

Verslag Distributed Databases

Dylan Cluyse, Laura Renders, Liam Goethals

11 december 2022

Inhoudsopgave

1	Inleiding	3
2	Tijd van een taxi-verplaatsing voorspellen met regressie.	4
3	Kredietkaartfraude bestrijden met binaire classificatie.	12
4	De inhoud van tweets detecteren met binaire NLP-classificatie.	15
5	Conclusie	19

Hoofdstuk 1

Inleiding

Tijdens het opleidingsonderdeel 'Distributed Databases' maakten wij kennis met Apache Spark. Spark biedt een data-analysetechnologie aan gericht op grootschalige gegevensverwerking. Deze technologie wordt aangeboden voor Java, Python, R en Scala. Voor dit opleidingsonderdeel gebruiken we Java als programmeertaal. Spark heeft pakketten die zich richten op andere aspecten binnen dataverwerking, één daarvan is *MLLib*. *MLLib* biedt *machine learning* technieken aan (ML). Wij kregen opdracht om via het platform Kaggle deel te nemen aan ML-gerichte competities, met het doel om zo meer vertrouwd te raken met ML en Spark.

In dit verslag nemen wij u mee in de wereld van Spark-ML. Wij tonen u met volle plezier drie concrete voorbeelden die verschillende regressie- en classificatiemethoden *highlighten*. U kan de gemaakte oefeningen digitaal terugvinden op deze link¹.

Deze opdracht was een proces van zelfstudie. De gehanteerde technieken zijn gebaseerd op de documentatie van Apache Spark² en van Spark MLLib³.

Allereerst willen wij de tijdsduur van een taxi-rit in downtown New York berekenen met regressie. Vervolgens detecteren wij kredietkaartfraude met behulp van binaire classificatie. Als derde oefening willen wij op basis van tekstinhoud achterhalen of een Tweet gerelateerd is aan een ramp. Tot slot geven wij u onze bevindingen mee van het MLLib pakket. Hier leggen we de aanpak van Spark en Sci-kit Learn, een pakket dat we in het opleidingsonderdeel Machine Learning zagen, parallel tegenover elkaar.

¹GitHub repository: https://github.com/Dyashen/Distributed_Databases/tree/main/Kaggle

²Apache Spark documentatie: <https://spark.apache.org/docs/latest/>

³Spark MLLib documentatie: <https://spark.apache.org/docs/latest/ml-guide.html>

Hoofdstuk 2

Tijd van een taxi-verplaatsing voorspellen met regressie.

Probleemstelling

Als eerste opdracht maakten wij de *New York City Taxi Trip Duration* wedstrijd. De competitie voorziet twee datasets: een training- en een testset. Het doel is om de totale tijdsduur van een taxirit in seconden te berekenen. Deze competitie kiest de inzending met de kleinste *Root Mean Square Log Error* (RMSLE). Volgende kolommen werden gegeven:

- Het ID en het verkopersnummer.
- Het aantal passagiers in een taxi.
- De longitude en latitude van het vertrek- en aankomstpunt.
- De pick-up en drop-off tijd.

De competitie kan u terugvinden op deze link. ¹. Wij namen inspiratie uit vooral de officiële documentatie en uit één notebook. De code, in *Python*, kan u hier ² terugvinden.

Algemeen

Wij declareren *final* variabelen. Deze worden meermaals gebruikt en worden niet aangepast doorheen de applicatie. Deze variabelen hebben een andere waarde bij de volgende twee oefeningen, maar ze komen in ieder model aan bod. De verhouding is de ratio waarop we de trainingset splitsen. Het label is de kolom waarop ons model zal voorspellen. De prediction houdt de naam van de prediction-kolom bij. We houden de metriek bij die we nodig hebben. Als laatste geven wij de metriek mee. Voorlopig kijken we naar het model met de beste RMSE waarde, daarop berekenen we de RMSLE.

```
1 private static final double[] verhouding = { 0.8, 0.2 };
2 private static final String label = "target";
3 private static final String prediction = "prediction";
```

Uiteraard starten wij de applicatie met een Spark-object. Wij passen één instelling aan om zo de irrelevante CLI-informatie achterwege te laten.

```
1 private static SparkSession spark = SparkSession.builder()
2     .appName("DisasterTweet")
3     .master("local[1]")
4     .getOrCreate();
5
6 spark.sparkContext().setLogLevel("ERROR");
```

¹Link naar competitie: <http://bit.ly/3umBKNS>

²Link naar gevolgde notebook: <http://bit.ly/3F2Aofj>

Data ophalen

De Kaggle-competitie bevat een training- en een testset. Dit in CSV-formaat. We kiezen ervoor om het schema vooraf op te bouwen. Dit geeft ons een snelheidsvoordeel, want met een *pre-built* schema moet Spark zelf niet de datatypes achterhalen. Voor de datums werken wij met een *Timestamp* type, want het bevat zowel de datum als het tijdstip wanneer een persoon werd opgehaald of afgezet. De code vindt u hieronder terug:

```
1 private static Dataset<Row> getTraining() {
2
3     List<StructField> fields = Arrays.asList(
4         DataTypes.createStructField("id", DataTypes.StringType, false),
5         DataTypes.createStructField("vendor_id", DataTypes.DoubleType, false),
6         DataTypes.createStructField("pickup_datetime", DataTypes.TimestampType, false),
7         DataTypes.createStructField("dropoff_datetime", DataTypes.TimestampType, false),
8         DataTypes.createStructField("passenger_count", DataTypes.DoubleType, false),
9         DataTypes.createStructField("pickup_longitude", DataTypes.DoubleType, false),
10        DataTypes.createStructField("pickup_latitude", DataTypes.DoubleType, false),
11        DataTypes.createStructField("dropoff_longitude", DataTypes.DoubleType, false),
12        DataTypes.createStructField("dropoff_latitude", DataTypes.DoubleType, false),
13        DataTypes.createStructField("store_and_fwd_flag", DataTypes.StringType, false),
14        DataTypes.createStructField("trip_duration", DataTypes.DoubleType, false)
15    );
16
17    StructType schemaData = DataTypes.createStructType(fields);
18
19    Dataset<Row> dataset = spark.read()
20        .option("header", true)
21        .schema(schemaData)
22        .csv("src/main/resources/train.csv");
23
24    return dataset
25        .withColumn("hour", hour(col("pickup_datetime")))
26        .withColumn("day", dayofweek(col("pickup_datetime")))
27        .drop("id", "pickup_datetime", "dropoff_datetime");
28 }
29
30 private static Dataset<Row> getTest() {
31     List<StructField> fields = Arrays.asList(
32         DataTypes.createStructField("id", DataTypes.StringType, false),
33         DataTypes.createStructField("vendor_id", DataTypes.DoubleType, false),
34         DataTypes.createStructField("pickup_datetime", DataTypes.TimestampType, false),
35         DataTypes.createStructField("passenger_count", DataTypes.DoubleType, false),
36         DataTypes.createStructField("pickup_longitude", DataTypes.DoubleType, false),
37         DataTypes.createStructField("pickup_latitude", DataTypes.DoubleType, false),
38         DataTypes.createStructField("dropoff_longitude", DataTypes.DoubleType, false),
39         DataTypes.createStructField("dropoff_latitude", DataTypes.DoubleType, false),
40         DataTypes.createStructField("store_and_fwd_flag", DataTypes.StringType, false)
41     );
42
43     StructType schema = DataTypes.createStructType(fields);
44
45     Dataset<Row> dataset = spark.read()
46         .option("header", true)
47         .schema(schema)
48         .csv("src/main/resources/test.csv");
49
50     return dataset.withColumn("hour", hour(col("pickup_datetime")))
51         .withColumn("day", dayofweek(col("pickup_datetime")))
52         .drop("id", "pickup_datetime");
53 }
```

