



# On The Problem of Software Quality In Machine Learning Systems

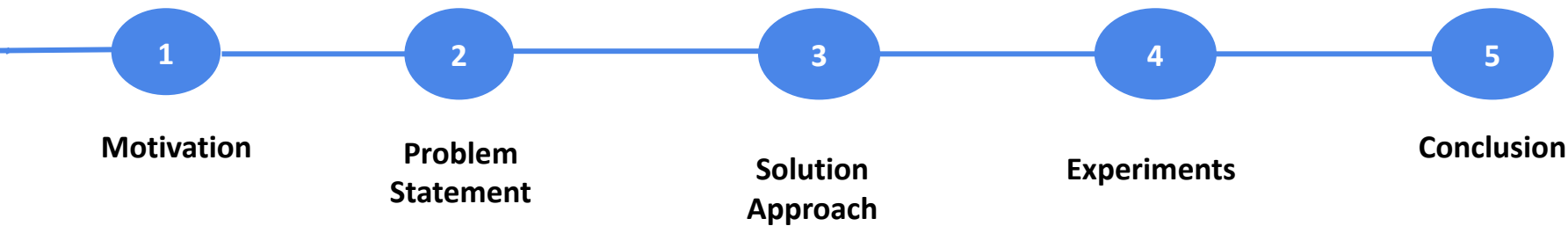
Haftamu Hailu Tefera | 21.11.2022

**Advisors:** Juan Soto | Prof. Dr. Volker Markl | Prof. Dr. Odej Kao

---



# Agenda





1

# Motivation



# Motivation

- Machine learning systems are deployed everywhere
- A ML bug refers to any imperfection that causes a inconsistency between existing and required conditions
- ML testing refers to any set of activities designed to reveal bugs



## Real World Facts

Although there is testing underway in the development of ML systems, there are still problems



The AI Incident Database [1]



Common Vulnerabilities and Exposures [2]

# THE RISKS DIGEST

Forum on Risks to the Public in Computers and Related Systems



2

# Problem Statement



## Problem Statement

How can users of machine learning systems, including practitioners and researchers be assured that Machine learning software is bug free?



