

Name: Nishit Shaileshbhai Rathod
Student Number: N01586439

Assignment

- Go to <https://www.openml.org/>
- Select a data set
- Analyze the data in DW and find some opportunities to improve the quality of data (like removing missing value or scaling data or one hot encoding, etc)
- Clean/transform data in Studio DW
- Upload the following items in BB
 - Your selected data set
 - DW .flow file (to learn how to download the .flow file see next slide)
 - A report that has covers at least the suggested ToC (please see next slides)

ToC for the report

- Data set fields descriptions
- The visualization and analysis that you have done in DW and what you learn out of those visualizations. Each analysis comes with a picture and your description beneath
- The transformations you have done, again you have to explain why you have chosen that transformation, include the picture and explain the results after transformation (with pictures)
- If you have added a code, include that in the report in the right spot

Step 1: Selected the dataset from : OpenML

Step 2: Knowing the dataset.

The dataset is called “basketball” which has 96 instances and 5 features. The dataset was originally taken from –

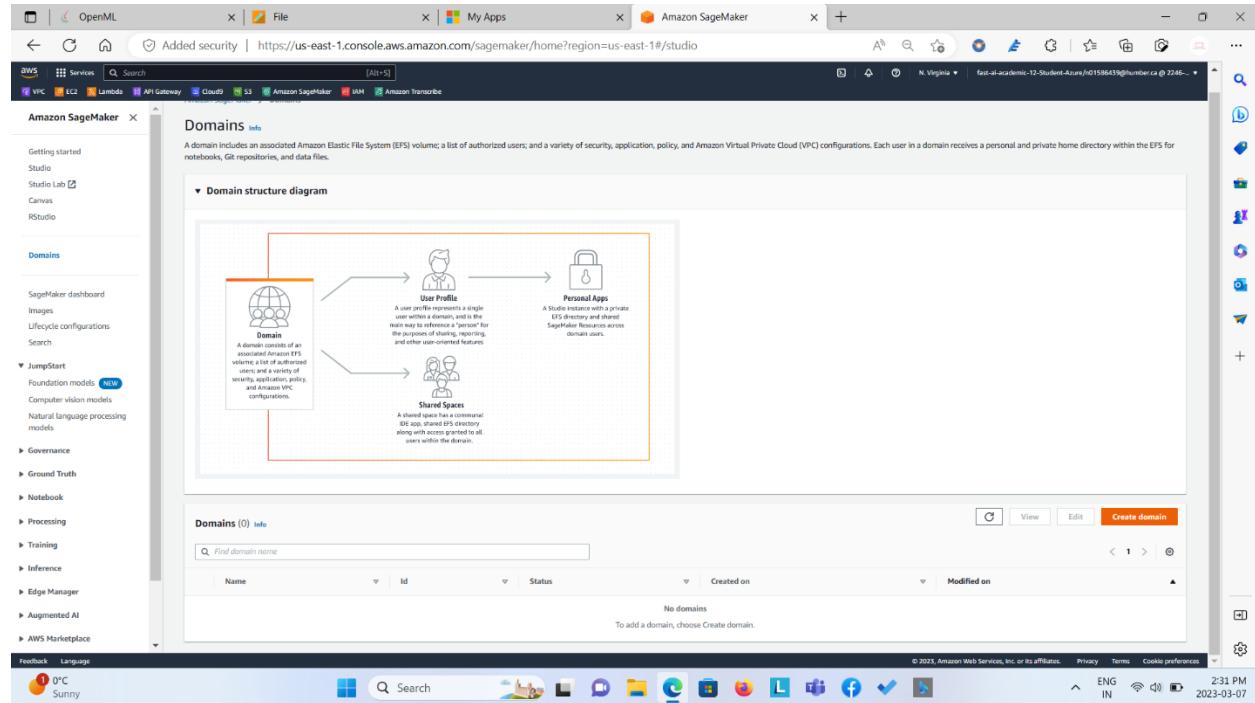
- Dataset from Smoothing Methods in Statistics (<ftp://stat.cmu.edu/datasets>)
- Simonoff, J.S. (1996). Smoothing Methods in Statistics. New York: Springer-Verlag.

