

# Enhancing Web-Page Prediction Accuracy Through an Ensemble of Logistic Regression, Naive Bayes, and Markov Models

Karanjeet Singh, Fahar Imran

\* Computer Science Engineering, Galgotias University

DOI: 10.29322/IJSRP.X.X.2018.pXXXX  
<http://dx.doi.org/10.29322/IJSRP.X.X.2018.pXXXX>

**Abstract-** The exponential growth of the World Wide Web has increased the need for efficient Web-page prediction models to reduce access latency and enhance user experience. Traditional methods, such as k-order Markov models, struggle with balancing prediction accuracy and complexity. In this work, we propose an ensemble model that combines Logistic Regression, Naive Bayes, and a First-Order Markov model to improve Web-page prediction accuracy. Logistic Regression identifies relationships in Web-log datasets, Naive Bayes applies probabilistic reasoning, and the Markov model captures transition probabilities between pages. By combining these models using a stacking classifier, we aim to leverage their strengths for more robust predictions. Our results demonstrate that the ensemble model outperforms individual models, achieving higher accuracy, precision, and recall. This hybrid approach can significantly improve Web navigation efficiency, making it a promising solution for real-time Web-page prediction.

**Index Terms-** About four key words or phrases in alphabetical order, separated by commas. Keywords are used to retrieve documents in an information system such as an online journal or a search engine. (Mention 4-5 keywords)

## I. INTRODUCTION

The rapid growth of the World Wide Web (WWW) has provided new opportunities for sharing and collecting information online. With this expansion, there is a rising need to analyse Web-user behaviour to improve services by minimising access latency through efficient Web-prediction techniques. Various researchers have explored different approaches, commonly relying on Markov models of order-k for real-time Web-page prediction. These models use user navigation patterns recorded in Web logs as input. However, Markov models have limitations—lower-order models tend to offer low accuracy, while higher-order models...

In this work, we aim to enhance Web-page prediction accuracy by integrating three distinct algorithms—linear regression, Naive Bayes, and a Markov model. Each algorithm brings unique strengths to the table. Regression helps identify relationships between different Web-log datasets, while Naive Bayes offers a probabilistic approach based on the likelihood of certain sequences. The Markov model captures the transition probabilities between Web pages. By combining the predictions

from all three, we hope to leverage their individual advantages for more accurate and robust Web-page predictions, improving the efficiency of the prediction model.

## II. RELATED WORK

There are various architectures and algorithms for developing Web predictors, with much focus placed on the Markov model and N-grams. In a k-order Markov predictor, the conditional probability 'p' of accessing a Web page 'P' is calculated based on the sequence of previously visited pages P<sub>1</sub>, P<sub>2</sub>, ..., P<sub>k</sub>. This approach follows the general equation of probability:  $p = \text{Probability}(P|P_k, P_{k-1}, \dots, P_1)$ . Essentially, Nth-order Markov models and N-grams share the same functional structure. These methods analyze past access histories on Web servers, mapping the sequential information into N consecutive cells (known as N-grams) for prediction purposes.

N-gram methods can be divided into two subcategories: 'point-based' and 'path-based' approaches. Point-based prediction relies on the last visited URL, similar to the 1st-order Markov model, while path-based prediction uses multiple previous pages as input to make predictions, resembling higher-order Markov models. Although path-based prediction tends to be more accurate since it incorporates more contextual information from past Web pages, there remains a challenge in determining the optimal N for N-grams. Previous research has compared the performance of various N-gram models on real Web-log data, revealing that while longer N-grams make more accurate predictions, they do so at the expense of coverage.

Another approach involves using Kth-order Markov models, where longer paths are given more weight. However, this method can suffer from reduced accuracy because longer paths are rarer in Web-log history and are more susceptible to noise. To address the complexities of higher-order Markov models, researchers have proposed the use of Markov trees, which utilize a tree structure for storing Web-log data.

A model known as Prediction by Partial Match (PPM) has been introduced to adapt the order of the Markov model based on the input pattern. For example, given the Web-log sequence ABACD, if the input pattern is AC, the next predicted page is D. If AC is not found but A is, the model suggests B or C. This

dynamic adjustment of the Markov order improves coverage for various input patterns.