Data Cleaning

De data is te *ruw*, en daarom maken wij een *clean*-functie aan. Zonder de *outliers* is het model nauwkeuriger. Wij passen het volgende stappenplan toe:

1. De *pick-up* en *drop-off* datetime wordt meegegeven als datetime-object. We veranderen de datetime-kolommen *pickuptime* naar *hour* en *day*. Zo krijgen we een pure numerieke waarde (bijvoorbeeld 1 tot en met 7 voor *day*).
2. We verwijderen de coördinaten-*outliers* uit de dataset. Hiervoor kijken we naar de rijen met een veel te hoge *longitude* en/of *latitude*.
3. We verwijderen de *distance-outliers* uit de dataset. Dit gebeurt na het berekenen van de distance. De functie hiervoor komt zo dadelijk aan bod.
4. We verwijderen alle rijen die géén passagiers meenemen. Alle rijen waar het aantal passagiers gelijk is aan 0 laten we vallen.
5. We *shufflen* de dataset. Dit wordt bij het splitsen nog eens gedaan, maar zo zijn we dubbel zeker dat er geen strikte volgorde is bij het splitsen.

```
1 private static Dataset<Row> clean(Dataset<Row> dataframe) {
2
3     dataframe = dataframe.na().drop(); // alle rijen met null-waarden verwijderen
4
5     double latMin = 40.6;
6     double latMax = 40.9;
7     double longMin = -74.25;
8     double longMax = -73.7;
9
10    // outliers verwijderen o.b.v. de longitude/latitude values
11    dataframe = dataframe
12        .where(col("pickup_longitude").$greaterEq(longMin))
13        .where(col("dropoff_longitude").$greaterEq(longMin))
14        .where(col("pickup_latitude").$greaterEq(latMin))
15        .where(col("dropoff_latitude").$greaterEq(latMin))
16        .where(col("pickup_longitude").$lessEq(longMax))
17        .where(col("dropoff_longitude").$lessEq(longMax))
18        .where(col("pickup_latitude").$lessEq(latMax))
19        .where(col("dropoff_latitude").$lessEq(latMax));
20
21    // outliers verwijderen o.b.v. de distance
22    double rightOuter = dataframe
23        .select(avg(col("distance")).plus(stddev(col("distance"))))
24        .first()
25        .getDouble(0);
26
27    double leftOuter = dataframe
28        .select(avg(col("distance")).minus(stddev(col("distance"))))
29        .first()
30        .getDouble(0);
31
32    dataframe = dataframe
33        .where(col("distance").$lessEq(rightOuter))
34        .where(col("distance").$greaterEq(leftOuter));
35
36    dataframe = dataframe
37        .where(col("passenger_count").$greater(0)); // minstens één passagier in de taxi
38
39    dataframe = dataframe
40        .orderBy(rand()); // willekeurige volgorde
41
42    return dataframe;
43 }
```

User-Defined Functions

De coördinaten van het vertrek- en aankomstpunt liggen te dicht bij elkaar. Daarnaast bieden deze vier waarden weinig inbreng aan het model. Daarom berekenen wij de afstand tussen de vertrek- en aankomstcoördinaten met de *haversine*-formule, ofwel de formule om de grootcirkelafstand te berekenen. Deze formule dient om de afstand op een bol oppervlak te berekenen. Hiervoor maken we een *user-defined-function* (UDF) aan dat vier parameters opvraagt: *de pickup longitude, pickup latitude, de dropoff-longitude en de dropoff-latitude*. We volgen de wiskundige formule, in drie aparte stappen. Wiskundige bewerkingen, indien mogelijk zoals de vierkantswortel, doen we met de Math-library van Java.

```
1 private static Dataset<Row> addDistance(Dataset<Row> dataframe) {
2     return dataframe
3         .withColumn("distance", call_udf("haversine",
4             col("pickup_latitude"), col("pickup_longitude"),
5             col("dropoff_latitude"), col("dropoff_longitude")))
6         .drop("pickup_longitude", "dropoff_longitude", "pickup_latitude", "dropoff_latitude");
7 }
8
9 UDF4<Double, Double, Double, Double, Double> haversine = new UDF4<>() {
10     public Double call(Double pickupLatitude, Double pickupLongitude, Double dropoffLatitude,
11         Double dropoffLongitude) throws Exception {
12
13         double radius = 6371 * 1.1; // Radius van de Aarde
14         double lat1 = pickupLatitude; double lon1 = pickupLongitude; // pick-up
15         double lat2 = dropoffLatitude; double lon2 = dropoffLongitude; // drop-off
16
17         double deltaLat = Math.toRadians(lat2 - lat1);
18         double deltaLon = Math.toRadians(lon2 - lon1);
19
20         double a = Math.sin(deltaLat / 2) * Math.sin(deltaLat / 2) + Math.cos(Math.toRadians(lat1))
21             * Math.cos(Math.toRadians(lat2)) * Math.sin(deltaLon / 2) * Math.sin(deltaLon / 2);
22         double c = 2 * Math.atan2(Math.sqrt(a), Math.sqrt(1 - a));
23         double d = radius * c;
24         return d;
25     }
26 };
27 spark.udf().register("haversine", haversine, DataTypes.DoubleType);
28 train = addDistance(train); test = addDistance(test);
```