Data Features:

- 1. points_per_minute (Target Variable) – Numeric type - 95 distinct values and 0 missing attributes.**
- 2. assists_per_minute – Numeric Type - 96 distinct values and 0 missing attributes.**
- 3. height – Numeric Type - 13 distinct values and 0 missing attributes.**
- 4. time_played – Numeric Type - 94 distinct values and 0 missing attributes.**
- 5. age – Numeric Type - 15 distinct values and 0 missing attributes.**

As the dataset didn't have Categorical feature as well as any missing values, I deliberately transformed the dataset and added "Sex" column as well as removed some of the values from the features. By doing this I can know have the hands-on practice in data pre-processing as well as in exploratory data analysis using AWS Sage maker Data Wrangler.

Step 3: Making the Domain, Creating the user in order to get excess to Studio.



The screenshot shows the 'Amazon SageMaker' service in the AWS console. The left sidebar includes sections like 'Getting started', 'Studio', 'Domains', 'JumpStart', 'Governance', and 'AWS Marketplace'. The main content area is titled 'Setup SageMaker Domain' and provides two options: 'Quick setup (1 min)' and 'Standard setup (10 min)'. The 'Standard setup' is described as controlling all aspects of account configuration, including permissions, integrations, and encryption. A note states: 'Perfect for single user domains and first time users looking to get started with SageMaker.' A 'Configure' button is at the bottom. The browser status bar shows the URL: <https://us-east-1.console.aws.amazon.com/sagemaker/home?region=us-east-1#studio/create-domain/standard-setup>.

This screenshot shows the 'General settings' step of the domain setup wizard. On the left, a sidebar lists steps: Step 1 General settings (selected), Step 2 Studio settings, Step 3 RStudio settings, and Step 4 Canvas settings. The main area has three tabs: 'Domain name' (with a 'Name' field containing 'baseball'), 'Authentication' (with a radio button selected for 'AWS Identity and Access Management (IAM)'), and 'Permission'. Under 'Permission', there are two sections: 'Default execution role' (with a dropdown for 'Enter a custom IAM role ARN' containing 'arn:aws:iam::224670572127:role/fast-ai-academic-12-Student-Azure') and 'Space default execution role' (with a dropdown for 'Enter a custom IAM role ARN' containing 'arn:aws:iam::224670572127:role/fast-ai-academic-12-Student-Azure'). The browser status bar shows the URL: <https://us-east-1.console.aws.amazon.com/sagemaker/home?region=us-east-1#studio/create-domain/standard-setup>.

Custom IAM role ARN
amawsiam:224670572127:role/fast-ai-academic-12-Student-Azure

If you leave space default execution role blank, you will not be able to use the shared spaces feature for collaboration. Once a role has been set for the space default execution role, it can never be deleted or unset. You can only change the default role from the domain settings page.

Network and Storage Section

VPC To enable internet access, make sure that your VPC has a NAT gateway and your security group allows outbound connections.
vpc-0765b259ff1bed1be (172.31.0.0/16) | aws-controltower-VPC

Subnet Choose one or more availability zones supported by Amazon SageMaker:
Choose one or more subnets
PrivateSubnet1A
PrivateSubnet2B

Security groups These security groups will also be associated with the iStudioServerPro App.
Choose one or more security groups
sg-07e9bd22fda815ca (default)

Encryption key - optional SageMaker uses AWS managed CMK to encrypt your EFS and EBS file systems by default. To use a customer managed CMK, enter its key ID or ARN. Learn more

No Custom Encryption

Cancel Next

Step 1 Studio settings
Step 2 Studio settings
Step 3 Default settings
Step 4 Canvas settings

Configure Studio IDE and Notebooks for your organization.

Jupyter Lab version info Default Jupyter Lab version Select the Jupyter Lab version by default for all users in the domain. Permissions to run Jupyter Lab versions are defined by an IAM policy. You must restart the Jupyter Server app to make the version change effective.
Jupyter Lab 5.0

Notebook Sharing Configuration Recommended defaults have been selected for you.

Enable notebook resource sharing Notebook resources include artifacts such as cell output and EBS Repositories. Learn more

SageMaker Projects and JumpStart (optional)

SageMaker Projects and JumpStart Note Create access and provisioning of AWS Service Catalog portfolios in Amazon SageMaker Studio for Amazon SageMaker Projects and Amazon SageMaker JumpStart.

Enable Amazon SageMaker project templates and Amazon SageMaker JumpStart for this account This account Amazon SageMaker Studio can view the Amazon SageMaker built-in project templates and Amazon SageMaker JumpStart notebooks published in AWS Service Catalog & Amazon Connect role and a CloudWatch Metrics role. If you are using a custom role, ensure the role has the necessary permissions to view the notebooks.