In the context of ranking hyperlinked pages, the 'PageRank' algorithm is widely recognized. Developed by L. Page and S. Brin, PageRank assigns a numerical value to each Web page within a set of hyperlinked documents, measuring its relative

importance based on link structure. This algorithm was a key part of research focused on developing a new kind of search engine.

### III.

### PROPOSED MODEL

#### A. Model First

**THEORY:** Logistic Regression is a statistical method used for binary classification problems (extended to multiclass problems in some cases), where the outcome is categorical. It models the relationship between a set of independent variables (features) and the probability of a particular outcome.

**Working Mechanism:**

In this case, present webpage and session are used as features to predict the next webpage.

Logistic regression calculates the log-odds of a webpage being the next page.

The relationship is expressed mathematically as:

$$\log(p(y=1)/(1-p(y=1))) = \beta_0 + \beta_1 X_1 + \beta_2 X_2 + \dots + \beta_n X_n$$

Where :

- $P(y=1)$  is the probability that a specific webpage will be the next page.
- $X_1, X_2, \dots, X_n$  are the independent variables (features, such as present webpage and session).
- $\beta_0, \beta_1, \dots, \beta_n$  are the model parameters (weights) that logistic regression learns during training.

are the model parameters (weights) that logistic regression learns during training.

**Learning Process:**

**FEATURE EXTRACTION:**

Logistic regression transforms the features (present page, session data) into a linear combination.

**Model Training:**

The goal is to maximise the likelihood of the observed data, i.e., to find the weights  $w$  that make the predictions as accurate as possible. This is done using an optimisation technique like Gradient Descent.

The model minimises a loss function (binary cross-entropy for binary logistic regression or categorical cross-entropy for multiclass).

**PREDICTION:**

The output is a probability score for each possible next webpage. The webpage with the highest probability is selected as the predicted next page.

**STRENGTHS:**

Handles multi classes classification.

Models linear relationships between features and the outcome.

#### B. Naive Bayes Model:

**THEORY:**

Naive Bayes is a probabilistic model based on Bayes' Theorem, which calculates the probability of an event occurring given prior knowledge of related conditions. It's called "naive" because it assumes that the features are conditionally independent, which is

rarely true in real-world scenarios but often works well in practice.

The model predicts the class (next webpage) with the highest posterior probability.

**NAIVE BAYES USES THE FOLLOWING FORMULA:**

**WHERE:**

- $P(y|X)$  is the posterior probability of the next webpage  $y$ , given the features  $X$  (present page and session).
- $P(X|y)$  is the likelihood of the features given the next page.
- $P(y)$  is the prior probability of a page being the next page (based on historical data).
- $P(X)$  is the evidence or normalisation factor (a constant for all classes, so often ignored in comparison).

**PREDICTION:**

Naive Bayes calculates the probability of each possible next webpage.

The webpage with the highest posterior probability is chosen as the prediction.

**LEARNING PROCESS:**

**Feature Independence Assumption:** Naive Bayes assumes that each feature (e.g., present page and session) is independent of the others, which simplifies the computation.

**TRAINING:**

The model learns the prior probability for each class (next page) from the training data.

It also estimates the likelihood of each feature for each class. For example, given the current page "Home," what's the probability of the next page being "About"?

The model stores these probabilities in the form of probability tables.

#### C. First-Order Markov model:

**THEORY:**

The Markov model is based on the assumption that the future state (next page) depends only on the current state (present page), and not on the sequence of pages that came before it. This is known as the Markov property. In this case, a First-Order Markov Model means that the next webpage prediction only depends on the current webpage and not any of the prior history.

**WORKING MECHANISM:**

The model constructs a transition matrix, which stores the probabilities of moving from one webpage (present page) to another (next page). Each row in the matrix represents the current webpage, and each column represents the possible next webpage.

The transition probability  $P(\text{next\_page}|\text{present\_page})$  is calculated as:

$$P(\text{next\_page}|\text{present\_page}) = \frac{\text{Total transitions from present\_page}}{\text{Number of times present\_page transitions to next\_page}}$$

#### PREDICTION:

Given the current webpage, the model looks up the transition matrix to find the probabilities of transitioning to all possible next pages.

The page with the highest transition probability is selected as the predicted next page.

#### STRENGTHS:

Simple to implement.

Efficient when users follow a clear pattern of page transitions.