3

# Solution Approach





## Solution Approach

1

2

3

### Study the past

- Stack Overflow posts
- ML and DL systems

### Familiarization

- Static analysis
- Smoke testing
- Metamorphic testing
- Model based testing

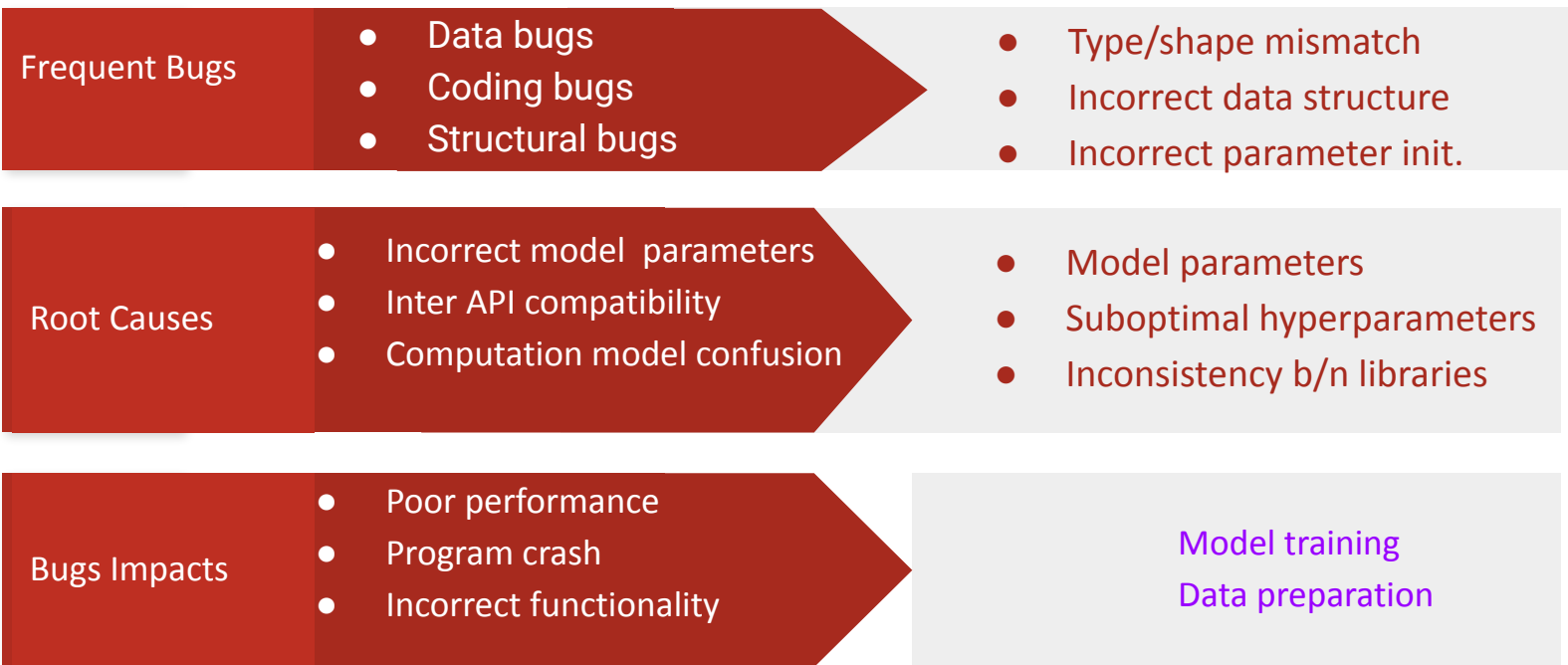
### Experiments

- Reproducibility
- Seeking to discover new bugs



# ML Systems Bug Analysis Study

We studied and analyzed frequent bugs, their root causes, and their impact [1,7].





## Terminologies

- **Static analysis** is the analysis of software programs performed without executing them
- **Smoke tests** assert that the most crucial functions of a program work
- **Metamorphic Testing** is founded on the principle of if the correct output for the input is not known, we can validate or verify the system based on the outputs of multiple related inputs
- **Model Based Testing** is a method where we test a piece of software based on its expected input and output.



4

# Experiments



# Experiment Design

1. Experiment class 1: Static analysis
2. Experiment class 2 : Smoke testing
3. Experiment class 3 : Model based testing
4. Experiment class 4 : Metamorphic testing
5. Experiment class 5 : Evaluation of classification algorithms



# Static Analysis

We analyzed open-source ML software repositories with SonarCloud

Scikit-learn

- Missing argument
- Unexpected argument
- Unreachable code

TensorFlow

- Method call error
- File not found
- Empty function body

Keras | PyTorch

- Unused variable or argument
- Code duplication with literals
- Image version tagging



# Smoke Testing

Herbold, Steffen. "Smoke testing for machine learning: simple tests to discover severe bugs"

## Scikit-learn

- Decision Tree
  - Random Forest
  - SVM | SGD | KNN
  - Naive Bayes
- Extreme class level imbalance
  - All data from the same class
  - Extremely large values
  - Attribute reference error

## Spark ML

- Decision Tree
  - Random Forest
  - Naive Bayes
  - Logistic Registic
- Unsupported arguments
  - Single training instance for a class
    - Division by zero (std)
  - Too many distinct categories



# Model Based Testing

A sentiment analysis model evaluated on Google Cloud NLP System

Please wait as we prepare the table data...

Capabilities	Minimum Functionality Test <i>failure rate % (over N tests)</i>	INVariance Test <i>failure rate % (over N tests)</i>	DIRectional Expectation Test <i>failure rate % (over N tests)</i>
<input type="checkbox"/> Vocabulary	48.6% (5)	16.2% (1)	34.6% (4)
<input type="checkbox"/> Robustness		13.6% (5)	
<input type="checkbox"/> NER		20.8% (3)	
<input type="checkbox"/> Fairness		1.6% (4)	
<input type="checkbox"/> Temporal	36.6% (1)		2.0% (1)
<input type="checkbox"/> Negation	100.0% (9)		
<input type="checkbox"/> SRL	90.8% (5)		





