

DSC 102: Systems for Scalable Analytics

Programming Assignment 2

1 Introduction

In this assignment we will conduct feature engineering for the Amazon dataset. The extracted features will be used for your next assignment, where you train a model (or models) to predict user ratings for a product.

We will use Apache Spark on AWS EMR, Amazon's managed service for Hadoop ecosystem applications.

You will spawn an EMR cluster (several EC2 instances with EMR services enabled) on AWS. You will then connect to the master node of the cluster and finish all the developments and tests there. You are not expected to code anything locally.

2 Dev-kit

A dev-kit consisting skeletons and other necessary files has been provided to you along with this document. When you spawn the cluster following the instructions in Section 6, this dev-kit should be prepared on your cluster's master node automatically, you do not need to manually download it.

Within the dev-kit there are several files:

```
assignment2.ipynb -- the task descriptions and a playground for your development
assignment2.py -- the deliverable of this assignment, your final file to submit
-----
log4j-spark.properties --
pa2_main.py --
utilities.py -- above files are necessary for your code to run. Do not modify any of them
```

3 Dataset Description

You are expected to extract features from three tables, their schemas and descriptions are listed below:

1. product
 - |-- asin: string, the product id, e.g., 'B00I8HVV6E'
 - |-- salesRank: map, a map between category and sales rank, e.g., {'Home & Kitchen': 796318}
 - | |-- key: string, category, e.g., 'Home & Kitchen'
 - | |-- value: integer, rank, e.g., 796318
 - |-- categories: array, list of list of categories, e.g., [['Home & Kitchen', 'Artwork']]
 - | |-- element: array, list of categories, e.g., ['Home & Kitchen', 'Artwork']
 - | | |-- element: string, category, e.g., 'Home & Kitchen'
 - |-- title: string, title of product, e.g., 'Intelligent Design Cotton Canvas'
 - |-- price: float, price of product, e.g., 27.9
 - |-- related: map, related information, e.g., {'also_viewed': ['B00I8HWOUK']}
 - | |-- key: string, the attribute name of the information, e.g., 'also_viewed'
 - | |-- value: array, array of product ids, e.g., ['B00I8HWOUK']
 - | | |-- element: string product id , e.g., 'B00I8HWOUK'
2. product_processed
 - |-- asin: string, same as above
 - |-- title: string, title column after imputation, e.g., 'Intelligent Design Cotton Canvas'
 - |-- category: string, category column after extraction, e.g., 'Home & Kitchen'
3. review

```
|-- asin: string, same as above
|-- reviewerID: string, the reviewer id, e.g., 'A1MIP8H7G33SHC'
|-- overall: float, the rating associated with the review, e.g., 5.0
```

The `review` table will be useful for extracting the rating information for each product in Task 1. We will be working primarily with `product` table throughout Task 1-4. `product_processed` is used for Task 5-6.

All the datasets required for this assignment can be found in S3 Link¹. We will be reading from this S3 bucket directly. You do not need to download any of them.

4 Tasks

You will be asked to complete six tasks in total. In each task you will need to implement a function `task.i()`. The function signatures and return types are fixed and provided to you in the dev-kit. Each function will take in several inputs and conduct the desired transformations. At the end of each task, you will be asked to extract several statistical properties (mean, variance, etc.) from the transformed data. You will need to programmatically put these properties in a python dictionary named `res`, the schema of which is also given.

Each of the tasks will be tested in unit. It means each function you write will be tested in isolation from the rest. We will award partial points even if some tasks failed.

You can use any combinations of the Spark APIs available in the environment. However, you can select (by setting a global variable called `INPUT_FORMAT`) one of the three APIs for inputs: `DataFrame`, `RDD`, `Koalas`. Inside your function body, you have the freedom to switch between them.

Important: Task 2, 3, 5, 6 **cannot** be solved solely with `Koalas`. Currently `Koalas` does not support nested types so Task 2, 3 are not doable with it, also it does not have the required ML support for Task 5, 6. You will need to switch to other APIs for these tasks.

4.1 Conventions

These rules apply to all the tasks.

4.1.1 Results format

Each task comes with a pre-defined schema for the output results. The result must be stored as python native dictionary and must contain all the keys and nested structures. You must only use python built-in datatypes. For instance, if your value is of datatype `np.float64()`, you must first cast it into python `float`.

