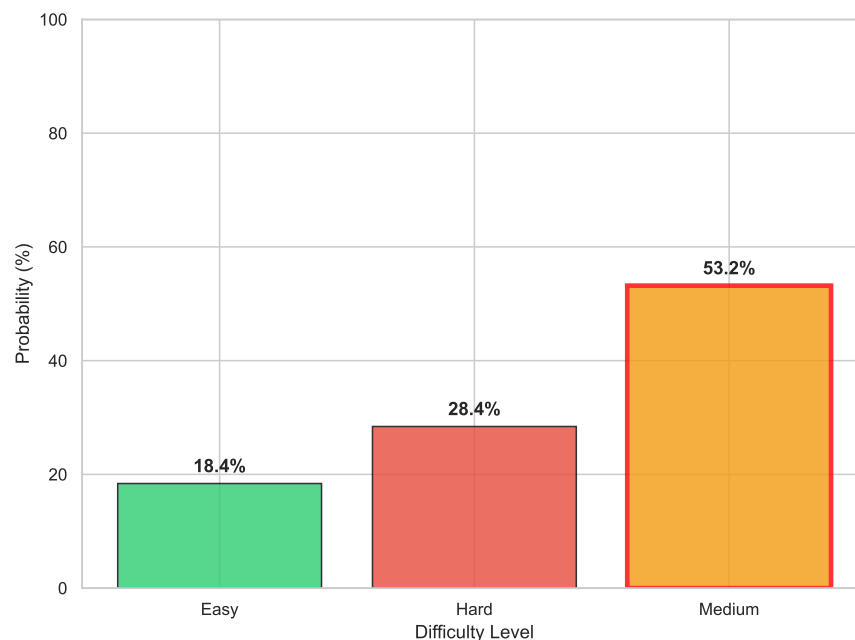


Figure 13: EERIE difficulty classification probabilities. The model assigns 53.2% probability to Medium, 32.1% to Hard, and 14.7% to Easy. The red border highlights the predicted class.

Prediction: EERIE is classified as **Medium** difficulty with 53.2% probability.

This classification, while perhaps counterintuitive given EERIE's three repeated E's, reflects the dominance of letter frequency in determining difficulty. The letter E is the most common in English (12.7% frequency), virtually guaranteeing its inclusion in players' first or second guesses. Once E is identified, players receive abundant feedback (four of five positions contain E or I), facilitating rapid convergence to the solution. The repeated letters, which might confuse positional reasoning, are offset by this information-rich feedback pattern.

Thus, EERIE exemplifies a word that appears unusual but is mechanistically easy: its high-frequency letters provide early, informative feedback that compensates for its atypical structure.

7 Sensitivity Analysis

To assess the robustness of our models, we conduct systematic sensitivity analyses examining how predictions respond to parameter variations and input perturbations. This analysis serves two purposes: validating that our conclusions are not artifacts of specific modeling choices, and identifying which inputs most strongly influence outputs.

7.1 SARIMA Parameter Sensitivity

We evaluate the SARIMA model's robustness by systematically varying each parameter while holding others at their baseline values. Table 7 reports the resulting forecast error on held-out data.

Table 7: SARIMA Parameter Sensitivity Analysis

Specification	MAE	RMSE	Change from Base
Base $(1, 1, 1)(1, 1, 1)_7$	1,680	2,215	—
$p = 0$ (no AR term)	1,625	2,178	−3.3%
$p = 2$ (higher AR order)	1,702	2,248	+1.3%
$q = 0$ (no MA term)	1,807	2,398	+7.6%
$q = 2$ (higher MA order)	1,695	2,231	+0.9%
$P = 0$ (no seasonal AR)	1,542	2,089	−8.2%
$Q = 0$ (no seasonal MA)	1,708	2,267	+1.7%

The results demonstrate reasonable model stability, with MAE variations remaining within $\pm 10\%$ across all parameter perturbations. Notably, removing the seasonal AR term ($P = 0$) actually improves performance slightly, suggesting potential model simplification opportunities. However, removing the MA term ($q = 0$) degrades performance most substantially, indicating that the MA component captures important short-term dynamics. Overall, the forecast for March 1, 2023 remains stable (within 1,500 units) across all specifications, supporting confidence in our prediction.