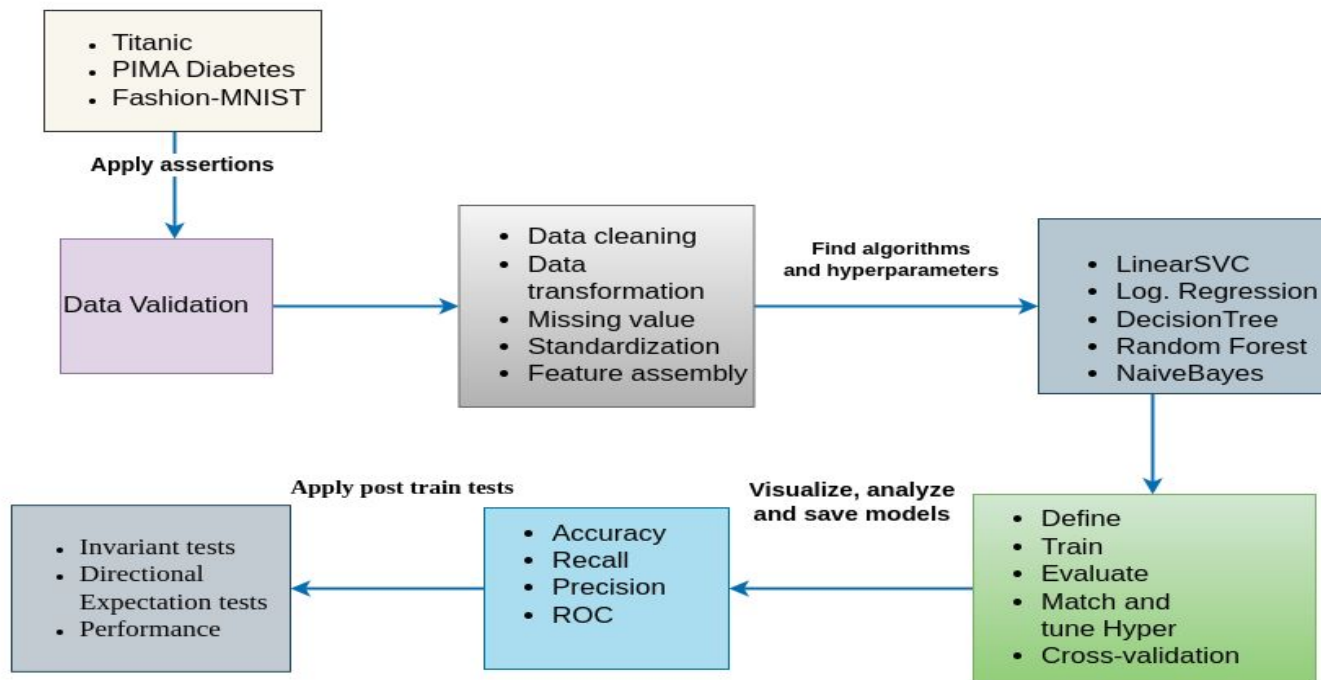
# Evaluation of Classification Algorithms

We evaluated classification algorithms from Spark ML, Scikit-learn and Keras network to detect or minimize bugs at an early stage.

Based on model behavioral testing for NLP models, we evaluated Scikit-learn models using invariant and directional expectation tests.



# Methodology



Scikit-learn

Spark ML

Keras

Jupyter Notebook

Python

Pandas DataFrame

Spark DataFrame



## Apply Pre-train Tests to Improve Data Quality

- To avoid unexpected column(s) , feature values and label values
- Ensure the range of values are in the specified range

---

```
1 assert set(titanic["Sex"].unique()) == set(("female", "male")), "Unknown gender"
2 assert set(titanic["Survived"].unique()) == set((0, 1)), "Unknown survived value"
3 assert set(titanic["Embarked"].unique()) == set(("S", "C", "Q")), "Unknown Embarked"
4 assert set(titanic["Pclass"].unique()) == set((1, 2, 3)), "Unknown Pclass value"
5 assert titanic["Age"].min() >= 0, "Age should be positive"
6 assert titanic["SibSp"].min() >= 0, "SibSp should be positive"
7 assert titanic["Fare"].min() >= 0, "Fare should be positive"
8 assert titanic["Parch"].min() >= 0, "Parch should be positive"
9 assert titanic["Parch"].max() <= 6, "Parch should be less than 7"
```

---



# Data Preprocessing and Feature Engineering

Missing value	Duplicate	Standardization	Encoding	Sampling	Vectorization
---------------	-----------	-----------------	----------	----------	---------------

```
1 round(titanic["Age"][titanic["Name"].str.contains("Mrs.")
2 & titanic["Sex"].str.contains("female")].mean())
3 titanic["Embarked"][titanic["Embarked"].isna()] = "S"
4 titanic.drop(index=titanic[titanic.drop(columns="Survived").duplicated()].index, inplace=True)
5 titanic["Sex"] = np.where(titanic["Sex"] == "male", 1, 0)
6 titanic = pd.concat([titanic[titanic['Survived'] == 0]
7     .sample(len(titanic[titanic['Survived'] == 1]), random_state=seed),
8 titanic[titanic['Survived'] == 1]],axis=0)
9
10 titanic = titanic.rename(columns={'Survived': 'label'})
11 va = VectorAssembler(inputCols = col, outputCol='features')
12 va_df = va.transform(data)
13 va_df = va_df.select(['features', 'label'])
14 return va_df.randomSplit([1-test_size, test_size], seed=seed)
```