Wij willen ook de gemiddelde snelheid van een taxi in *miles* per hour (mph) bekijken. Voor deze berekening maken wij een tweede UDF aan met als input-parameters: de totale afgelegde afstand en de totale tijdsduur van een rit. De output is de gemiddelde snelheid van een taxi per rit als kommagetal. De testset ontbreekt de verlopen tijd van een taxirit, dus we gebruiken deze kolom enkel als extra statistiek bij het analyseren van de data.

```
1 private static Dataset<Row> addSpeed(Dataset<Row> dataframe) {
2     return dataframe.withColumn("speed", call_udf("speed", col("distance"), col("label")));
3 }
4
5 UDF2<Double, Double, Double> speed = new UDF2<>() {
6     public Double call(Double distance, Double time) throws Exception {
7         return distance / (time / 3600);
8     }
9 };
10 spark.udf().register("speed", speed, DataTypes.DoubleType);
```

Als tussentijdse controle tonen wij, in de terminal, de gemiddelde snelheid en gemiddelde afstand per dag en uur. Zo controleren we de transformaties voor: *distance, speed, hour en day*. Hieronder ziet u de gemiddelde afstand/snelheid per dag of per uur.

```
1 train.groupBy("hour").mean("speed").orderBy("hour").show();
2 train.groupBy("day").mean("speed").orderBy("day").show();
3 train.groupBy("hour").sum("distance").orderBy("hour").show();
4 train.groupBy("day").sum("distance").orderBy("day").show();
```

Opbouwen van de pipeline

Voor de pipeline werken we met vier stappen:

1. We hebben een kolom met een vlag, ofwel een categorische tekstwaarde. Dit moeten we *one-hot-encoden* of *labelen*. In Spark gebruiken we een *StringIndexer*.
2. Spark MLLib verplicht dat we werken met een vector van feature-waarden. Alle kolommen omzetten naar één featurekolom doen we met een *VectorAssembler*.
3. De waarden verschillen sterk qua grootte, en dit heeft een nadelig effect op ons model bij regressie. We moeten de waarden schalen. In Spark zijn er twee manieren: *MinMaxScaling* en *StandardScaling*. Wij gebruiken hier *MinMaxScaler* wat de waarden tussen 0 en 1 zal plaatsen. De hoogste waarde zal dicht bij één aanleunen. De kleinste waarde zal dicht bij nul aanleunen.
4. Als laatste object geven wij het regressiemodel mee. Wij testen vier verschillende modellen: lineaire regressie, *Random Forest Regressor* (RFR), *Gradient Boost Regressor* (GBR) en een gegeneraliseerd lineaire regressie (GLR). Wij vonden enkel het lineaire regressiemodel terug in de notebooks, maar uit interesse waagden wij een poging met drie extra modellen.

Hieronder vindt u de pipeline voor het lineaire regressiemodel. De volgende code vindt u terug in de *main*-methode.

```
1 StringIndexer indexer = new StringIndexer ()
2   .setHandleInvalid ("keep")
3   .setInputCol ("store_and_fwd_flag")
4   .setOutputCol ("flag_ind");
5
6 VectorAssembler assembler = new VectorAssembler ()
7   .setInputCols (new String [] {
8     "vendor_id", "passenger_count", "hour", "day", "flag_ind", "distance"
9   })
10  .setOutputCol ("features");
11
12 MinMaxScaler minMax = new MinMaxScaler ()
13   .setInputCol (assembler.getOutputCol ())
14   .setOutputCol ("scaledFeatures");
15
16 LinearRegression linReg = new LinearRegression ()
17   .setLabelCol (label)
18   .setFeaturesCol (minMax.getOutputCol ());
19
20 Pipeline pipelineLinReg = new Pipeline ()
21   .setStages (new PipelineStage [] {
22     indexer, assembler, minMax, linReg
23   });
24
25 PipelineModel model = pipelineLinReg.fit (trainSetSplit);
26 Dataset<Row> trainedLinReg = model.transform (trainTestSplit);
```

Bij *machine learning* splitsen we de trainingsset in twee delen op. We stellen een vaste verhouding in door middel van een *final* variabele. Wij kiezen voor de verhouding [80% - 20%]. Met de ingebouwde *randomSplit* methode splitsen we de dataframe.

```
1 Dataset<Row>[] datasets = train.randomSplit (verhouding, 42);
2 Dataset<Row> trainSetSplit = datasets[0];
3 Dataset<Row> trainTestSplit = datasets[1];
```


Finetuning

Om het meest optimale model te vinden, moeten we finetunen. In Spark pakken wij dit aan met een *ParamGridBuilder*. Dit is het equivalent van werken met een *GridSearch* in *Sci-kit Learn*. Wij geven verschillende waarden per hyperparameter mee. Iedere combinatie zal worden getest. Hoe meer parameters, hoe langer de uitvoertijd van de applicatie.

We passen crossvalidatie toe op ons model. Het beste model kiezen gebeurt op basis van een metriek. Zoals eerder gezegd willen we de beste RMSE hebben. We geven mee dat we voor de pipeline dezelfde objecten gebruiken als hierboven. Het *tweaken* gebeurt binnenin de pipeline. Hieronder een voorbeeld van een RFR-pipeline.

```
1 RandomForestRegressor rfr = new RandomForestRegressor ()
2   .setLabelCol (label)
3   .setFeaturesCol (minMax.getOutputCol ()) ;
4
5 ParamMap[] paramGridRFR = new ParamGridBuilder ()
6   .addGrid (rfr.maxDepth () , new int [] { 5, 10, 15 })
7   .addGrid (rfr.numTrees () , new int [] { 10, 15, 20 })
8   .build () ;
9
10 CrossValidator cvRFR = new CrossValidator ()
11   .setEstimator (rfr)
12   .setEvaluator (regEval)
13   .setEstimatorParamMaps (paramGridRFR) ;
14
15 Pipeline pipelineRFR = new Pipeline ()
16   .setStages (new PipelineStage [] { indexer , assembler , minMax , cvRFR }) ;
17
18 PipelineModel pipelineModelRFR = pipelineRFR.fit (trainSetSplit) ;
19 Dataset<Row> trainedRFR = pipelineModelRFR.transform (trainTestSplit) ;
```