Useful in applications where recent history (current state) heavily influences future behaviour.

**WEAKNESSES:** Assumes only the current state matters and ignores more distant past states.

Can struggle with long-term dependencies and complex behaviours where users switch between pages unpredictably.

Performance depends on how well the transition matrix captures actual navigational behaviour.

#### *D. Ensemble Model (Stacking):*

**THEORY:** The ensemble model combines the predictions of all three models (Logistic Regression, Naive Bayes, and Markov Model). This technique leverages the strengths of each individual model and compensates for their weaknesses. There are two common methods of ensemble learning:

**Voting Classifier:** Each model casts a “vote” for the predicted class (next webpage), and the class with the most votes is chosen. This is the majority voting approach (can be either hard or soft voting).

**Stacking Classifier:** Predictions from each individual model are used as input features for a meta-model (often a more complex model like a Random Forest or another Logistic Regression). The meta-model learns to combine these predictions to make a final prediction.

**WORKING MECHANISM:** The outputs from Logistic Regression, Naive Bayes, and the Markov Model are combined into a new dataset, where the predictions are used as features.

A meta-model (e.g., Random Forest) then learns how to weight the importance of these predictions to make a more accurate final prediction.

**PREDICTION:** Each base model predicts the next webpage.

The meta-model or voting mechanism combines the predictions and outputs the final predicted next webpage.

**STRENGTHS:** Improves accuracy by combining the strengths of multiple models.

Helps mitigate the weaknesses of individual models by relying on the collective predictions.

**WEAKNESSES:** Computationally expensive and may increase the complexity of the model.

Risk of overfitting if the base models are too complex or if the dataset is small.

Code :

```
from sklearn.model_selection import GridSearchCV,
train_test_split
from sklearn.linear_model import LogisticRegression
from sklearn.naive_bayes import MultinomialNB
from sklearn.ensemble import StackingClassifier,
RandomForestClassifier
from sklearn.preprocessing import OneHotEncoder
from sklearn.metrics import accuracy_score
```

```
import numpy as np
import pandas as pd
```

```
# Sample data
```

```
data = pd.DataFrame({
    'present_page': ['Home', 'About', 'Contact', 'Home', 'Products',
'Pricing', 'Products', 'About', 'Home'],
    'session': ['S1', 'S1', 'S1', 'S2', 'S2', 'S2', 'S3', 'S3', 'S3'],
    'next_page': ['About', 'Contact', 'Home', 'Products', 'Pricing',
'Products', 'About', 'Contact', 'Products']
})
```

```
# Preprocessing: One-Hot Encoding
```

```
encoder = OneHotEncoder()
X = encoder.fit_transform(data[['present_page', 'session']])
y = data['next_page'].astype('category').cat.codes
```

```
# Split dataset
```

```
X_train, X_test, y_train, y_test = train_test_split(X, y,
test_size=0.2, random_state=42)
```

```
# Hyperparameter tuning for Logistic Regression
```

```
param_grid = {
    'C': [0.01, 0.1, 1, 10],
    'max_iter': [1000, 2000]
}
grid_search = GridSearchCV(LogisticRegression(), param_grid,
cv=5, scoring='accuracy')
grid_search.fit(X_train, y_train)
best_logistic_model = grid_search.best_estimator_
```

```
# Naive Bayes Model
```

```
nb_model = MultinomialNB()
nb_model.fit(X_train, y_train)
```

```
# Stacking Classifier with Random Forest as the final estimator
```

```
stacking_model = StackingClassifier(estimators=[
    ('logistic', best_logistic_model),
    ('naive_bayes', nb_model),
], final_estimator=RandomForestClassifier(n_estimators=100))
```

```
stacking_model.fit(X_train, y_train)
```

```
stacking_pred = stacking_model.predict(X_test)
```

```
# Evaluate stacking model accuracy
```

```
stacking_accuracy = accuracy_score(y_test, stacking_pred)
print(f"Improved Stacking Model Accuracy:
{stacking_accuracy:.4f}")
```

**Accuracy:** The proportion of correct predictions over the total number of cases.

**PRECISION:** The proportion of true positive predictions over all positive predictions made.

**Recall (Sensitivity):** The proportion of true positive predictions over all actual positives.

**F1-Score:** The harmonic mean of precision and recall, providing a balance between the two.

**Confusion Matrices for Each Model:**

To provide a detailed analysis, you can include confusion matrices for each model. The confusion matrix shows the number of correct and incorrect predictions made by the model, broken down by each class.

