

Model Evaluation and Cross-Testing Report

Evaluation Metrics Summary

1. Model A (Balanced Data) → Tested on Balanced Data

Logistic Regression:				
Accuracy: 0.4473150962512665				
Classification Report:				
	precision	recall	f1-score	support
1.0	0.50	0.57	0.53	389
2.0	0.37	0.32	0.34	397
3.0	0.35	0.31	0.33	395
4.0	0.40	0.42	0.41	397
5.0	0.58	0.62	0.60	396
accuracy			0.45	1974
macro avg	0.44	0.45	0.44	1974
weighted avg	0.44	0.45	0.44	1974

- Accuracy: **44.73%**
- Macro Avg Precision: 0.44
- Macro Avg Recall: 0.45
- F1-Score (Macro): 0.44
- Algorithm: Logistic regression
- Model size: 353 KB

Code:

```
import pickle

# Save the model

with open('balance.pkl', 'wb') as f:

    pickle.dump(lr, f)

print(" Model saved using pickle.")

with open("vecbalance.pkl", "wb") as f:

    pickle.dump(tv, f)
```

2. Model B (Imbalanced Data) → Tested on Imbalanced Data

```
🔧 Best Parameters Found: {'bootstrap': False, 'max_depth': None, 'min_samples_leaf': 1, 'min_samples_split': 5, 'n_estimators': 200}
✅ Accuracy: 0.6729426055268752
📊 Classification Report:
```

	precision	recall	f1-score	support
1.0	0.63	0.40	0.49	323
2.0	0.97	0.81	0.88	490
3.0	0.86	0.74	0.79	832
4.0	0.53	0.80	0.64	994
5.0	0.61	0.42	0.50	654
accuracy			0.67	3293
macro avg	0.72	0.64	0.66	3293
weighted avg	0.70	0.67	0.67	3293

- Accuracy: **67.29%**
- Macro Avg Precision: 0.72
- Macro Avg Recall: 0.64
- F1-Score (Macro): 0.66
- Algorithm: Random forest
- Model size: 156 MB

Code:

```
import pickle

# Save the model

with open('random_forest_model.pkl', 'wb') as f:
    pickle.dump(rf, f)

print(" Model saved using pickle.")

with open("vectorizer.pkl", "wb") as f:
    pickle.dump(tv, f)
```

Cross-Testing Results

3. Model A → Tested on Imbalanced Test Set

Accuracy: 0.429121041728322				
	precision	recall	f1-score	support
1.0	0.32	0.61	0.42	334
2.0	0.28	0.27	0.28	509
3.0	0.44	0.30	0.35	848
4.0	0.51	0.42	0.46	1012
5.0	0.51	0.64	0.57	676
accuracy			0.43	3379
macro avg	0.41	0.45	0.42	3379
weighted avg	0.44	0.43	0.42	3379

- Accuracy: **42.91%**
- Macro Avg Precision: 0.41
- Macro Avg Recall: 0.45
- F1-Score (Macro): 0.42
- Algorithm: Logistic Regression

Code:

```
import pickle

# Load the model

with open("balance.pkl", "rb") as f:
    loaded_model = pickle.load(f)

# Load the vectorizer

with open("vecbalance.pkl", "rb") as f:
    loaded_vectorizer = pickle.load(f)

# Transform new data using the loaded vectorizer
X_train_vec_new = loaded_vectorizer.transform(X_train)
X_test_vec_new = loaded_vectorizer.transform(X_test)

# Predict using the loaded model
y_pred = loaded_model.predict(X_test_vec_new)

print("Accuracy:", accuracy_score(y_test, y_pred))

print(classification_report(y_test, y_pred))
```

4. Model B → Tested on Balanced Test Set

Accuracy: 0.4199594731509625				
	precision	recall	f1-score	support
1.0	0.50	0.68	0.58	389
2.0	0.48	0.03	0.05	397
3.0	0.37	0.14	0.21	395
4.0	0.30	0.70	0.42	397
5.0	0.60	0.56	0.58	396
accuracy			0.42	1974
macro avg	0.45	0.42	0.37	1974
weighted avg	0.45	0.42	0.37	1974

- Accuracy: **41.99%**
- Macro Avg Precision: 0.45
- Macro Avg Recall: 0.42
- F1-Score (Macro): 0.37
- Algorithm: Random forest

Code:

```
import pickle

# Load the model
with open("random_forest_model.pkl", "rb") as f:
    loaded_model = pickle.load(f)

# Load the vectorizer
with open("vectorizer.pkl", "rb") as f:
    loaded_vectorizer = pickle.load(f)

# Transform new data using the loaded vectorizer
```

```
#X_train_vec_new = loaded_vectorizer.transform(X_train)
X_test_vec_new = loaded_vectorizer.transform(X_test)

# Predict using the loaded model
y_pred = loaded_model.predict(X_test_vec_new)
print("Accuracy:", accuracy_score(y_test, y_pred))
print(classification_report(y_test, y_pred))
```

TABLE

Scenario	Model	Algorithm	Training Dataset	Test Dataset	Accuracy
1. Trained and tested on balanced	A	Logistic Regression	Balanced	Balanced	0.4473
2. Trained and tested on imbalanced	B	Random Forest	Imbalanced	Imbalanced	0.6729
3. Model A tested on imbalanced	A	Logistic Regression	Balanced	Imbalanced	0.4291
4. Model B tested on balanced	B	Random Forest	Imbalanced	Balanced	0.4199

Observations

- **Model_B performed better overall**, especially on its native imbalanced test set, showing strong predictive power with **67.3% accuracy** and **F1 of 0.66**.
- **Model_A showed better balance** across all classes when tested on its own (balanced) test set but performed poorly on the imbalanced test set (accuracy dropped to **42.91%**).
- **Model_B struggled to generalize** when tested on the balanced test set (F1 dropped to **0.37**), indicating it **overfit the dominant classes** in the imbalanced data.

Recommendation

Recommended Model: Model B (Random Forest trained on imbalanced data)

1. Overall Accuracy Performance

- Model B achieves **highest accuracy (67.29%)** when trained and tested on the imbalanced dataset.
- In contrast, Model A, trained on balanced data, achieved only **44.73% accuracy** on the balanced test set, and **42.91% on the imbalanced set**.
- Even in cross-testing, Model B performs nearly the same (41.99%) on the balanced set as Model A does on the imbalanced one (42.91%).

2. Better Real-World Representation

- **Model B is trained on the imbalanced dataset**, which mirrors real-world class distributions more accurately.
- Since most production environments do **not have balanced class distributions**, Model B's training aligns better with deployment realities.

3. Class-wise Performance (Precision, Recall, F1-Score)

- Model B, though trained on imbalanced data, maintains **strong F1-scores for majority classes (2.0, 3.0, 4.0)**.
- Especially in its own domain (imbalanced test set), it produces **high recall and F1 for the dominant classes**, making it more reliable in high-frequency categories.
- Model A underperforms especially for minority and majority classes in cross-domain tests, showing poor generalization.

4. Robustness Across Domains

- While both models show performance degradation when tested outside their training distribution, **Model B's drop is more graceful**.
 - Drop from 67.29% → 41.99% (~25.3% decrease)
 - Model A drops from 44.73% → 42.91% (~1.8% decrease) but starts from a much lower base.
- This suggests **Model B has better generalization capacity** even when used on data distributions it wasn't trained on.

5. Algorithmic Strength

- Model B uses **Random Forest**, an ensemble-based non-linear method which naturally handles complex data distributions and class imbalance better.
- Model A uses **Logistic Regression**, which is linear and often struggles when classes are not separable linearly or when data imbalance is severe.