Een alternatief voor crossvalidatie is *trainvalidationsplit* (TVS). Het verschil is dat TVS hier één datasetpaar maakt, terwijl crossvalidatie afhankelijk is van het aantal *folds*. Dit gebruiken we bij GLR en GBR. Hieronder ziet u de pipeline voor het GLR-model.

```
1 GeneralizedLinearRegression glr = new GeneralizedLinearRegression ()
2   .setFamily ("gaussian")
3   .setLink ("identity")
4   .setLabelCol (label)
5   .setFeaturesCol (minMax.getOutputCol ())
6   .setMaxIter (10) .setRegParam (0.3) ;
7
8 ParamMap[] paramGridGLR = new ParamGridBuilder ()
9   .addGrid (glr.maxIter () , new int [] { 15, 20, 25 })
10  .addGrid (glr.regParam () , new double [] { 0.6, 0.9 })
11  .build () ;
12
13 TrainValidationSplit trainValidationSplitGLR = new TrainValidationSplit ()
14   .setEstimator (glr)
15   .setEvaluator (regEval)
16   .setEstimatorParamMaps (paramGridGLR)
17   .setTrainRatio (verhouding [0]) ;
18
19 Pipeline pipelineGLR = new Pipeline () .setStages (new PipelineStage [] { indexer , assembler ,
20   minMax , trainValidationSplitGLR }) ;
21
22 PipelineModel pipelineModelGLR = pipelineGLR.fit (trainSetSplit) ;
23 Dataset<Row> trainedGLR = pipelineModelGLR.transform (trainTestSplit) ;
```

Evaluatie en controle

Het achterhalen van de interessante features doen we met een *correlation matrix*.

```
1 private static void printCorrelation (Dataset<Row> dataframe) {
2     Row r1 = Correlation.corr(dataframe, "features").head();
3     System.out.printf("\n\nCorrelation Matrix\n");
4     Matrix matrix = r1.getAs(0);
5     for (int i = 0; i < matrix.numRows(); i++) {
6         for (int j = 0; j < matrix.numCols(); j++) {
7             System.out.printf("%.2f \t", matrix.apply(i, j));
8         }
9         System.out.println();
10    }
11 }
```

Met Crossvalidatie voorspelt het model met de meest optimale combinatie hyperparameters. Wij vragen de volgende metrieke op: *Mean Squared Error* (MSE), *Root Mean Squared Error* (RMSE), *Mean Absolute Error* (MAE) en de R^2 -waarde. Als laatst berekenen we de RMSLE door de log-functie te gebruiken op de RMSE. Met een *for*-lus kunnen wij voor iedere gevraagde metriek informatie ophalen.

```
1 private static void printRegressionEvaluation (Dataset<Row> dataframe) {
2     String[] metricTypes = { "mse", "rmse", "r2", "mae" };
3     System.out.printf("\n\nMetrics:\n");
4     for (String metricType : metricTypes) {
5         RegressionEvaluator evaluator = new RegressionEvaluator()
6             .setLabelCol("label")
7             .setPredictionCol("prediction")
8             .setMetricName(metricType);
9
10        double calc = evaluator.evaluate(dataframe);
11        System.out.printf("%s: \t%.5f \n", metricType, calc);
12    }
13
14    RegressionEvaluator evaluator = new RegressionEvaluator()
15        .setLabelCol("label")
16        .setPredictionCol("prediction")
17        .setMetricName("rmse");
18
19    double calc = evaluator.evaluate(dataframe);
20    double log = Math.log(calc);
21    System.out.printf("%s: \t%.5f \n", "rmsle", log);
22 }
```

Tussen de *fit* en de *transform* printen we de meest optimale hyperparameters uit. De beste combinatie wordt achterhaald door of het *crossvalidatiemodel*, ofwel de *traintestvalidation*. Hieronder een voorbeeld van de *traintestvalidation*. We moeten rekening houden met de volgorde waarin we de parameters toevoegen aan de ParamGridBuilder. Die volgorde komt overeen met de index van de parameters in de *ParamMap*-array.

```
1 System.out.printf("Ideale parameters:\nMax Iter: %s\nReg params: %s",
2     trainValidationSplitGLR.getEstimatorParamMaps()[0].get("glr.maxIter").toString(),
3     trainValidationSplitGLR.getEstimatorParamMaps()[1].get("glr.regParam").toString()
4 );
```

Hieronder een voorbeeld van hoe we de optimale hyperparameters ophalen bij een cross-validatiemodel.

```
1 System.out.printf("Ideale parameters:\nMax depth: %s\nNumber of trees: %s",
2     cvRFR.getEstimatorParamMaps()[0].get("rfr.maxDepth").toString(),
3     cvRFR.getEstimatorParamMaps()[1].get("rfr.numTrees").toString()
4 );
```

Als laatste evaluatie willen we kijken naar de marge tussen de verschillende voorspelde waarden en de effectieve waarden. Zo weten we hoe sterk het model afwijkt en ook in hoeveel gevallen het in die mate afwijkt van de echte waarden. We werken met de volgende bereiken:

1. 0-50 seconden
2. 50-200 seconden
3. 200-500 seconden
4. 500-5000 seconden
5. 5000+ seconden

```
1 private static Dataset<Row> getRangeDataFrame(Dataset<Row> dataframe) {
2     dataframe = dataframe.withColumn("margin",
3         abs(col("prediction").minus(col("trip_duration"))));
4
5     Dataset<Row> range = dataframe
6         .withColumn("range",
7             when(col("margin").leq(50), "0-50")
8             .when(col("margin").leq(200), "50-200")
9             .when(col("margin").leq(500), "200-500")
10            .when(col("margin").leq(5000), "500-5000")
11            .otherwise("5000+"));
12     return range.groupBy("range").count();
13 }
```

