



SMS AND EMAIL SPAM DETECTION SYSTEM

Presented by

Gajanand Saini

MSA24008

Problem Statement

MILLIONS OF SPAM SMS/EMAILS ARE SENT DAILY, CAUSING CONFUSION FOR USERS. MANY PEOPLE CANNOT DISTINGUISH BETWEEN LEGITIMATE MESSAGES (E.G., BANK ALERTS) AND SPAM (E.G., FAKE LOTTERY SCAMS FINANCIAL FRAUD, PHISHING ATTACKS) AND PRIVACY BREACHES OCCUR DUE TO DECEPTIVE SPAM.

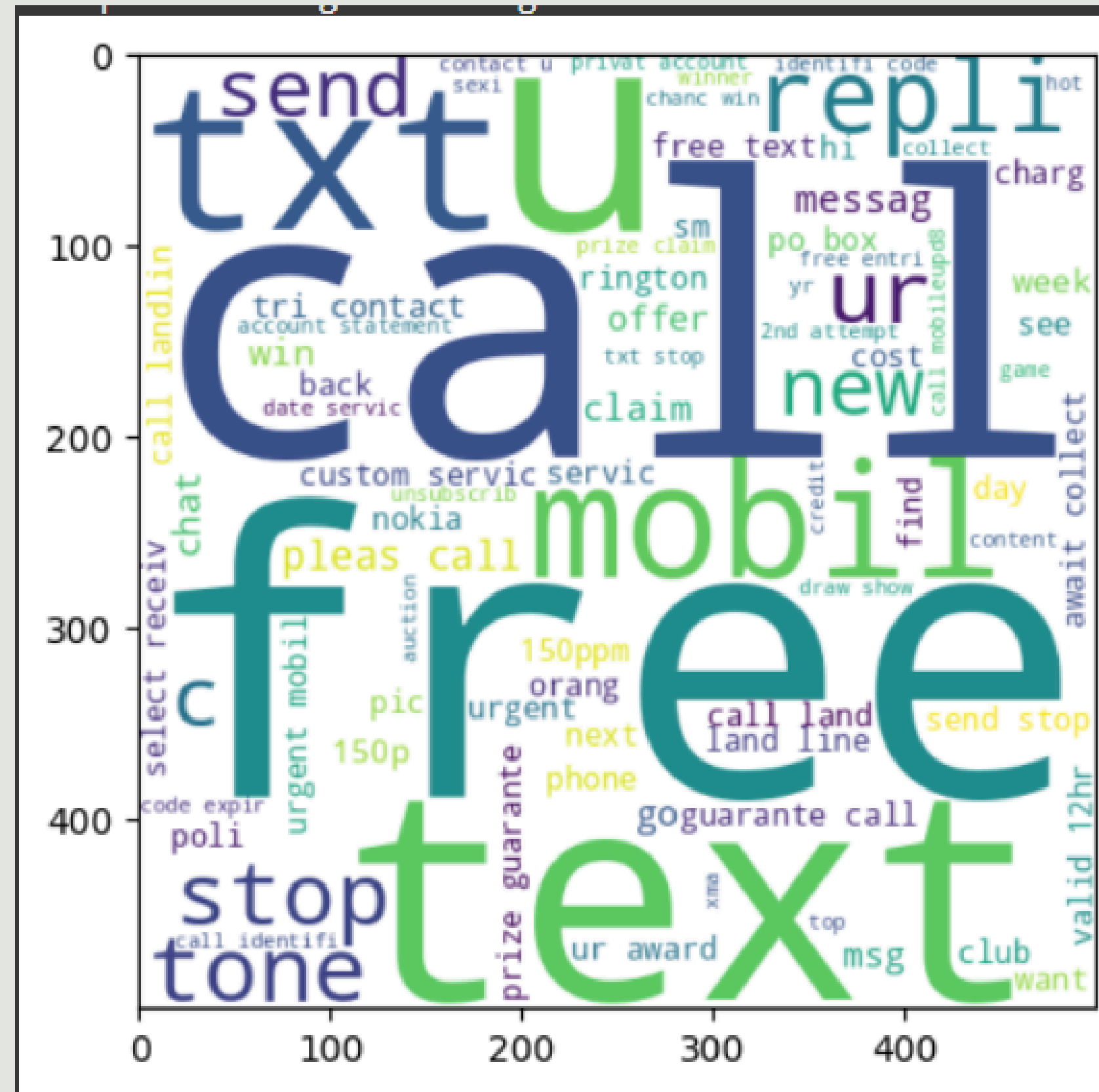
MACHINE LEARNING MODELS

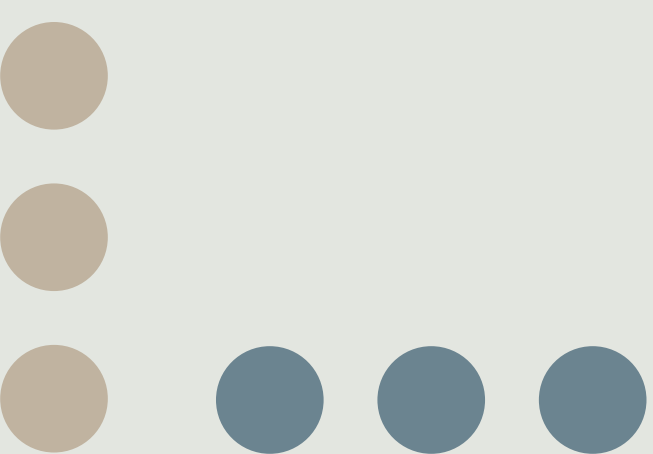
- Classifiers Used
- Naive Bayes
- Logistic Regression
- K-Nearest Neighbors
- Decision Trees
- Random Forest
- Support Vector Machine

Implementation Workflow

```
graph LR; A["[Data]"] --> B["[Preprocess]"]; B --> C["[Feature Extract]"]; C --> D["[Balance Data]"]; C --> E["[Train ML Models]"]; E --> F["[Deploy]"]
```

INTRODUCTION THE SPAM WORD

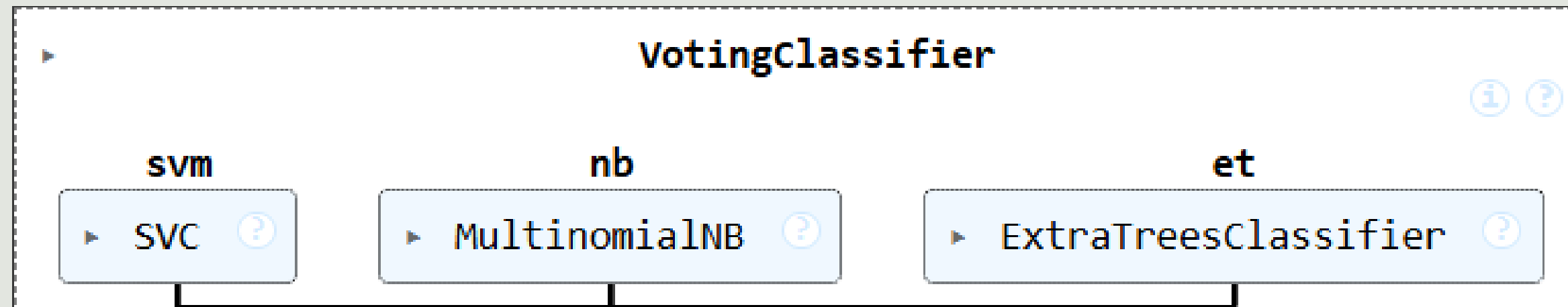


[illegible]

EXPERIMENTAL RESULTS WITH COUNT VECTORIZER AND SMOTE

	Algorithm	Accuracy	Precision
1	KN	0.910853	1.000000
8	ETC	0.972868	1.000000
5	RF	0.967054	1.000000
4	LR	0.974806	0.980198
6	AdaBoost	0.936047	0.938462
10	xgb	0.972868	0.935780
9	GBDT	0.950581	0.928571
3	DT	0.936047	0.913043
2	NB	0.977713	0.909836