Enable Amazon SageMaker project templates and Amazon SageMaker JumpStart for studio users Studio users Studio users can view the domain execution role to create projects using template and template solutions published by Amazon SageMaker in AWS Service Catalog & Amazon Connect role. If you are using a custom role, ensure the role has the necessary permissions to enable them on the user profile page.

Create the roles which are needed to use the latest updated AWS Service catalog products in Amazon SageMaker Studio for Amazon SageMaker Projects and Amazon SageMaker JumpStart. Roles to be created: AWSServiceRole, CloudWatchLogs, CloudBuild, CloudFront, Events, Forecast, Glue, Lambda, Sagemaker, Segment.

Cancel Back Next

The screenshot shows the 'Amazon SageMaker > Setup SageMaker Domain' wizard. The current step is 'Step 1 General settings'. The main panel displays the 'General settings' configuration, which includes shared configurations across the entire SageMaker Domain. A note at the top states: 'Note: A service role for AWS License Manager is needed if you want to configure RStudio on SageMaker. Please refer to AWS License Manager Getting Started documentation to set up this role: <https://docs.aws.amazon.com/sagemaker/latest/dg/rstudio-license.html>'. Below this, the 'RStudio Workbench - optional' section is shown, containing a 'License' field with the message: 'A saved license from AWS License Manager will be automatically detected once a license has been added. If a license is not detected, you must add one to AWS License Manager to activate, and use RStudio Workbench, and other RStudio tools.' A warning message '⚠️ RStudio Workbench license not detected.' is displayed. At the bottom right are 'Cancel', 'Back', and 'Next' buttons.

The screenshot shows the 'Amazon SageMaker Canvas settings' configuration page. The left sidebar lists steps: Step 1 General settings, Step 2 Studio settings, Step 3 RStudio settings, and Step 4 Canvas settings. The current step is 'Step 4 Canvas settings'. The main panel contains three sections: 'Canvas base permissions configuration', 'Time series forecasting configuration', and 'Local file upload configuration'. Under 'Canvas base permissions configuration', there is a note about enabling base permissions for users to build models in Canvas. Under 'Time series forecasting configuration', there is a note about enabling time series forecasting for users to use time series forecasting in Canvas. Under 'Local file upload configuration', there is a note about enabling local file upload for users to upload local files in Canvas. At the bottom right are 'Cancel', 'Back', and 'Submit' buttons.

The screenshot shows the completed 'Amazon SageMaker Canvas settings' configuration page. The left sidebar shows the steps completed: Step 1 General settings, Step 2 Studio settings, Step 3 RStudio settings, and Step 4 Canvas settings. The main panel displays the configuration details for Canvas settings, including the base permissions configuration and the time series forecasting configuration. At the bottom right are 'Cancel', 'Back', and 'Submit' buttons.

The screenshot shows the Amazon SageMaker console with the URL <https://us-east-1.console.aws.amazon.com/sagemaker/home?region=us-east-1#/studio>. The left sidebar is collapsed, and the main content area displays the 'Domains' page. The title bar says 'Amazon SageMaker > Domains'. A sub-header 'Domains Info' provides a brief description of what a domain is. Below it is a 'Domain structure diagram' button. The main table lists one domain: 'baskball' (d-qe6qfc8igww), which is 'inService' and was created on Mar 07, 2023, at 20:10 UTC, last modified on Mar 07, 2023, at 20:13 UTC. Action buttons for 'View', 'Edit', and 'Create domain' are available.

This screenshot shows the 'Domain details' page for the 'baskball' domain. The title bar includes a warning message: '⚠ No new Jupyter Lab 1 version apps can be created from March 30, 2023 onwards, with only Jupyter Lab 3 version app creation being supported. All existing apps running on Jupyter Lab 1 version will be removed on April 30, 2023.' A 'Learn more' button is present. The sub-header is 'Amazon SageMaker > Domains > Domain: baskball'. The 'User profiles' tab is selected, showing a table with one entry: 'No users'. A note below the table says 'To add a user, choose Add user and enter a user name.' Other tabs include 'Space management', 'Environment', and 'Domain settings'. The bottom of the screen shows the standard Windows taskbar with various pinned icons.

Screenshot of the "Add user profile" wizard Step 1: General settings.

General settings
User profile and details.