Inzending

Voor de opdracht gaan we alle voorspellingen opslaan en koppelen aan het ID. We hergebruiken deze code, dus we plaatsen dit in een functie. Bij het opslaan geven we een save mode mee. De *Overwrite* zorgt ervoor dat we geen foutmelding krijgen als de map al bestaat. Het resultaat wordt telkens overgeschreven. We geven u mee dat deze code herhaalt wordt voor alle oefeningen. In het verslag vermelden we enkel hier deze functie.

```
1 private static void createSubmission(Dataset<Row> dataframe, String fnaam) {
2     String output = "src/main/resources/" + fnaam + ".csv";
3     dataframe = dataframe.select("id", "prediction");
4     dataframe.write().mode(SaveMode.Overwrite).csv(output);
5 }
6
7 // In de main-methode
8 createSubmission(predictedLogReg, "submissionLogReg");
```

Conclusie

De opdracht vergde een stuk denkwerk. Om te oefenen, probeerden wij werkwijzen uit de documentatie toe te passen op vrij beschikbare datasets op Kaggle. Dit om te wennen aan de syntax en structuur. Enkele zaken uit de lessen kwamen hier aan bod. Zo moesten we twee *UDF's* implementeren om twee kolommen te berekenen, waaronder een kolom met de grootcirkelafstand. Op basis van *groupBy's* kunnen we regelmatig data-analyse uitvoeren. Deze werkwijze hielp ons om outliers terug te vinden.

De resultaten van ons model vielen in lijn met andere *notebooks* die deze technieken gebruikten. Notebooks dat regressiemodellen gebruikten, zoals *XGBoost* wordt momenteel niet door Spark aangeboden. Daarnaast beschikten wij niet over een *built-in* methode om de RMSLE te berekenen. Wij losten dit tijdelijk op door de log van de RMSE te berekenen. Na zoekwerk leidden wij af dat het logaritme van de RMSE berekenen een deelse oplossing was. *Sci-kit Learn* biedt hiervoor een *built-in* function aan.

Hoofdstuk 3

Kredietkaartfraude bestrijden met binaire classificatie.

Probleemstelling

Tijdens het zoeken van een classificatieprobleem stuiten wij ons op deze dataset. Deze dataset bestaat uit circa 25.000 rijen. We hebben hier twee klassen: echte en fraudulente kredietkaarttransacties. We pakken dit probleem aan met een binair classificatiemodel. De dataset kan u hier¹ terugvinden. Voor deze oefening baseerden wij ons vooral op de nodige officiële documentatie om een classificatieprobleem in Spark aan te pakken. Wij namen inspiratie uit het notebook dat u hier² kan terugvinden.

Data ophalen

Voor deze opdracht krijgen we één dataset, namelijk *creditcard.csv*. We hebben hier meer dan 30 features. Daarom kiezen we ervoor om geen apart schema aan te maken, maar we gebruiken de *inferSchema* optie van Spark. Zo achterhaalt Spark zelf het datatype van een kolom. Spark herkent hier alle nodige features als een *double*.

```
1 private static Dataset<Row> getData() {  
2     return spark.read()  
3         .option("header", true)  
4         .option("inferSchema", true)  
5         .csv("src/main/resources/creditcard.csv");  
}
```

Data Cleanen

De oorspronkelijke dataset is *skewed*. Bij het bestuderen van de dataset merken we op dat 99.83% van alle transacties echt zijn, terwijl enkel 0.71% wél fraudulent is. Dit pakken we aan met een *clean*-methode. Zonder deze methode hebben we een grotere kans op *overfitting*, want het model neemt aan dat de meeste transacties echt zullen zijn. We lossen dit op door een *sub-sample* te maken van de volledige dataframe. We voorzien een 50:50 ratio tussen echte en fraudulente transacties.

```
1 private static Dataset<Row> createSubSample(Dataset<Row> dataframe) {  
2     dataframe = dataframe.sample(1.0);  
3     dataframe = dataframe.orderBy(rand());  
4     Dataset<Row> nietFraudulent = dataframe.where(col("label").equalTo(1));  
5     long lengte = nietFraudulent.count();  
6     Dataset<Row> fraudulent = dataframe.where(col("label").equalTo(0)).limit((int) lengte);  
7     dataframe = nietFraudulent.union(fraudulent);  
8     return dataframe.sample(1.0);  
}
```

¹Link naar competitie: <http://bit.ly/3iM6zIB>

²Link naar gevolgd notebook: <https://bit.ly/3Y9Kmo6>

Pipeline

Alle nodige features starten met een 'V'-teken. Met behulp van een *for-loop* en de *startswith* methode kunnen wij de features zo aan een array van Strings toevoegen. Het alternatief is 28 features manueel ingeven bij het opstellen van de assembler. De *assembler* zal de 28 kolommen omzetten naar één vector van features.

```
1 String[] arrFeatures = new String[28];
2 int teller = 0;
3 for (String kolom : creditCardSubSample.columns()) {
4     if (kolom.startsWith("V")) {
5         arrFeatures[teller] = kolom;
6         teller++;
7     }
8 }
```