Pandas

Spark



## Define and Match Hyperparameter

```
1 grid = {"LinearSVC": {"max_iter": [10, 50, 100, 500], "tol":  
2 [1e-2, 1e-3, 1e-4, 1e-5], "C": [0.1, 0.01, 0.001, 0.0001],  
3 "fit_intercept": [True, False]},  
4 "LogisticRegression": {"max_iter": [10, 50, 100, 500],  
5 "tol": [1e-2, 1e-3, 1e-4, 1e-5], "C": [0.1, 0.01, 0.001, 0.0001],  
6 "fit_intercept": [True, False]},
```

← Scikit-learn hyperparameters

Match hyperparameter names

```
1 para = { "LinearSVC": {"max_iter": ["maxIter"], "tol": ["tol"], "C": ["regParam"],  
2 "fit_intercept": ["fitIntercept"]},  
3 "LogisticRegression": {"max_iter": ["maxIter"], "tol": ["tol"], "C": ["regParam"],  
4 "fit_intercept": ["fitIntercept"]},
```

← PySpark hyperparameters



## 4. Mapping Algorithm Names

```
1  spark_grid = {}
2  spark_model_name_mapping = {
3      "GaussianNB": "NaiveBayes",
4      "GradientBoostingClassifier": "GBClassifier",
5      "MLPClassifier": "MultilayerPerceptronClassifier",
6      "OneVsRestClassifier": "OneVsRest"}
7
8  for sk_model_name, param_mapping in para.items():
9      spark_model_name = sk_model_name
10     if sk_model_name in spark_model_name_mapping.keys():
11         spark_model_name = spark_model_name_mapping.get(sk_model_name)
12     spark_grid[spark_model_name] = {}
13     for skt_name, pyspark_name in param_mapping.items():
14         spark_grid[spark_model_name][
15             pyspark_name[0]] = grid[sk_model_name][skt_name]
16  grid = spark_grid
```

(Scikit-learn: PySpark)

There are naming variations





## Model Training and Evaluation : Scikit-learn

Define the models (1)

Train the models (2)

Evaluate the models (2)

---

```
1 sklearn_classifiers = [  
2     sklearn_models.LinearSVC(random_state=seed),  
3     sklearn_models.LogisticRegression(random_state=seed),  
4  
5 for clf in sklearn_classifiers:  
6     clf.fit(x_train, y_train)  
7     accuracy, confusion, roc, precision, recall = eval_methods.eval(  
8         clf, y_test, clf.predict(x_test), ax, bx)
```

---



# Model Training and Evaluation : PySpark

```
1 pyspark_classifiers = [  
2     pyspark_models.LinearSVC(labelCol="label"),  
3     pyspark_models.LogisticRegression(labelCol="label"),  
4  
5 for clf in pyspark_classifiers:  
6     clf = clf.fit(train)  
7     real = np.array([  
8         1 if "1" in str(x) else 0  
9         for x in clf.transform(test).select("label").collect()])  
10    pred = np.array([  
11        1 if "1" in str(x) else 0  
12        for x in clf.transform(test).select("prediction").collect()])  
13    accuracy, confusion, roc, precision, recall = eval_methods.eval(  
14        clf, real, pred, ax, bx)
```

Define models

Train models

Evaluate models





# Apply Post-train Test to Ensure Expected Learned Behavior

```
1 path = "./models/titanic/sklearn/"
2 for i in os.listdir(path):
3     model = pickle.load(open(path+i, 'rb'))
4     X = datasets.get_titanic().iloc[291]
5     y = X["Survived"]
6     X = X[1:]
7     p2_prob = model.predict(np.array(X).reshape(1, -1))[0] #1.0
8     X['Embarked'] = 2.47593535
9     assert p2_prob == model.predict(np.array(X).reshape(1, -1))[0] #1.0
10    p2_prob = model.predict(np.array(X).reshape(1, -1))[0] # 1.0
11    X['Sex'] = 0.83739228
12    p2_male_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.56
13    assert p2_prob > p2_male_prob,
14    'Changing gender from female to male should decrease survival probability.'
15    X['Pclass'] = 0.95828974
16    p2_class_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.0
17    assert p2_prob > p2_class_prob,
18    'Changing class from 1 to 3 should decrease survival probability.'
```

Pclass, Sex and Fare are relevant features

```
1 titanic = datasets.get_titanic()
2 titanic.iloc[291]
```