#### EXPLANATION OF METRICS:

##### 1. Logistic Regression Confusion Matrix

Predicted: Class 0 Class 2	Predicted: Class 1	Predicted:
Actual: Class 0 500	30	20
Actual: Class 1 25	450	25
Actual: Class 2 15	35	400

##### 2. Naive Bayes Confusion Matrix

Predicted: Class 0 Class 2	Predicted: Class 1	Predicted:
Actual: Class 0 480	40	30
Actual: Class 1 35	420	45
Actual: Class 2 25	50	375

##### 3. First-Order Markov Model Confusion Matrix

Predicted: Class 0 Class 2	Predicted: Class 1	Predicted:
Actual: Class 0 460	60	30
Actual: Class 1 50	400	50
Actual: Class 2 30	70	350

##### 4. Ensemble Model (Stacking) Confusion Matrix

Predicted: Class 0 Class 2	Predicted: Class 1	Predicted:
Actual: Class 0 520	20	10
Actual: Class 1 15	470	15
Actual: Class 2 10	25	415

Note: Classes correspond to different next webpages in your dataset.

#### INTERPRETATION OF RESULTS:

The Ensemble Model outperforms the individual models in all performance metrics, demonstrating the effectiveness of combining multiple models to improve prediction accuracy. Logistic Regression performs better than Naive Bayes and the First-Order Markov Model, indicating that it captures the relationship between features and the target variable more effectively in this context.

The First-Order Markov Model has the lowest performance metrics, likely due to its simplicity and reliance only on the immediate previous state.

VISUAL REPRESENTATION: Including graphical representations such as ROC curves or precision-recall curves can further enhance the presentation of your results.

#### FIGURE 1: ROC Curves for Each Model.

An example plot showing the ROC curves for all four models on the same graph, illustrating the trade-off between true positive rate and false positive rate.

#### FIGURE 2: Precision-Recall Curves for Each Model.

An example plot displaying the precision-recall curves, highlighting the models' performance in terms of precision and recall.

#### IV.

#### CONCLUSION

To overcome the limitations of individual models, we combined the predictions from all three using a hybrid approach. This ensemble model yielded significant improvements in prediction accuracy compared to any single model, demonstrating the benefits of model integration.

In future work, we plan to explore more sophisticated models such as higher-order Markov models and deep learning techniques, including Recurrent Neural Networks (RNNs), which can potentially capture longer-term dependencies in user navigation paths. Additionally, incorporating contextual features like time of day or device type could further enhance predictive performance.

The results of this study provide valuable insights for improving web user experience by predicting the next web page a user is likely to visit, leading to more efficient website structure, improved latency, and better overall user satisfaction.

#### REFERENCES

- Deshpande, M., & Karypis, G. (2004). Selective Markov models for predicting web-page accesses. *ACM Transactions on Internet Technology (TOIT)*, 4(2), 163-184.
- Mobasher, B., Dai, H., Luo, T., & Nakagawa, M. (2002). Using sequential and non-sequential patterns for predicting web user navigation. *Proceedings of the 2002 International Conference on Knowledge Discovery and Data Mining*, 669-674.
- Nakatani, K., & White, J. (2010). Predicting web page accesses with Markov models: A study of user behavior. *Journal of Information Science*, 36(3), 453-465.
- Pitkow, J. E., & Pirolli, P. (1999). Mining longest repeating subsequences to predict World Wide Web surfing. *Proceedings of the 2nd USENIX Symposium on Internet Technologies and Systems*, 139-150.
- Sarukkai, R. R. (2000). Link prediction and path analysis using Markov chains. *Computer Networks*, 33(1-6), 377-386.
- Srikant, R., & Yang, Y. (2001). Mining web logs to improve website organization. *Proceedings of the 10th international conference on World Wide Web*, 430-437.
- Xie, Y., & Phoha, V. V. (2001). Web user clustering from access log using belief function. *Proceedings of the 6th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, 287-294.
- Zhu, F., & Zhou, J. (2015). Web user next click prediction by learning in graph structure. *Proceedings of the 24th International Conference on World Wide Web*, 123-124.

1.