Het model bevat 28 features. Zelf weten we niet welke er noodzakelijk zijn. Als experiment, en toepassing van een techniek uit de lessen *Business Intelligence*, passen wij PCA toe aan de pipeline. Zo voorzien we *dimensionality reduction* op onze features. Zo behouden we enkel de features nodig voor het model. Voor de binaire classificatie keken wij naar drie verschillende classificatiemodellen: *Random Forest*, *Lineaire SVM* en *logistische regressie*. Hyperparameters tunen doen we met een *ParamGridBuilder*.

```
1 VectorAssembler assembler = new VectorAssembler()
2     .setInputCols(arrFeatures)
3     .setOutputCol("features");
4
5 MinMaxScaler minmax = new MinMaxScaler()
6     .setMax(1.0).setMin(0.0)
7     .setInputCol(assembler.getOutputCol())
8     .setOutputCol("scaledFeatures");
9
10 PCA pca = new PCA().setInputCol(minmax.getOutputCol()).setOutputCol("pcaFeatures").setK(3);
11
12 RandomForestClassifier rfc = new RandomForestClassifier()
13     .setLabelCol("label")
14     .setFeaturesCol(pca.getOutputCol());
15
16 ParamMap[] paramGridRFC = new ParamGridBuilder()
17     .addGrid(rfc.maxDepth(), new int[] { 3, 7, 9 })
18     .addGrid(rfc.numTrees(), new int[] { 40, 60, 80 })
19     .addGrid(pca.k(), new int[] { 3, 6, 9 })
20     .build();
21
22 CrossValidator cvRFC = new CrossValidator().setEstimator(rfc)
23     .setEvaluator(new BinaryClassificationEvaluator().setLabelCol("label"))
24     .setEstimatorParamMaps(paramGridRFC)
25     .setNumFolds(3);
26
27 Pipeline pipelineRFC = new Pipeline().setStages(new PipelineStage[]
28     { assembler, minmax, pca, cvRFC });
29
30 PipelineModel pipelineModelRFC = pipelineRFC.fit(trainSetSplit);
31
32 System.out.printf("Ideale parameters:\nMax depth: %s\nNum trees: %s\nk: %s",
33     cvRFC.getEstimatorParamMaps()[0].get(rfc.maxDepth()).toString(),
34     cvRFC.getEstimatorParamMaps()[1].get(rfc.numTrees()).toString(),
35     cvRFC.getEstimatorParamMaps()[2].get(pca.k()).toString());
36
37
38 Dataset<Row> predictionsRFC = pipelineModelRFC.transform(trainTestSplit);
39 createSubmission(predictionsRFC, "submissionRFC");
```

Evaluatie

We hebben een confusionmatrix nodig om het classificatiemodel te evalueren. Zo zien wij hoe goed het model de klassen kan voorspellen. Met de confusion matrix kunnen wij ook de precision en recall berekenen. Deze functies zijn ingebouwd in Spark, dus dit hoeven wij niet manueel te berekenen. Als laatste berekenen we de nauwkeurigheid van het model. Hieronder vindt u de twee methoden om de *confusion matrix* weer te geven, en ook de methode om de nauwkeurigheid van het model op te halen.

```
1 private static double getAreaROCCurve(Dataset<Row> dataframe) {
2     return new BinaryClassificationEvaluator()
3         .setLabelCol(label)
4         .evaluate(dataframe);
5 }
6
7 private static void printCorrelation(Dataset<Row> dataframe) {
8     Row r1 = Correlation.corr(dataframe, "pcaFeatures").head();
9     System.out.printf("\n\nCorrelation Matrix\n");
10    Matrix matrix = r1.getAs(0);
11    for (int i = 0; i < matrix.numRows(); i++) {
12        for (int j = 0; j < matrix.numCols(); j++) {
13            System.out.printf("%.2f \t", matrix.apply(i, j));
14        }
15        System.out.println();
16    }
17 }
18
19 private static void printConfusionMatrixMetrics(Dataset<Row> dataframe) {
20     dataframe = dataframe
21         .select(prediction, label)
22         .orderBy(prediction)
23         .withColumn(label, col(label).cast("double"));
24
25     MulticlassMetrics metrics = new MulticlassMetrics(dataframe);
26     System.out.printf("Precision: %.5f \n", metrics.weightedPrecision());
27     System.out.printf("Recall: %.5f \n", metrics.weightedRecall());
28     System.out.printf("Recall: %.5f \n", metrics.accuracy());
29 }
```

Met een *group-by* kunnen we een *confusion matrix* opstellen. In de *main*-methode spreken we de functies als volgt aan. Hieronder tonen we de evaluatie bij het logistische regressiemodel. De code is gelijkaardig aan de andere modellen.

```
1 System.out.printf("Ideale parameters:\nThreshold: %s\nMax Iterations: %s\nk: %s",
2     cvLogReg.getEstimatorParamMaps()[0].get(lr.threshold()).toString(),
3     cvLogReg.getEstimatorParamMaps()[1].get(lr.maxIter()).toString(),
4     cvLogReg.getEstimatorParamMaps()[2].get(pca.k()).toString()
5 ); // beste hyperparameters opvragen
6
7 System.out.printf("\n\nLogistische Regressie\n");
8 printConfusionMatrixMetrics(predictionsLogReg);
9 System.out.printf("Area ROC curve: %.4f\n", getAreaROCCurve(predictionsLogReg));
10 predictionsRFC.groupBy(col(label), col(prediction)).count().show(); // confusion matrix
```

Bevindingen

Wij merken op dat zowel de Random Forest als de Lineaire SVM even sterk scoort. Het logistische regressiemodel daarentegen komt net iets te kort. Zo zijn er véél meer *False Positives* vergeleken met de False Negatives. Dit wilt zeggen dat het model veel sneller een transactie goed zal keuren, terwijl het eigenlijk een fraudulente transactie is.

Het logistische regressiemodel is optimaal op een threshold van 75%, hoogstens vijf iteraties en een k van 3. Het lineaire SVM-model is optimaal op hoogstens vijf iteraties met een *reg params* van 0.2 en een k van 3. Het RFC-model is optimaal op een maximale diepte van 3, 60 bomen en een k van 3. Onze eigen inbreng, PCA, heeft een kleine invloed gelaten op het model. Zo is de *precision* met een kleine fractie verhoogd.

Hoofdstuk 4

De inhoud van tweets detecteren met binaire NLP-classificatie.

Probleemstelling

Sociale media is een actueel onderwerp. Dubbelzinnige teksten spelen een grote rol. Voor deze opdracht moeten we, op basis van gegeven Engelstalige Tweets, achterhalen of een Tweet gerelateerd is aan een ramp of niet. De dubbelzinnigheid van een zin is een obstakel voor ons model. Zo kan een tweet zeggen *'look at the sky it was ablaze'* terwijl het woord *ablaze* een andere context heeft. In deze zin is er een spreekwoordelijke betekenis. De tekst is ruw. Zo zijn er woorden en symbolen dat ons model kan hinderen. De link naar de competitie kan u hier¹ terugvinden. Wij baseerden ons vooral op de documentatie.