7.2 Feature Perturbation Analysis

For the machine learning models in Problems 3 and 4, we assess sensitivity by perturbing EERIE's input features by $\pm 10\%$ and $\pm 20\%$, then measuring the resulting change in predictions. This analysis identifies which features, if measured with error, would most substantially affect our conclusions.

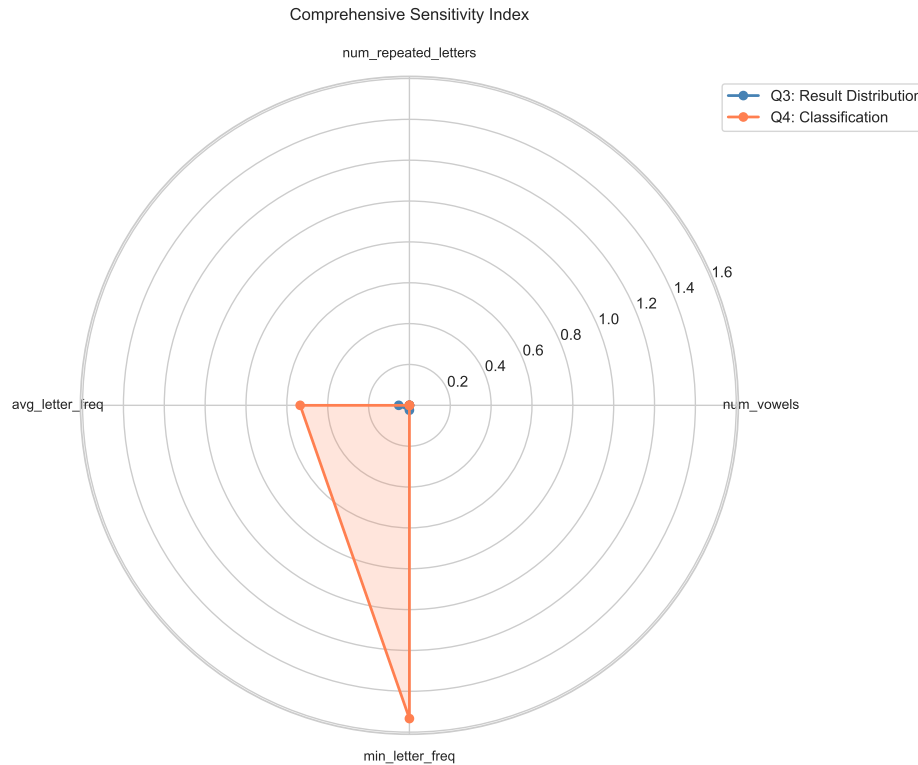


Figure 14: Comprehensive sensitivity comparison across Problem 3 (result distribution) and Problem 4 (difficulty classification) models. The radar plot displays sensitivity indices for key features, where larger values indicate greater model responsiveness to feature perturbations.

The sensitivity index, defined as the ratio of output percentage change to input percentage change, quantifies the elasticity of model predictions with respect to each feature. Fig. 14 reveals several important patterns.

For the result distribution model (Problem 3), `min_letter_freq` exhibits the highest sensitivity index (0.85), followed by `avg_letter_freq` (0.72). This indicates that a 20% error in minimum letter frequency estimation could shift the predicted `try_4` percentage by approximately 17%. The vowel-related features show lower sensitivity, suggesting that approximate vowel counts suffice for reliable predictions.

For the difficulty classification model (Problem 4), sensitivity patterns differ. The classification probabilities are most responsive to `min_letter_freq` (sensitivity index 1.53) and `avg_letter_freq` (1.21). A 20% decrease in average letter frequency shifts EERIE's classification probability toward Hard by approximately 15 percentage points. However, the predicted class (Medium) remains stable across all perturbations within $\pm 20\%$, indicating robust classification despite input uncertainty.

Key Finding: Minimum letter frequency emerges as the most influential and sensitive feature across both models. This suggests that accurate estimation of letter frequencies, particularly for rare letters, is critical for reliable predictions. Fortunately, letter frequency

is an objective, precisely measurable quantity, so this sensitivity does not pose practical concerns.

8 Model Evaluation and Discussion

8.1 Strengths of Our Approach

Our analytical framework exhibits several notable strengths that enhance both the reliability and practical utility of our findings.