For the following schema:

```
res
| -- single_value: int -- an integer number
| -- list_of_values: list -- a list of values
|   -- element: float -- a float number
```

A desired python code snippet to compose up the dictionary would be similar to:

```
1 ...
2 data = ... # Your transformed data
3 res = {
4     'single_value': None,
5     'list_of_values': [None]
6 } # Skeleton given for the result
7 res['single_value'] = int(data.some_op())
8 res['list_of_values'] = [float(data.some_op()), float(data.some_op()), float(data.some_op())]
9 ...
```

¹<s3://dsc102-pa2-public/dataset>

4.1.2 Dealing with null, None and NaN

The input tables contain `null` (or `None` as in RDD/python, or `NaN` as in Koalas/pandas, we will be using these notations interchangeably) and dangling references. You do not need to deal with dangling reference unless instructed. For `null` values we will follow the common practice in SQL world: unless instructed otherwise, you need to ignore all `nulls` when calculating statistics such as count, mean and variance. Of course, do not ignore `null` when you are explicitly asked to count the number of `null` entries.

4.2 Task1: mean and count of ratings

First you will aggregate and extract some information from the user review table. We want to know for each product, what are the mean rating and the number of ratings it received. Implement a function `task_1` that does the following:

1. For each product ID `asin` in `product_data`, calculate the average rating it received. The ratings are stored in column `overall` of `review_data`.
2. Similarly, put the count of ratings for each product in a new column named `countRating`.
3. You need to conduct the above operations, then extract some statistics out of the generated columns. You need to put the statistics in a python dictionary named `res`. The description and schema of it are as follows:

```
res
| -- count_total: int -- count of total rows of the entire table after your operations
| -- mean_meanRating: float -- mean value of column meanRating
| -- variance_meanRating: float -- variance of meanRating
| -- numNulls_meanRating: int -- count of nulls of meanRating
| -- mean_countRating: float -- mean value of countRating
| -- variance_countRating: float -- variance of countRating
| -- numNulls_countRating: int -- count of nulls of countRating
```

If for a product ID, there is not a single reference in `review`, meaning it was never reviewed, you should put `null` in both `meanRating` and `countRating`.

4.3 Task 2: flatten categories and salesRank

Implement a function `task_2()` to conduct the following operations:

1. For the `product` table, each item in column `categories` contains an array of arrays of hierarchical categories. The schema is `ArrayType(ArrayType(StringType))`. We are only going to use the most general category, which is the first element of the nested array: `array[0][0]`. For each row, put the the first element of `categories` in a new column `category`. If `categories` is `null` or empty, put a `null` in your new column.
2. On the other hand, each entry in column `salesRank` is a key-value pair: `(bestSalesCategory, rank)`. Your task is to flatten it into two columns. Put the key in a new column named `bestSalesCategory` and the value in `bestSalesRank`. Put `null` if the original entry was `null` or empty.
3. The schema of output is as follows:

```
res
| -- count_total: int -- count of total rows of the entire table
| -- mean_bestSalesRank: float -- mean value of bestSalesRank
| -- variance_bestSalesRank: float -- variance of bestSalesRank
| -- numNulls_category: int -- count of nulls of category
| -- countDistinct_category: int -- count of all distinct values of category
| -- numNulls_bestSalesCategory: int -- count of nulls of bestSalesCategory
| -- countDistinct_bestSalesCategory: int -- count of distinct values of bestSalesCategory,
|                                     excluding nulls
```

4.4 Task 3: flatten related

Values of `related` column are maps with four keys/attributes: `also_bought`, `also_viewed`, `bought_together`, and `buy_after_viewing`. Each value of these maps contains an array of product IDs. We call them attribute arrays. You need to calculate the length of the arrays and find out the average prices of the products in these arrays.