User profile

Name: default-1678221157457
The name can have up to 63 characters. Valid characters: A-Z, a-z, 0-9, and - (hyphen)

Execution role: fast-ai-academic-12-Student-Azure
The default execution role for both users and spaces in the domain. The execution role must have the [AmazonSageMakerFullAccess](#) policy attached.
Create role using the role creation wizard

Tags - optional
Add tag
You can attach up to 50 tags

Cancel Next

Screenshot of the "Add user profile" wizard Step 2: Studio settings.

Studio settings
Configure Studio IDE and Notebooks for your organization.

Jupyter Lab version Info
Default Jupyter Lab version
The Jupyter Server runs with the selected version by default for the user. Permissions to run Jupyter Lab versions are defined by an [IAM policy](#). You must restart the Jupyter Server app to make the version changes effective.
Jupyter Lab 3.0

SageMaker Projects and JumpStart - optional
Enable access and provisioning of AWS Service Catalog Portfolio of products in Amazon SageMaker Studio for Amazon SageMaker Projects and JumpStart. Learn more [\[?\]](#)

Enable Amazon SageMaker project templates and Amazon SageMaker JumpStart for Studio users
If enabled, this setting allows users who are currently using the domain execution role to create projects using templates and JumpStart solutions published by Amazon SageMaker in AWS Service Catalog. If there are individual users using custom execution roles in your organization, you need to enable them on the user profile page.

Cancel Back Next

OpenML Class Collaborate Class Collaborate - Machine My Apps Amazon SageMaker New tab

Added security https://us-east-1.console.aws.amazon.com/sagemaker/home?region=us-east-1#/studio/d-qe6qfc8lqwv/add-user-pro... N. Virginia fast-ai-academic-12-Student-Azure/n01586439@humber.ca @ 2246...

aws Services Search [Alt+S] VPC EC2 Lambda API Gateway Cloud9 S3 Amazon SageMaker IAM Amazon Transcribe

Amazon SageMaker Domains Domain: basketball Add user profile

Add user profile

Step 1 General settings

Step 2 Studio settings

Step 3 RStudio settings

Step 4 Canvas settings

RStudio settings

Configure RStudio IDE for your organization.

RStudio Workbench - optional

License

A saved license from AWS License Manager will be automatically detected once a license has been added. If a license is not detected, you must add one to AWS License Manager to activate, and use RStudio Workbench, and other RStudio tools.

⚠️ RStudio Workbench license not detected.

Cancel Back Next

Feedback Language 0°C Sunny 3:33 PM 2023-03-07

OpenML Class Collaborate Class Collaborate - Machine My Apps Amazon SageMaker New tab

Added security https://us-east-1.console.aws.amazon.com/sagemaker/home?region=us-east-1#/studio/d-qe6qfc8lqwv/add-user-pro... N. Virginia fast-ai-academic-12-Student-Azure/n01586439@humber.ca @ 2246...

aws Services Search [Alt+S] VPC EC2 Lambda API Gateway Cloud9 S3 Amazon SageMaker IAM Amazon Transcribe

Amazon SageMaker Domains Domain: basketball Add user profile

Add user profile

Step 1 General settings

Step 2 Studio settings

Step 3 RStudio settings

Step 4 Canvas settings

Amazon SageMaker Canvas settings Info

Configure Canvas for your organization.

Canvas base permissions configuration

Enable Canvas base permissions

If you enable Canvas base permissions, your users will have the necessary permissions to build models in Canvas. If you disable Canvas base permissions, your users won't have the necessary permissions to use Canvas, and you must manually configure IAM permissions for full Canvas functionality.

The [AmazonSageMakerCanvasFullAccess](#) policy will be attached to the default Sagemaker execution role that you specified in General settings.

Time series forecasting configuration

Enable time series forecasting

Enable time series forecasting to allow users to use time series forecasting in Canvas.

Your users won't be able to use time series forecasting in Canvas because you have disabled the time series forecasting permission for Canvas.

Cancel Back Submit

The screenshot shows the Amazon SageMaker console with a green success message: "User profile was successfully created." Below it, a warning message states: "No new Jupyter Lab 1 version apps can be created from March 30, 2023 onwards, with only Jupyter Lab 3 version app creation being supported. All existing apps running on Jupyter Lab 1 version will be removed on April 30, 2023." The navigation bar includes tabs for Services, Search, and the current region N. Virginia. The left sidebar lists various services like VPC, EC2, Lambda, API Gateway, Cloud9, S3, Amazon SageMaker, IAM, and Amazon Transcribe.