Data ophalen

Deze oefening verschilt niet in aanpak. Wij krijgen de training- en testset mee in CSV-formaat. Dit lezen wij in met het Spark-object. Net zoals bij de regressie-oefening bouwen wij een schema op.

```
1 private static Dataset<Row> getTraining() {
2     List<StructField> fields = Arrays.asList(
3         DataTypes.createStructField("id", DataTypes.IntegerType, false),
4         DataTypes.createStructField("keyword", DataTypes.StringType, false),
5         DataTypes.createStructField("location", DataTypes.StringType, false),
6         DataTypes.createStructField("text", DataTypes.StringType, false),
7         DataTypes.createStructField("target", DataTypes.DoubleType, false)
8     );
9     StructType schemaData = DataTypes.createStructType(fields);
10    return spark.read()
11        .option("header", true)
12        .schema(schemaData)
13        .csv("src/main/resources/train.csv");
14 }
15
16 private static Dataset<Row> getTest() {
17     List<StructField> fields = Arrays.asList(
18         DataTypes.createStructField("id", DataTypes.IntegerType, false),
19         DataTypes.createStructField("keyword", DataTypes.StringType, false),
20         DataTypes.createStructField("location", DataTypes.StringType, false),
21         DataTypes.createStructField("text", DataTypes.StringType, false)
22     );
23     StructType schemaData = DataTypes.createStructType(fields);
24    return spark.read()
25        .option("header", true)
26        .schema(schemaData)
27        .csv("src/main/resources/test.csv");
28 }
```

¹Link naar competitie: <https://www.kaggle.com/c/nlp-getting-started>

Data Cleaning

Voor we beginnen met teksttransformaties moeten we eerst de *dataset* cleanen.

1. We droppen alle rijen die *null*-waarden bevatten.
2. We voegen een *string-only* kolom toe aan het dataframe. We behouden enkel de alfabetische karakters. We gaan er van uit dat we enkel met het Latijns alfabet te maken hebben. Arabische en Kanji-symbolen worden bijvoorbeeld weggefilterd.
3. De data is licht *skewed*. We willen een 50:50 verhouding hebben. Dit passen we aan met de *createSubSample* methode. Deze is gelijkaardig aan de gelijknamige methode uit de vorige oefening.
4. We splitsen de trainingset in twee: een *trainSetSplit* en een *trainTestSplit*. Ter controle printen we de lengte uit van beide dataframes.

```
1 train = train.select(col("id"), col("text"), col("label"));
2 train = train.na().drop();
3 train = train.withColumn("text", regexp_replace(col("text"), "\\d+", ""));
4
5 test = test.select(col("id"), col("text"));
6 test = test.na().drop();
7 test = test.withColumn("text", regexp_replace(col("text"), "\\d+", ""));
8
9 train = createSubSample(train);
10 train.groupBy("label").count().show();
11
12 Dataset<Row>[] datasets = train.randomSplit(verhouding, 42);
13 Dataset<Row> trainSetSplit = datasets[0];
14 Dataset<Row> trainTestSplit = datasets[1];
15
16 System.out.printf("Lengte trainsetsplit: %s\n", trainSetSplit.count());
17 System.out.printf("Lengte testsetsplit: %s\n", trainTestSplit.count());
```

Pipeline

We krijgen rauwe tekstdata binnen. Deze bevat niet-alfanumerieke karakters, overbodige spaties en stopwoorden. Deze woorden willen wij mijden in ons model. Om de tekstdata om te zetten naar bruikbare data voor ons model maken wij gebruik van een pipeline. Op het einde van de pipeline komt het classificatiemodel aan bod.

1. We behouden enkel de alfanumerieke karakters. Alles met spaties en symbolen verwijderen we. Hiervoor gebruiken we een *RegexTokenizer* object.
2. Vervolgens willen we de stopwoorden verwijderen. Dit doen we met een *StopWordsRemover*. Als de tweets in een andere taal waren, bijvoorbeeld Nederlands, dan hadden wij eerst de stopwoorden uit die taal moeten opladen. Vervolgens moeten wij die set van stopwoorden koppelen aan het model. Dit hoeven wij niet te doen.
3. De gefilterde woorden moeten we, net zoals aparte featurekolommen bij het regressiemodel, gaan omzetten naar één featurekolom. Dit zal een vector zijn van alle woorden. In dit verslag gebruiken wij twee verschillende transformers: *CountVectorizer* of *HashingTF*. Hieronder ziet u een voorbeeld van een pipeline met *HashingTF*.
4. Het einde van de rit: een classificatiemodel trainen en testen. Bij deze oefening testen we twee verschillende modellen uit: logistische regressie én een *Random Forest Classifier*. Voor beide classificatiemodellen werken we met crossvalidatie.

Dit is de eerste pipeline. Hier werken we met Logistische regressie, inclusief crossvalidatie. De features worden naar een vector omgezet met *HashingTF*.

```
1 RegexTokenizer tokenizer = new RegexTokenizer()
2   .setInputCol("text")
3   .setOutputCol("words")
4   .setPattern("\\W");
5
6 StopWordsRemover remover = new StopWordsRemover()
7   .setInputCol(tokenizer.getOutputCol())
8   .setOutputCol("filtered");
9
10 HashingTF htf = new HashingTF()
11   .setInputCol(remover.getOutputCol())
12   .setOutputCol("hashedFeatures");
13
14 LogisticRegression lr = new LogisticRegression()
15   .setFeaturesCol(hashingTF.getOutputCol())
16   .setLabelCol(label)
17   .setMaxIter(10).setRegParam(0.001);
18
19 ParamMap[] paramGridLogReg = new ParamGridBuilder()
20   .addGrid(lr.maxIter(), new int[] { 30, 40, 50 })
21   .addGrid(lr.threshold(), new double[] { 0.8, 0.85 })
22   .build();
23
24 CrossValidator cvLogReg = new CrossValidator()
25   .setEstimator(lr)
26   .setEvaluator(new BinaryClassificationEvaluator().setLabelCol(label))
27   .setEstimatorParamMaps(paramGridLogReg);
28
29 Pipeline pipelineLogReg = new Pipeline()
30   .setStages(new PipelineStage[]
31     { tokenizer, remover, htf, cvLogReg }
32   );
```