The logic for all four attributes are identical. For the sake of simplicity, you are only required to flatten the `also_viewed` attribute. Your task is to implement function `task_3()` that does the following:

1. For each row of `related`, you need to :
 1. Calculate the mean price of all products from the `also_viewed` attribute array. Put it in a new column `meanPriceAlsoViewed`. Remember to ignore the products if they do not match any record in `product`, or if they have `null` in price. Do **not** ignore the product if it has `price=0`
 2. Similarly, put the length of that array in a new column `countAlsoViewed`. You do not need to check if the product IDs in that array are dangling references or not. Put `null` (instead of zero) in the new column, if the attribute array is `null` or empty
2. The schema of output is as follows:

```
res
| -- count_total: int -- number of rows of the entire processed table
| -- mean_meanPriceAlsoViewed: float -- mean value of meanPriceAlsoViewed
| -- variance_meanPriceAlsoViewed: float -- variance of meanPriceAlsoViewed
| -- numNulls_meanPriceAlsoViewed: int -- count of null-value entries of meanPriceAlsoViewed
| -- mean_countAlsoViewed: float -- mean value of countAlsoViewed
| -- variance_countAlsoViewed: float -- variance of countAlsoViewed
| -- numNulls_countAlsoViewed: int -- count of null-value entries of countAlsoViewed
```

4.5 Task 4: data imputation

You may have noticed that there are lots of `nulls` in the table. Now your task is to impute them with meaningful values that can be used for ML.

Since the schema is already flattened, now we only have two datatypes in our table: numerical (including integer and floating numbers) and string. Now you need to impute a numerical column `price`, as well as a string column `title`.

1. Please implement a function `task_4()`. For column `price`, first cast it to float type. Then impute the `nulls` with the mean value of all the rest values. Store the outputs in a new column `meanImputedPrice`.
2. Same as above, but this time impute with the **median** value. Store the outputs in a new column `medianImputedPrice`.
3. As for the string-typed columns, we want to simply impute **nulls and empty strings** with a special string `'unknown'`. Store the outputs in a new column `unknownImputedTitle`.
4. The schema of output is as follows:

```
res
| -- count_total: int -- count of total rows of the entire table after above operations
| -- mean_meanImputedPrice: float or None -- mean
| -- variance_meanImputedPrice: float -- variance
| -- numNulls_meanImputedPrice: int -- count of null-value entries
| -- mean_medianImputedPrice: float or None -- mean
| -- variance_medianImputedPrice: float -- variance
| -- numNulls_medianImputedPrice: int -- count of null-value entries
| -- numUnknowns_unknownImputedTitle: float -- count of 'unknown' value entries
```

4.6 Task 5: embedding title with word2vec

This task assumes the `title` column is already imputed with **unknown**. We have provided the imputed data table `product_processed_data`.

In this task we want to transform `title` into a fixed-length vector via word2vec.

1. You need to implement function `task_5()`. For each row, convert `title` to lowercase, then split it by whitespace (' ') to an array of strings, store the array in a new column `titleArray`
2. Train a word2vec model out of column `titleArray`. Do not try to implement word2vec yourself. Instead, use `M.feature.Word2Vec`. See instructions below.
3. For each of the three words inputed as `<word_0>`, `<word_1>`, and `<word_2>`, use your obtained word2vec model to get the 10 closest synonyms along with similarity scores (cosine similarity of word vectors). `M.feature.Word2Vec` also has built-in method for this task.
4. The schema of output is as follows:

```
res
| -- count_total: int -- count of total rows of the entire table
| -- size_vocabulary: int -- the size of the vocabulary of your word2vec model
| -- word_0_synonyms: list -- synonyms tuples of word_0
|   | -- element: tuple -- tuple of format (synonym, score)
|   |   | -- element: string -- synonym
|   |   | -- element: float -- score
| -- word_1_synonyms: list
|   | -- element: tuple
|   |   | -- element: string
|   |   | -- element: float
| -- word_2_synonyms: list
|   | -- element: tuple
|   |   | -- element: string
|   |   | -- element: float
```

word2vec instructions:

1. Set `minCount`, the minimum number of times a token must appear to be included in the word2vec model's vocabulary, to 100.
2. Set the dimension of output word embedding to 16.
3. You need to set the random seed as `SEED`, this is a global variable defined to be 102.
4. Set `numPartitions` to 4.
5. You should keep all other settings as default.
6. `M.feature.Word2Vec` is not fully reproducible (although we have set the seed here). We are aware of the issue and your score will not be affected by its internal randomness.

Reference for word2vecs on Spark ML², on Spark MLlib³.

4.7 Task 6: one-hot encoding category and PCA

Assume `categories` is already flattened and **unknown** imputed for the input data. We have provided you with the preprocessed table

Now you need to one-hot encode the categorical features. Meanwhile, they may be correlated. So as a practice, we would like to run PCA on these categories.