Survived	1.000000
Pclass	-1.371412
Sex	-1.201153
Age	0.601520
SibSp	-0.470631
Parch	-0.428538
Ticket	-0.753102
Fare	0.715063
Embarked	-0.610667
Name: 61, dtype: float64	

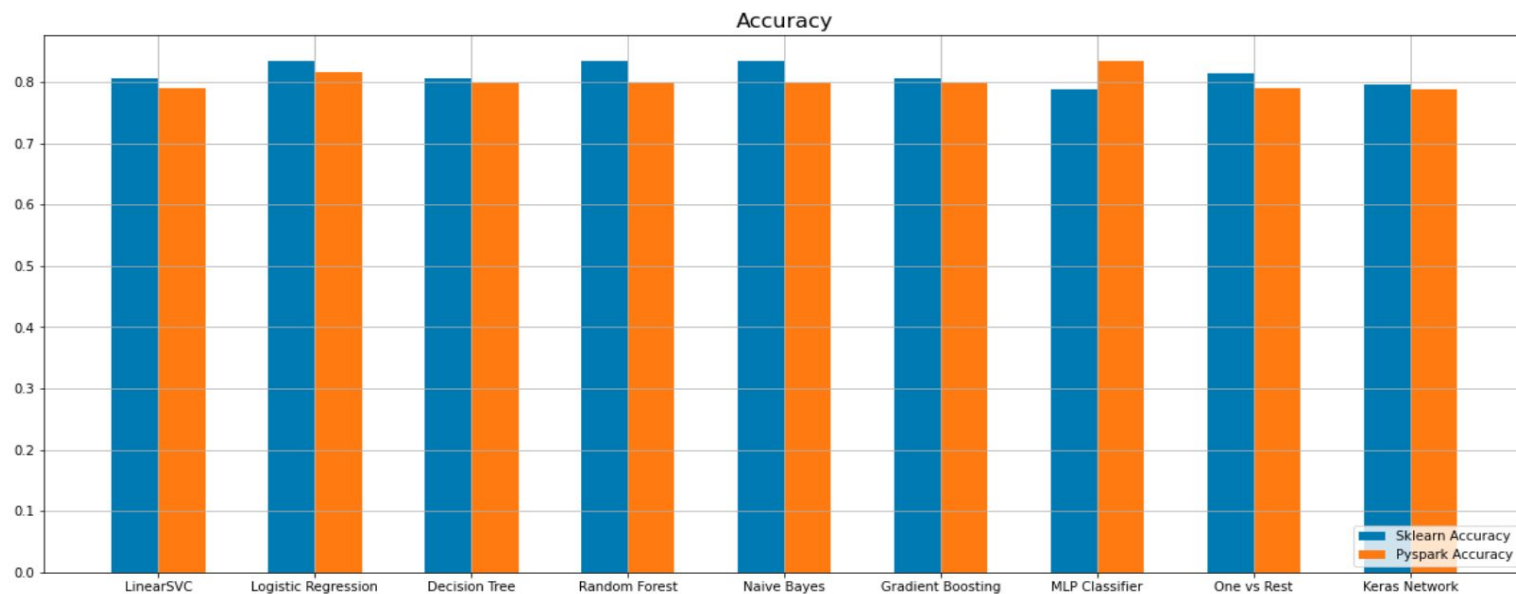


# Evaluation



# Accuracy

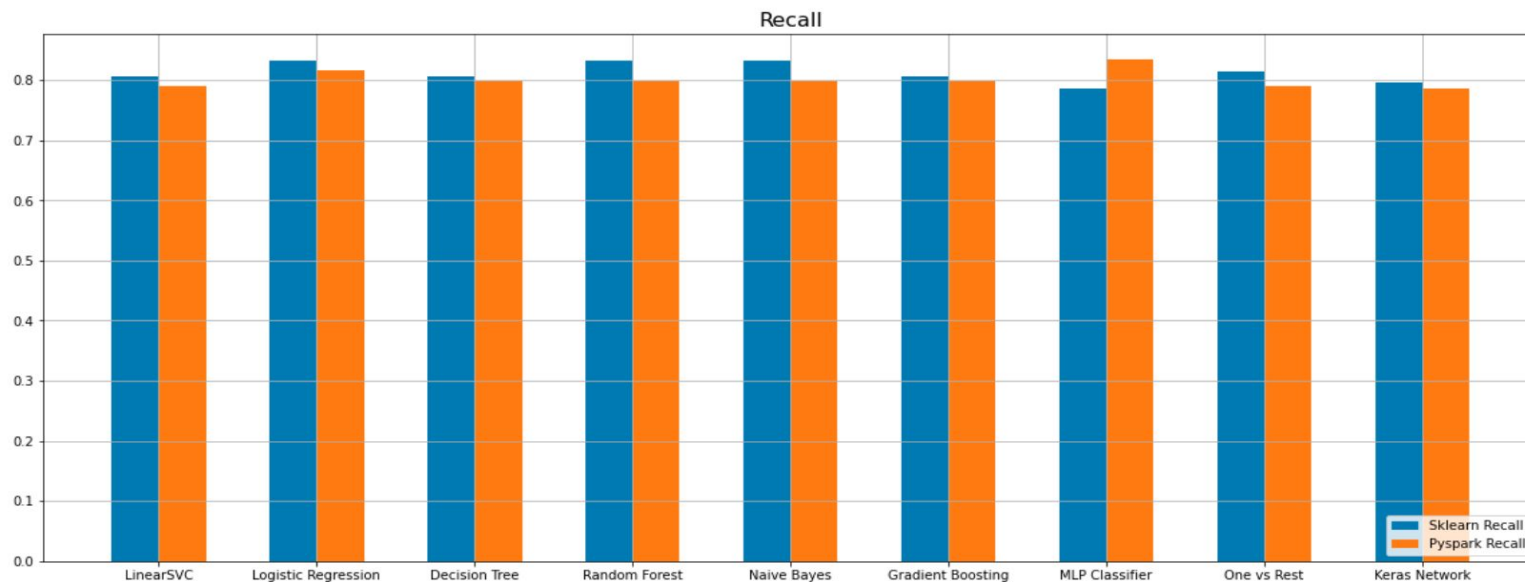