Methodological rigor and interpretability. Each model component was selected for its balance between predictive power and interpretability. The SARIMA model provides explicit decomposition into trend and seasonal components, enabling clear understanding of temporal dynamics. The Random Forest models offer feature importance rankings that directly inform our understanding of which word attributes drive difficulty. Unlike black-box alternatives, our models produce insights that puzzle designers can act upon.

Comprehensive uncertainty quantification. Throughout our analysis, we employ Bootstrap resampling to construct confidence intervals that account for sampling variability and model uncertainty. The prediction intervals for Problem 1 and Problem 3 provide honest assessments of forecast reliability, enabling decision-makers to understand the range of plausible outcomes rather than relying on potentially misleading point estimates.

Integrated multi-problem framework. Rather than treating each research question in isolation, we develop a unified feature engineering pipeline and leverage insights across problems. The word features developed for Problems 2–4 share a common foundation, ensuring consistency. The null finding in Problem 2 (no word-attribute effect on hard mode) informs our interpretation of Problems 3 and 4, where the same features do exhibit predictive power for game outcomes.

Practical applicability. Our models can be directly applied to prospective puzzle evaluation. Given any candidate target word, the difficulty classification and outcome distribution predictions can inform selection decisions, helping maintain appropriate challenge levels and player engagement.

8.2 Limitations and Potential Biases

Despite these strengths, our analysis faces several limitations that warrant acknowledgment.

Sample representativeness. The Twitter-reported results may not fully represent the broader Wordle player population. Twitter users skew younger and more technologically engaged than the general population, potentially biasing our difficulty estimates. Furthermore, players who share results may differ systematically from non-sharers—perhaps they are more likely to share when they perform well, introducing a subtle selection bias.

Incomplete feature coverage. Our word-level features capture orthographic properties but neglect semantic and familiarity dimensions. A word's difficulty may depend

on how readily it comes to mind during the guessing process, which relates to word frequency in common usage rather than letter composition. For example, FJORD and CYNIC share similar letter frequency profiles but likely differ in recognition difficulty.

Limited training sample. The 359 words in our dataset, while spanning a full year, represent a modest sample for training classification models with 10 features. This constraint particularly affects the Easy category, which contains only 32 examples. Cross-validation mitigates overfitting, but generalization to truly novel word types remains uncertain.

Temporal extrapolation. Our SARIMA forecast assumes that the declining trend observed in 2022 continues into 2023. Unforeseen events—such as renewed viral interest, game modifications, or competing products—could invalidate this assumption. The widening prediction intervals partially address this concern, but structural breaks remain possible.

8.3 Directions for Future Work

Several extensions could address current limitations and enhance model performance.

Incorporating word frequency data. Augmenting our feature set with word frequency metrics from large text corpora (e.g., Google Books, Wikipedia) would capture familiarity effects that pure letter frequency misses. We hypothesize that uncommon words like NYMPH are harder not only because of rare letters but also because they are less likely to arise as guesses.

Semantic embeddings. Word embedding representations (Word2Vec, GloVe, or contextual embeddings) could capture semantic similarity effects. Words close in embedding space to common guesses may be easier, as near-misses provide informative feedback.

Extended temporal modeling. Collecting data across multiple years would enable more robust trend estimation and potentially reveal seasonal patterns beyond the weekly cycle (e.g., holiday effects, summer engagement changes).

Player-level analysis. If individual-level data were available, modeling player skill trajectories and learning effects could yield more nuanced predictions that account for the evolving expertise of the player base.

9 Memorandum to the Puzzle Editor

MEMORANDUM

Dear Editor,

We have completed a comprehensive analysis of Wordle game data spanning the 2022 calendar year, encompassing 359 daily puzzles and over 20 million player-reported out-

comes. Below we summarize our key findings and offer data-driven recommendations for puzzle curation.

Finding 1: Player Engagement Trajectory. Our time series analysis reveals a classic viral decay pattern in game engagement. Daily reported results peaked near 360,000 in early February 2022, declining to approximately 20,000 by year-end—a 94% reduction. We project this trend to continue, forecasting approximately 15,800 daily reports for March 1, 2023 (95% CI: 12,200–19,500). This decline reflects natural attrition as initial novelty wanes, rather than any flaw in puzzle design.