Hieronder ziet u ons tweede model: een *RFC*-model dat *CountVectorizer* gebruikt om de features naar één vector om te zetten.

```
1 RandomForestClassifier rfc = new RandomForestClassifier()
2   .setLabelCol(label)
3   .setFeaturesCol(vectorizer.getOutputCol()).setSeed(42);
4
5 CountVectorizer vectorizer = new CountVectorizer()
6   .setInputCol(remover.getOutputCol())
7   .setOutputCol("features")
8   .setVocabSize(10000)
9   .setMinDF(5);
10
11 ParamMap[] paramGridRFC = new ParamGridBuilder()
12   .addGrid(rfc.maxDepth(), new int[] { 5, 7, 10 })
13   .addGrid(rfc.numTrees(), new int[] { 40, 80 })
14   .build();
15
16 CrossValidator cvRFC = new CrossValidator()
17   .setEstimator(rfc)
18   .setEvaluator(new BinaryClassificationEvaluator().setLabelCol(label))
19   .setEstimatorParamMaps(paramGridRFC).setNumFolds(5);
20
21 Pipeline pipelineRFC = new Pipeline()
22   .setStages(new PipelineStage[] { tokenizer, remover, vectorizer, cvRFC });
```

Evaluatie

We evalueren het model op basis van de volgende metrieke:

- Nauwkeurigheid
- De *precision* en *recall*.
- *Area under ROC curve*
- De confusion matrix printen we uit met door een *groupby* uit te voeren op het dataframe met de voorspelde waarden.

Net zoals bij de voorbije twee oefeningen printen we de ideale hyperparameters uit.

```
1 private static double getAreaROCCurve(Dataset<Row> dataframe) {
2     return new BinaryClassificationEvaluator().setLabelCol(label).evaluate(dataframe);
3 }
4
5 private static void printConfusionMatrixMetrics(Dataset<Row> dataframe) {
6     System.out.println();
7     dataframe = dataframe.select(prediction, label).orderBy(prediction)
8     .withColumn("target_d", col(label).cast("double")).drop(col(label));
9
10    MulticlassMetrics metrics = new MulticlassMetrics(dataframe);
11    System.out.printf("Precision: %.5f \n", metrics.weightedPrecision());
12    System.out.printf("Recall: %.5f \n", metrics.weightedRecall());
13    System.out.printf("Nauwkeurigheid: %.5f \n", metrics.accuracy());
14 }
15
16 // in main-methode
17 PipelineModel pipelineModelLogReg = pipelineLogReg.fit(trainSetSplit);
18
19 System.out.printf("Ideale parameters:\nMax Iter: %s\nReg params: %s",
20     cvLogReg.getEstimatorParamMaps()[0].get(lr.maxIter()).toString(),
21     cvLogReg.getEstimatorParamMaps()[1].get(lr.threshold()).toString()
22 );
23
24 Dataset<Row> predictedLogReg = pipelineModelLogReg.transform(trainTestSplit);
25 printConfusionMatrixEssence(predictedLogReg);
26 System.out.printf("Area ROC Curve: %.4f\n", getAreaROCCurve(predictedLogReg));
27
28 predictedLogReg.groupBy(col(label), col(prediction)).count().show();
```

Bevindingen

Het logistische regressiemodel is optimaal op 30 maximale iteraties en een *reg params* van 0.85. De diepte van het RFC-model moet hoogstens 4 zijn met 30 bomen om een optimaal resultaat te bekomen. Het finetunen deden we door *trial & error*. We vergroten of verkleinen de parameters stelselmatig.

Beide modellen scoren hier even sterk. Noch het verschil tussen de *HashingTF* en *Count-Vectorizer*, noch het verschil tussen logistische regressie en *Random Forest Classification* valt hier op. Er is een aanvaardbaar evenwicht tussen de *false negatives* en *false positives*.

Hoofdstuk 5

Conclusie

Algemene bevindingen

In het tweede jaar Data Engineering doken wij voor het eerst in de wereld van *machine learning*. Tijdens de lessen werkten wij exclusief met de *Python-library Scikit-Learn*. Spark MLlib was voor onze groep een nieuwe ervaring. Tijdens het maken van deze oefeningen merkten wij vier algemene punten op, waaronder drie gebreken en één voordeel.

De eerste kwaal is het debuggen. De luxe van de Python notebooks hebben wij hier niet meer. Een klein deel van de code uitvoeren, bijvoorbeeld het herladen van een dataset of het aanpassen van een model, is niet mogelijk. Alles moet vanaf nul worden opgestart.

Als tweede punt ontbreekt Spark de nodige aanschouwelijkheid. Data visualiseren is niet mogelijk op een terminal en dit is nadelig om de ruwe data te analyseren. Tijdens de opdrachten gebruikten we de analytische functies van Spark, bijvoorbeeld het tellen van het aantal per klasse, om te kijken of de data *skewed* is. Wij benaderen dit in de regressieoefening met een *case-when*. Daarnaast werkten wij met het gemiddelde en de standaardafwijking om alles buiten een bereik weg te laten. Alles binnen Spark met Java blijft *terminal*-gebaseerd. Er bestaan alternatieven om grafieken te genereren. Na onderzoek kwamen wij uit op *GraphFrames*. Met dit pakket kan je grafieken creëren binnen Java. Dit is een deelse oplossing, want het is niet even eenvoudig als bij *Matplotlib* met Python.

Toch komen de sterktepunten van Apache ook boven water. Spark en Python Pandas staan op een relatief gelijk niveau. Ze bieden een even krachtige oplossing voor data-analyse en data-manipulatie. Na verder onderzoek blijkt Spark toch een krachtige tool voor *Machine Learning*. Zo is dit pakket ideaal om meerdere statistische berekeningen parallel op gigantische datasets uit te voeren. *Sci-kit Learn* is eerder geschikt voor middelgrote datasets. Dit is iets wat wij als toekomstige *data-engineers* best in het achterhoofd houden.

Conclusie

Met het maken van drie distincte ML-oefeningen hebben wij een beter zicht gekregen op Spark MLlib. Dit pakket is een must wanneer je op de volgende twee vragen positief antwoord:

1. Wil ik zware rekenkrachtige toepassingen draaien op een gigantische hoeveelheid data?
2. Is alle data verspreid over een *distributed system*?

We kunnen besluiten dat Spark MLlib een krachtig en veelzijdig platform is voor data-verwerking en *machine learning*. Het biedt gebruikers de mogelijkheid om grote hoeveelheden gegevens snel en efficiënt te verwerken. Daarnaast biedt het een rijk aanbod aan machine learning-algoritmen en -hulpmiddelen. Hoewel er enkele uitdagingen en beperkingen zijn, het kan een waardevol hulpmiddel zijn wanneer je gegevens op grote schaal wil verwerken en analyseren.