The ratio of correctly predicted instances to the total observations





# Recall

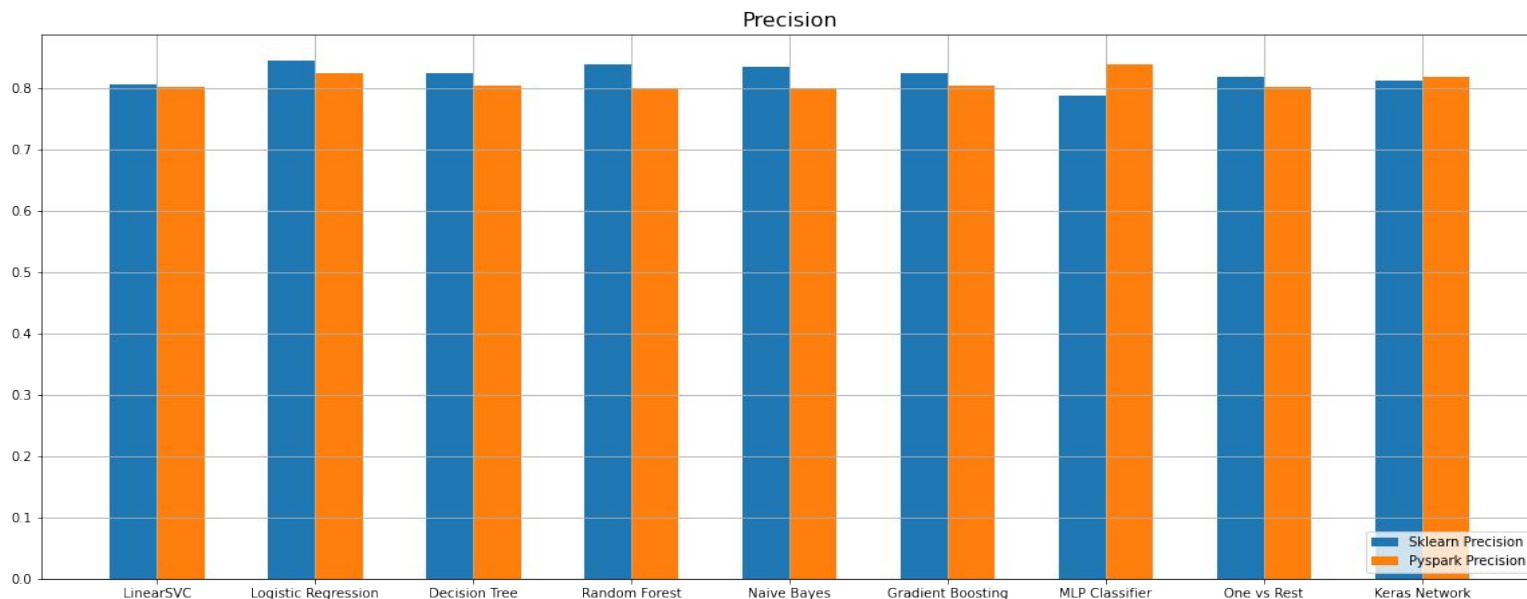
The ratio of **relevant instances** that has been retrieved over the total number of instances