²<https://spark.apache.org/docs/latest/ml-features.html#word2vecs>

³<https://spark.apache.org/docs/latest/mllib-feature-extraction.html>

1. Implement function `task_6()`. First one-hot encode `category` and put the resulted vectors in a new column `categoryOneHot`. Ensure the dimension of generated vectors equals to the size of domain. For example, if we have three categories in total: $V = \{\text{'Electronics'}, \text{'Books'}, \text{'Appliances'}\}$. Then the encoding for 'Electronics' can be $[1, 0, 0]$ or $[0, 1, 0]$ or $[0, 0, 1]$, but the dimension of this vector must be 3.

Hint: For DataFrame, before encoding a string-typed column, you may have to first convert it to a column of numerical indices with `M.feature.StringIndexer`. Then use `M.feature.OneHotEncoderEstimator` to do the encoding. Set `dropLast` argument to false.

For RDD, you may need to implement the one-hot-encoding logic yourself. Consider to build the one-hot mapping locally and broadcast it.

2. Apply PCA on the one-hot-encoded column. Reduce the dimension of each one-hot vector to 15, put the transformed vectors in a new column `categoryPCA`. On DataFrame, use `M.feature.PCA`. On RDD, see instructions⁴.
3. Column `categoryOneHot` and `categoryPCA` will be of `VectorType`. You do not need to worry if the vectors are sparsely or densely represented.
4. The schema of output is as follows::

```
res
| -- count_total: int -- count of total rows of the entire transformed table
| -- meanVector_categoryOneHot: list -- mean vector of transformed one-hot-encoding vectors
|   | -- element: float -- element of the mean vector, from first to last dimension
| -- meanVector_categoryPCA: list -- mean vector of the PCA-transformed vectors
|   | -- element: float
```

5 Deliverables

Code up all the tasks in the designated places in `assignment2.py`. Then rename the file to `assignment2-<your team id>.py`. For instance, if your team id is 18, then your filename would be `assignment2_18.py`. Submit this file on Canvas, only one team member needs to do so.

6 Getting started

6.1 Prerequisite

You need to download Docker on your own computer and prepare the key pairs on AWS.

6.1.1 Docker

The only dependency of this assignment is docker, which can be found in Link⁵. A container image named `yuhzhang/dsc102-pa2` will be used for all AWS-related operations.

Windows users: Your system may not fulfill the requirements for the newest Docker Desktop. In that case, please install Docker Toolbox⁶ instead.

Within the container you are provided with the following utilities

```
1 s3-init
2 emr-launch
3 emr-list
4 emr-dns
5 emr-terminate
```

The entry point of these utilities is

```
1 docker run --env-file <path> yuhzhang/dsc102-pa2 <utility>
```

⁴<https://spark.apache.org/docs/latest/mllib-dimensionality-reduction>

⁵<https://docs.docker.com/install>

⁶https://docs.docker.com/toolbox/toolbox_install_windows

<path> is the path to your credential file, see below for instructions. <utility> is any utility listed above. The AWS CLI tool is also available as `aws2`. For example, to list your s3 bucket, you can run:

```
1 docker run --env-file <path> yuhzhang/dsc102-pa2 aws2 s3 ls
```

6.1.2 Prepare key pairs on AWS console

If you have not done so. Go to your AWS console⁷ and click key pairs. Follow the instructions to create or upload a public key to AWS. Note down the name of the key as displayed on your AWS console.

6.2 Store AWS credentials

Create a new text file named `credentials.list` under any directory and put the following content in it:

```
1 PID=...
2 AWS_ACCESS_KEY_ID=...
3 AWS_SECRET_ACCESS_KEY=...
4 AWS_SESSION_TOKEN=...
```

Fill in the blanks (placeholderd by ...) with corresponding values. PID is your UCSD pid (e.g., a13230999, with 'a' in lower case). You can find the rest three AWS credentials from Link⁸. Note these credentials are only temporary. You will need to update this file if the token expires. To obtain a new token, simple revisit the link.

