

CSC 369

# Loan Acceptance Predictor

Anson Yamvinij, Ishaan Sathaye,  
Michael Montemurno

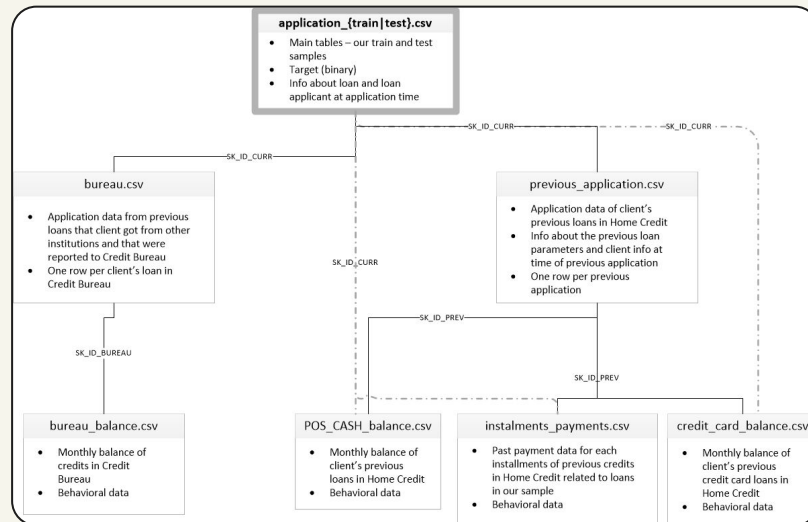
Confidential

Copyright ©



# Our Data

- Home Credit Default Risk dataset from Kaggle
- Loan Applicant as our observational unit
- TARGET: whether an applicant was able to repay the loan without trouble
- 307,511 rows by 122 columns



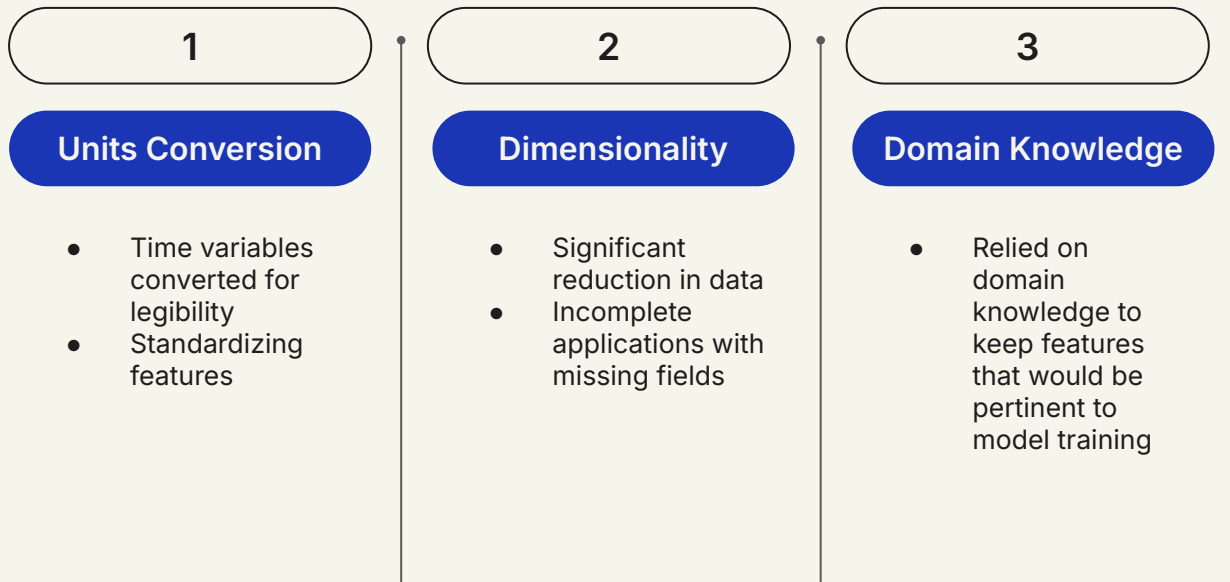
# Our Objective

Play the role of a bank and create a machine learning model to determine if we should give an applicant a loan.

We want our model to be accurate since we don't want to lose money on defaulted loans.



# Preprocessing



# Algorithms Implemented



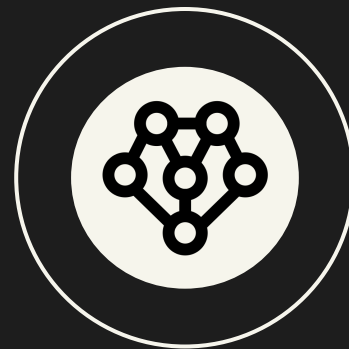
## Logistic Regression

With Regularization



## K-Nearest Classifier

With Standardization  
and Euclidean  
Distance



## Naive Bayes

# Logistic Regression

- Data split based on training and validation
- Features scaled for normalization
- Gradient Descent
- Class weights added for imbalance
- L2 Regularization for overfitting
- Sigmoid Activation with threshold at 0.5

```
1 // Train the model using gradient descent
2 for (_ <- 1 to maxIterations) {
3     val gradients = trainingRDD.map { case (features, label) =>
4         val linearCombination = features.zip(weights).map {
5             case (f, w) => f * w }.sum
6         val prediction = sigmoid(linearCombination)
7         val error = prediction - label
8         val weight = classWeights(label)
9         features.map(_ * error * weight)
10    }.reduce((g1, g2) => g1.zip(g2).map { case (x, y) => x + y })
11
12    // Update weights
13    val regularization = 0.315
14    weights = weights.zip(gradients).map { case (w, g) => w -
15        learningRate * g + regularization * w }
16 }
```