# Precision

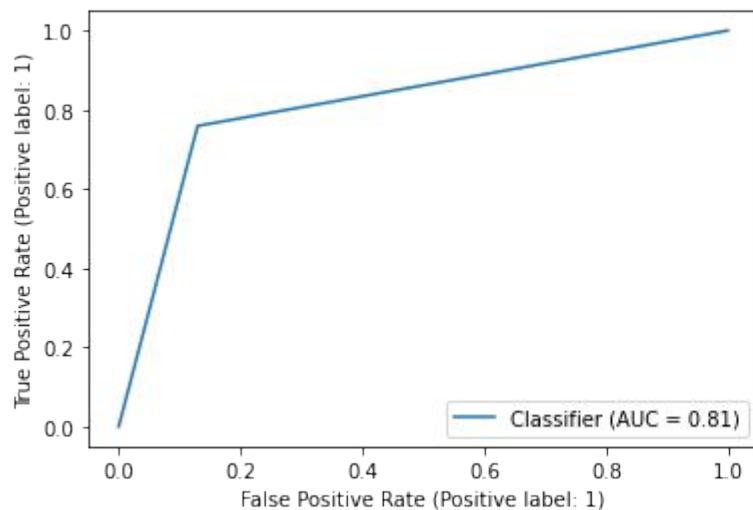
The fraction of **survived instances** among the **retrieved instances** (Survived + False Positives)