6.3 Initialize assignment-related S3 buckets

Use the following command to initialize the S3 buckets needed for this assignment:

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 s3-init
```

This command will setup a EMR log bucket named <your pid>-emr-logs, and a bucket storing your scripts named <your pid>-pa2. It will also copy the assignment dev-kit to the latter bucket.

6.4 Launch an EMR cluster

1. Use `emr-launch` utility built in the docker image `yuhzhang/dsc102-pa2`:

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-launch -k <key name>
```

The above utility will spawn a EMR cluster named <pid>-emr-cluster. Available flags/arguments are:

- k <key name>, the name of your secrets to access the cluster via ssh
- b, optional, if set, your CORE instances will be run with Spot pricing
- n <number of workers>, optional, the number of workers you want to have, default: 4
- c, optional, if set, use a cheaper hardware m4.large, otherwise use m5.xlarge
- t, optional, if set, setup Theia IDE on master node port 3000
- f, optional, if set, force starting the second cluster, as by default only one cluster is permitted

2. Example usage

Spawn a cluster with 4 workers and Spot pricing:

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-launch -k <key name>
  -b
```

This environment will be used to evaluate all of your submissions (except we will not use spot instances).

⁷<<https://ets-apps.ucsd.edu/dsc102-custom-aws/>>

⁸<<https://ets-apps.ucsd.edu/dsc102-custom-aws/?mode=env>>

6.5 Access your cluster and the assignment

1. Use the following command to list the cluster IDs:

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-list
```

2. Then query the DNS name of cluster's master node using:

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-dns <cluster ID>
```

This will return the DNS name such as:

```
ec2-###-##-##-###.compute-1.amazonaws.com
```

3. (Optional) You can now try to SSH into your master node via:

```
1 ssh -i path/to/key hadoop@ec2-###-##-##-###.compute-1.amazonaws.com
```

Log out the ssh session before proceeding to the next steps.

4. In your browser, access

```
ec2-###-##-##-###.compute-1.amazonaws.com:8888
```

this will direct you to the jupyter notebook running on the master node. The password would be your pid. The working directory of this jupyter notebook is S3 bucket <your pid>-pa2 mounted. So all your modifications to the assignment files will be reflected to that bucket.

5. (Optional) and if you used -t flag when launching the cluster, you can access Theia IDE running on

```
ec2-###-##-##-###.compute-1.amazonaws.com:3000
```

6. In Jupyter notebook, rename `assignment2.ipynb` to `assignment2-<your pid>.ipynb` and continue the assignment by following the instructions written in the notebook.

6.6 Testing and submission

You will **not** submit the notebook. Instead, you need to put your implementations of `task_1` to `task_6`, along with all the dependencies you imported and helper functions you defined, in the file co-located with the notebook: `assignment2.py`.

If you are collaborating in team, please combine your work into one single file. Only **one** person needs to submit the final file. Do **not** modify the filename yet.

6.6.1 Test your file

Before submitting the file, you need to make sure your script runs under the deployment environment, otherwise you may lose points.

1. If your cluster is not in deployment mode, meaning if you don't have 4 m5.xlarge workers, terminate your current cluster (see instructions below), then launch a deployment cluster via

```
1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-launch -k <key name> -b
```

2. SSH into the master node by

```
1 ssh -i path/to/key hadoop@ec2-###-##-##-###.compute-1.amazonaws.com
```

3. Go to your root directory of scripts

```
1 cd /mnt/<your pid>-pa2/src
```

4. Run PA2 with the following command, do not modify anything except <your pid>:


```

1 spark-submit \
2 --py-files utilities.py,assignment2.py \
3 --files log4j-spark.properties \
4 --master yarn \
5 --deploy-mode client \
6 --conf spark.memory.fraction=0.8 \
7 --conf spark.dynamicAllocation.enabled=false \
8 --conf spark.sql.crossJoin.enabled=true \
9 --driver-java-options "-Dlog4j.configuration=file:log4j-spark.properties" \
10 --conf "spark.executor.extraJavaOptions=-Dlog4j.configuration=file:log4j-spark.properties"
11 \
    pa2_main.py --pid <your pid>

```

Make sure your script can execute and try to pass as many tests as you can.

6.6.2 Submit your file

Go to your AWS console, navigate to S3 buckets and find the bucket named <your pid>-pa2. Download the assignment2.py file to your own machine.

Then rename the file to assignment2-<your team id>.py. For instance, if your team id is 18, then your filename would be assignment2_18.py.

Upload this file to canvas, only one of the team members needs to do so.

6.7 Terminate your cluster

Don't forget to terminate the cluster when you are done:

```

1 docker run --env-file path/to/credentials.list yuhzhang/dsc102-pa2 emr-terminate <cluster id>

```