7	BgC	0.955426	0.888889
0	SVC	0.967054	0.886957

VOTING CLASSIFIER



FINAL RESULT

	Algorithm	Variable	Value	Accuracy_max_ft_3000	Precision_max_ft_3000
0	KN	Accuracy	0.910853	0.910853	1.000000
1	ETC	Accuracy	0.972868	0.972868	1.000000
2	RF	Accuracy	0.967054	0.967054	1.000000
3	LR	Accuracy	0.974806	0.974806	0.980198
4	AdaBoost	Accuracy	0.936047	0.936047	0.938462
5	xgb	Accuracy	0.972868	0.972868	0.935780
6	GBDT	Accuracy	0.950581	0.950581	0.928571
7	DT	Accuracy	0.936047	0.936047	0.913043
8	NB	Accuracy	0.977713	0.977713	0.909836
9	BgC	Accuracy	0.955426	0.955426	0.888889
10	SVC	Accuracy	0.967054	0.967054	0.886957
11	KN	Precision	1.000000	0.910853	1.000000
12	ETC	Precision	1.000000	0.972868	1.000000
13	RF	Precision	1.000000	0.967054	1.000000
14	LR	Precision	0.980198	0.974806	0.980198
15	AdaBoost	Precision	0.938462	0.936047	0.938462
16	xgb	Precision	0.935780	0.972868	0.935780
17	GBDT	Precision	0.928571	0.950581	0.928571

CONCLUSION

- Using multiple vectorization methods and classifiers model achieved the highest accuracy in SMS spam detection.

Limitations

- Experiments limited to a single UCI dataset which may not represent diverse SMS spam patterns globally.
-

THANK YOU