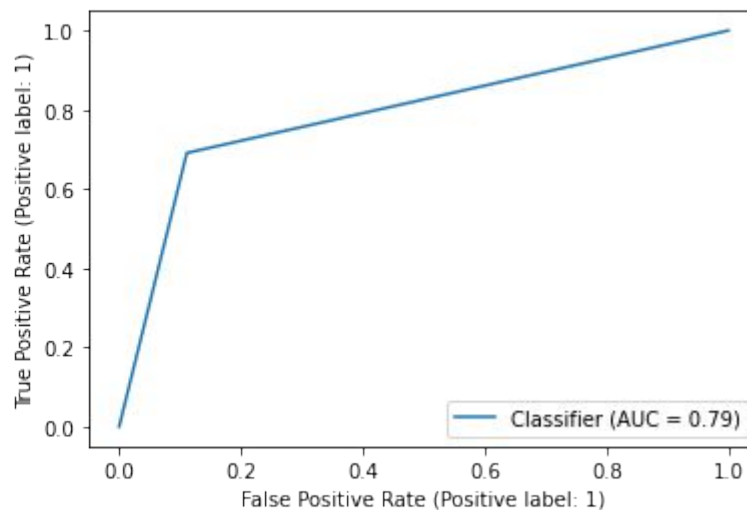


## ROC for Linear SVC Models

The LinearSVC from scikit-learn better identifies the survived passengers



a) Scikit-learn

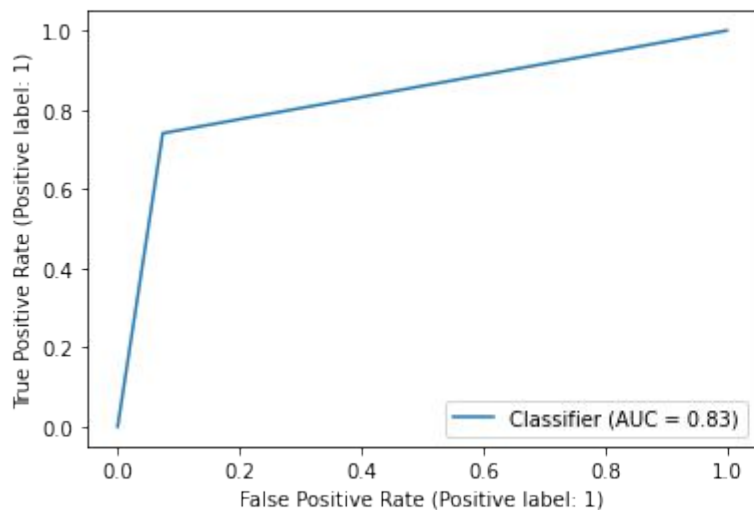


b) PySpark

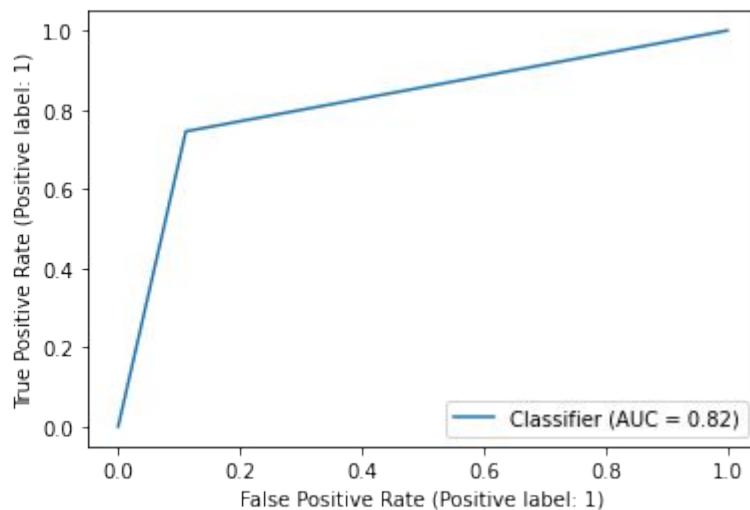


## ROC of Logistic Regression

The two logistic regressions identifies the survived passengers almost equally



a) Scikit-learn



b) PySpark



## Directional Expectation Tests for LinearSVC

The LinearSVC scikit-learn model did not pass the test

```
Testing model: LinearSVC(C=0.0001, max_iter=10, random_state=0, tol=0.01)
```

```
-----  
AssertionError                                Traceback (most recent call last)  
/tmp/ipykernel_38917/3038047054.py in <cell line: 2>()  
     9     X['Sex'] = 0.83739228 #Change gender  
    10     p2_male_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.56  
--> 11     assert p2_prob > p2_male_prob, 'Changing gender from female to male should decrease survival probab  
ility.'  
    12     X['Pclass'] = 0.95828974 # Change class  
    13     p2_class_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.0
```

```
AssertionError: Changing gender from female to male should decrease survival probability.
```

Change the gender value  
from **-1.2011** to **0.8373**





# Directional Expectation Tests for Gradient Boosting Classifier

The Gradient Boosting Classifier model did not pass the test

```
Testing model: GradientBoostingClassifier(learning_rate=0.01, max_depth=2, max_features='auto',
                                          random_state=0)
```

```
-----
AssertionError                                Traceback (most recent call last)
/tmp/ipykernel_38917/1720687962.py in <cell line: 2>()
    12     X['Pclass'] = 0.95828974 # Change class
    13     p2_class_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.0
--> 14     assert p2_prob > p2_class_prob, 'Changing class from 1 to 3 should decrease survival probability.'
    15     X['Fare'] = -0.575978 # Lower fare
    16     p2_fare_prob = model.predict(np.array(X).reshape(1, -1))[0] # 0.85
```

```
AssertionError: Changing class from 1 to 3 should decrease survival probability.
```

Change passenger class  
from **-1.37** to **0.9582**



5

# Conclusion



## Summary

- We study the past and perform static analysis
- We examined and experimented with three novel testing techniques
- Unable to detect new bugs and the performance difference is negligible
- We implemented pre-train and post train tests tests to improve data quality and ensure model learned behavior
- GBTClassifier and LinearSVC from PySpark supports binary classification



## Future Work

- Extend to model post-train testing for PySpark models
- Experiment tracking for best hyperparameters, performance scores, visualizations and other model artifacts
- Provide holistic benchmark with state-of-the-art by providing algorithms specific capabilities
- Apply Metamorphic testing and smoke testing to classification algorithms



## References

1. [AI Incident Database](#)
2. [CVE.org](#)
3. Islam, Md Johirul, et al. "What do developers ask about ml libraries? a large-scale study using stack overflow." *arXiv preprint arXiv:1906.11940* (2019).
4. Ribeiro, Marco Tulio, et al. "Beyond accuracy: Behavioral testing of NLP models with CheckList." *arXiv preprint arXiv:2005.04118* (2020).
5. Zhang, Jie M., et al. "Machine learning testing: Survey, landscapes and horizons." *IEEE Transactions on Software Engineering* (2020).
6. Breck, Eric, et al. "The ML test score: A rubric for ML production readiness and technical debt reduction." *2017 IEEE International Conference on Big Data (Big Data)*. IEEE, 2017.
7. Humbačová, Nargiz, et al. "Taxonomy of real faults in deep learning systems." *Proceedings of the ACM/IEEE 42nd International Conference on Software Engineering*. 2020.