Step 4: Creating S3 Bucket and Uploading the dataset (.csv file) in it.

The screenshot shows the AWS S3 console with the "Create bucket" wizard. The "General configuration" step is active, showing a "Bucket name" field with "basketballbucket" and an "AWS Region" dropdown set to "US East (N. Virginia) us-east-1". The "Object Ownership" section shows "ACLs disabled (recommended)" selected. The status bar at the bottom indicates "Upcoming permission changes to disable ACLs". The browser toolbar at the top shows the URL as https://s3.console.aws.amazon.com/s3/bucket/create?region=us-east-1.

The screenshot shows the AWS Management Console interface for creating a new S3 bucket. The URL in the address bar is <https://s3.console.aws.amazon.com/s3/bucket/create?region=us-east-1>. The main form includes sections for:

- Tags (0) - optional**: A note about using tags to track storage costs and organize buckets, with a link to learn more.
- Default encryption**: A note that server-side encryption is automatically applied to new objects stored in this bucket.
- Encryption key type**: A radio button group where "Amazon S3 managed keys (SSE-S3)" is selected.
- Bucket Key**: A note about KMS encryption, with "Enable" selected.
- Advanced settings**: A link to view additional configuration options.

At the bottom right of the form is a large orange **Create bucket** button. Below the form is a status bar showing the date and time (3:37 PM, 2023-03-07).

The screenshot shows the AWS Management Console interface for the "baskballbucket" S3 bucket. The URL in the address bar is <https://s3.console.aws.amazon.com/s3/buckets/baskballbucket?region=us-east-1&tab=objects>. The page displays the following:

- Buckets** sidebar: Shows access points, multi-region access points, batch operations, and IAM access analyzer for S3.
- Objects** tab: Shows a table with no objects. It includes columns for Name, Type, Last modified, Size, and Storage class. Buttons for Actions, Create folder, and Upload are available.
- Find objects by prefix** search bar.
- No objects** message: "You don't have any objects in this bucket."
- Upload** button.

At the bottom right of the page is a status bar showing the date and time (3:38 PM, 2023-03-07).

Screenshot of the AWS S3 Management console showing the upload process for a CSV file.

The browser address bar shows: `https://s3.console.aws.amazon.com/s3/upload/basketballbucket?region=us-east-1`

The AWS navigation bar includes: Services, Search, [Alt+S], Global, fast-ai-academic-12-Student-Azure/n01586439@humber.ca @ 2246...

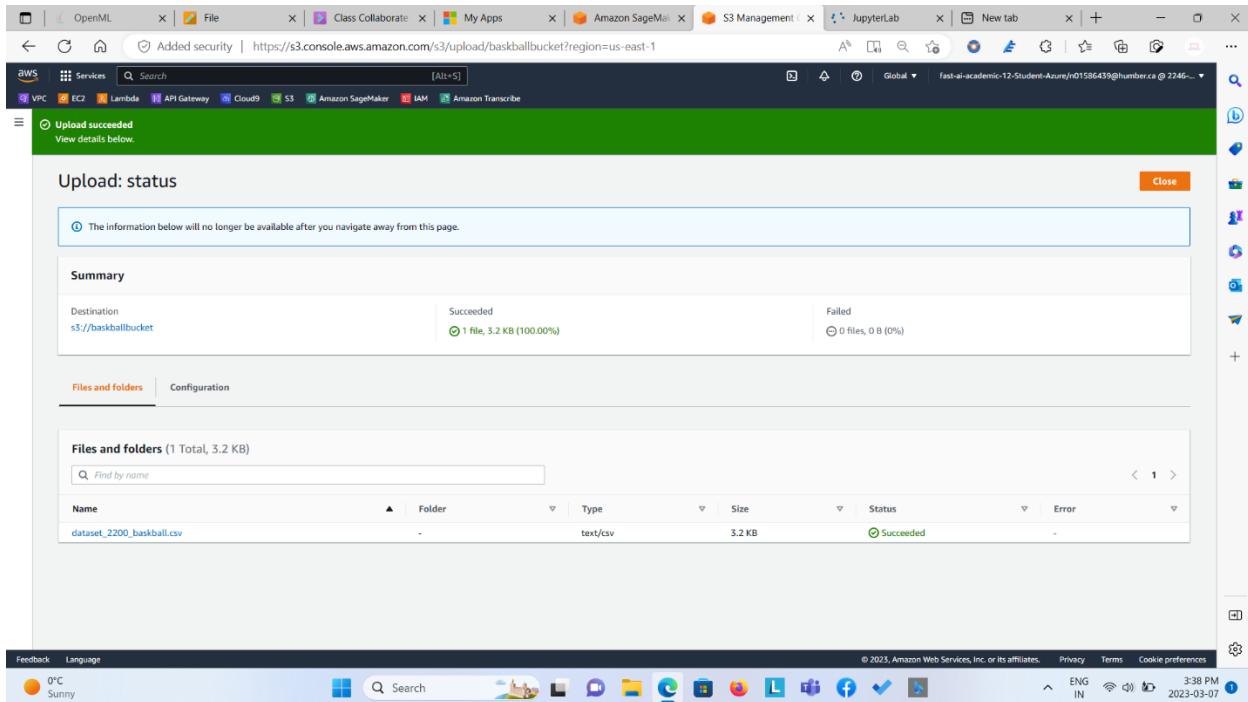
The main interface shows the "Upload" step of the wizard. The "Destination" field is set to `s3://basketballbucket`. A file selection dialog is open, showing the local file `dataset_2200_basketball.csv` selected for upload.