# K-Nearest Classifier

- No real “model”, used a 80/10/10 split for training data, validation data, and testing data
- Features scaled with Standardization
- To avoid too high dimensionality, used only four features per featureset.
- Multiple different configurations tested
  - Different values of K, with ~20 nearest neighbors being ideal
  - Different featuresets, with Amount of income, years employed, age, and amount of credit being the winners
  - Threshold for being considered positive at 0.55

```
def predict(train_data : RDD[(Double, (Double, Double, Double, Double, Double))],  
            test_data : RDD[(Double, (Double, Double, Double, Double, Double))],  
            k : Int) : RDD[(Double, Double)] = {  
  // print(k + "\n")  
  val revised_data = test_data  
    // do a cartesian product over all of the data  
    .cartesian(train_data)  
    // group by the key of the testing data  
    .groupBy(x => x._1._1)  
    // find the distance for each pair  
    .map(x => (x._1, x._2.map(pair => euclideanDistance(pair._1, pair._2))))  
  // output an rdd with each prediction  
  revised_data.map(x => (x._1, x._2.toList.sorted.take(k).fold(0.0)((total, n) => total + n)))  
    .map(x => (x._1, if(x._2 / k.toDouble < 0.5) 0.0 else 1.0 ))  
}
```

# Naive Bayes Classifier

- 80/20 train/test split
- Standardized features
- 3 Steps
  - Calculate mean and variance for each feature
  - Calculate the probability for each class
  - Calculate Log Likelihood
  - Highest = Prediction

$$P(A|B) = \frac{P(B|A) P(A)}{P(B)}$$

```
1 def predict(model: NaiveBayesModel, loanApplicant: LoanApplicant): Double = {
2   val featureValues = Array(
3     loanApplicant.NAME_CONTRACT_TYPE,
4     loanApplicant.FLAG_OWN_CAR,
5     loanApplicant.FLAG_OWN_REALTY,
6     loanApplicant.AMT_INCOME_TOTAL,
7     loanApplicant.AMT_CREDIT,
8     loanApplicant.AMT_ANNUITY,
9     loanApplicant.AMT_GOODS_PRICE,
10    loanApplicant.REGION_POPULATION_RELATIVE,
11    loanApplicant.CNT_FAM_MEMBERS,
12    loanApplicant.REGION_RATING_CLIENT,
13    loanApplicant.EXT_SOURCE_2,
14    loanApplicant.EXT_SOURCE_3,
15    loanApplicant.YEARS_BEGINEXPLUATATION_AVG,
16    loanApplicant.OBS_60_CNT_SOCIAL_CIRCLE,
17    loanApplicant.DEF_60_CNT_SOCIAL_CIRCLE,
18    loanApplicant.AMT_REQ_CREDIT_BUREAU_YEAR,
19    loanApplicant.AGE,
20    loanApplicant.YEARS_EMPLOYED,
21    loanApplicant.YEARS_REGISTERED,
22    loanApplicant.YEARS_ID_PUBLISH
23  )
24
25  val logLikelihoods = model.classProbs.keys.map { targetClass =>
26    val classProb = log(model.classProbs(targetClass))
27    val featureLikelihoods = featureValues.zipWithIndex.map { case (featureValue, featureIndex) =>
28      val (mean, variance) = model.featureProbs(targetClass)(featureIndex)
29      val likelihood = (-0.5 * log(2 * Pi * variance)) - (0.5 * pow(featureValue - mean, 2) / variance)
30      likelihood
31    }
32    val logLikelihood = classProb + featureLikelihoods.sum
33    (targetClass, logLikelihood)
34  }
35
36  val predictedClass = logLikelihoods.maxBy(_._2)._1
37  predictedClass
38 }
```



# Results

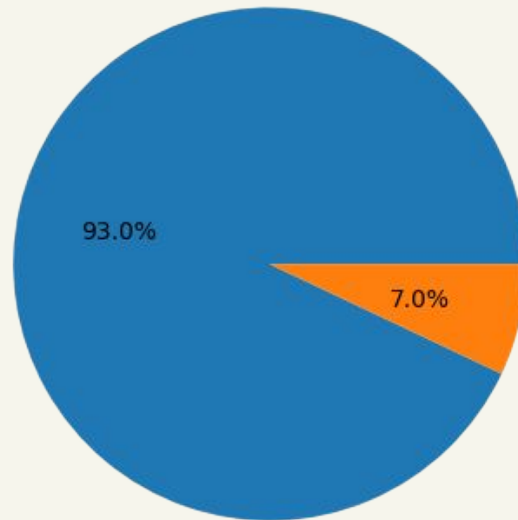
Accuracy, Precision, Recall, F1 Score

# Metrics

Model	Accuracy	Precision	Recall	F1-Score
Logistic Regression	0.75	0.11	0.29	0.16
KNN	0.90	0.037	0.068	0.024
Naive Bayes	0.876	0.231	0.235	0.233

# What we learned

- Working with RDDs and ML algorithms
- Logistic Regression vs Naive Bayes vs K-Nearest Neighbors Classifier
- Impact of Class Imbalance
- Feature Scaling
- Choosing correct metrics



*Majority are all other cases of loan statuses vs Minority of loan repayment difficulties*

# Challenges

## Scala Syntax

General challenges that were posed when developing the 3 algorithms

## Hyperparameters

Tuning hyperparameters took time and computation power, especially for convergence

## Model Selection

Deciding on suitable model among many algorithms required time, experimentation and other multiple factors

# Thank you.