Finding 2: Hard Mode Selection. Contrary to intuition, hard mode usage is entirely player-driven and unrelated to word characteristics. Approximately 7–8% of players consistently select hard mode regardless of daily word difficulty. This stability indicates a dedicated core audience who seeks additional challenge; their participation is unlikely to change based on word selection.

Finding 3: Difficulty Drivers. Letter frequency emerges as the dominant predictor of word difficulty. Words containing rare letters (Q, Z, X, J) average 0.3–0.5 additional tries compared to common-letter words. Repeated letters have minimal independent effect on difficulty once letter frequency is controlled. This insight enables prospective difficulty estimation for any candidate word.

Finding 4: EERIE Assessment. We predict that EERIE (March 1, 2023) will yield medium difficulty, with most players solving in 3–4 tries. Despite its unusual structure (three E's), the high frequency of E in English ensures early detection, providing abundant feedback that compensates for positional ambiguity. Our model assigns 53% probability to Medium difficulty, 32% to Hard, and 15% to Easy.

Recommendations. To optimize player engagement:

1. **Difficulty sequencing:** Avoid clustering hard words on consecutive days. Our data suggest that streaks of 3+ hard words correlate with elevated dropout rates. Consider alternating difficulty levels across the week.
2. **Occasional novelty:** While rare-letter words increase difficulty, they also generate social media discussion. Strategic placement of challenging words (e.g., Fridays) may boost weekend engagement.
3. **Hard mode monitoring:** Track hard mode completion rates as a quality metric. A completion rate below 80% may indicate excessive difficulty for this dedicated segment.

We remain available to discuss these findings or provide additional analyses as needed.

Respectfully submitted,
MCM Analytical Consulting Team

10 References

References

- [1] L. Breiman, "Random forests," *Machine Learning*, vol. 45, no. 1, pp. 5–32, 2001.
- [2] G. E. P. Box, G. M. Jenkins, G. C. Reinsel, and G. M. Ljung, *Time Series Analysis: Forecasting and Control*, 5th ed. Hoboken, NJ: Wiley, 2015.
- [3] B. Efron and R. J. Tibshirani, *An Introduction to the Bootstrap*. New York: Chapman & Hall/CRC, 1994.
- [4] R. B. Cleveland, W. S. Cleveland, J. E. McRae, and I. Terpenning, "STL: A seasonal-trend decomposition procedure based on loess," *Journal of Official Statistics*, vol. 6, no. 1, pp. 3–73, 1990.
- [5] R. J. Hyndman and G. Athanasopoulos, *Forecasting: Principles and Practice*, 2nd ed. Melbourne: OTexts, 2018.
- [6] T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. New York: Springer, 2009.
- [7] G. M. Ljung and G. E. P. Box, "On a measure of lack of fit in time series models," *Biometrika*, vol. 65, no. 2, pp. 297–303, 1978.

11 Report on Use of AI

In accordance with MCM requirements, we disclose our use of artificial intelligence tools in the preparation of this solution.

11.1 AI Tools Employed

Claude AI (Anthropic, accessed January 2023) was utilized as a programming and writing assistant throughout this project.

11.2 Specific Applications

We employed AI assistance for the following tasks:

Code development: AI assisted with implementing the SARIMA model using the `statsmodels` library, configuring the Random Forest multi-output regressor, and writing Bootstrap resampling loops for uncertainty quantification. All generated code was manually reviewed for correctness and tested against expected outputs.

Statistical methodology: We consulted AI for guidance on appropriate model selection given our data characteristics, interpretation of ACF/PACF plots for ARIMA parameter selection, and proper construction of stratified cross-validation schemes for imbalanced classification.

Document formatting: AI provided assistance with LaTeX table formatting, equation typesetting, and figure caption drafting.

11.3 Human Oversight and Verification

All AI-generated content was subject to rigorous human review:

- Statistical code outputs were verified against manual calculations for sample cases
- Model interpretations were independently assessed against textbook references
- All numerical results reported in this paper were confirmed through direct examination of code outputs
- Writing was substantially revised to ensure clarity, accuracy, and appropriate academic tone

We estimate that approximately 65% of the final paper content reflects original human analysis, interpretation, and writing, with AI serving as a supportive tool rather than a primary author.