Screenshot of the AWS S3 Management console showing the upload process for a CSV file.

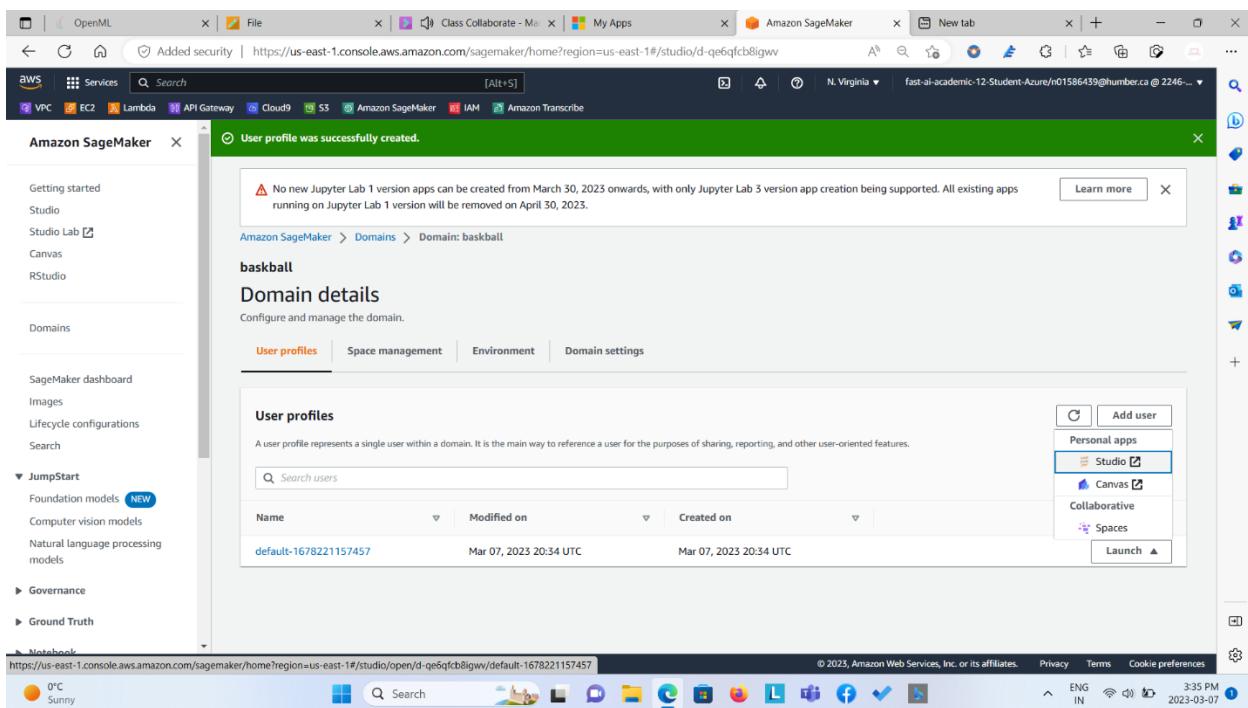
The browser address bar shows: `https://s3.console.aws.amazon.com/s3/upload/basketballbucket?region=us-east-1`

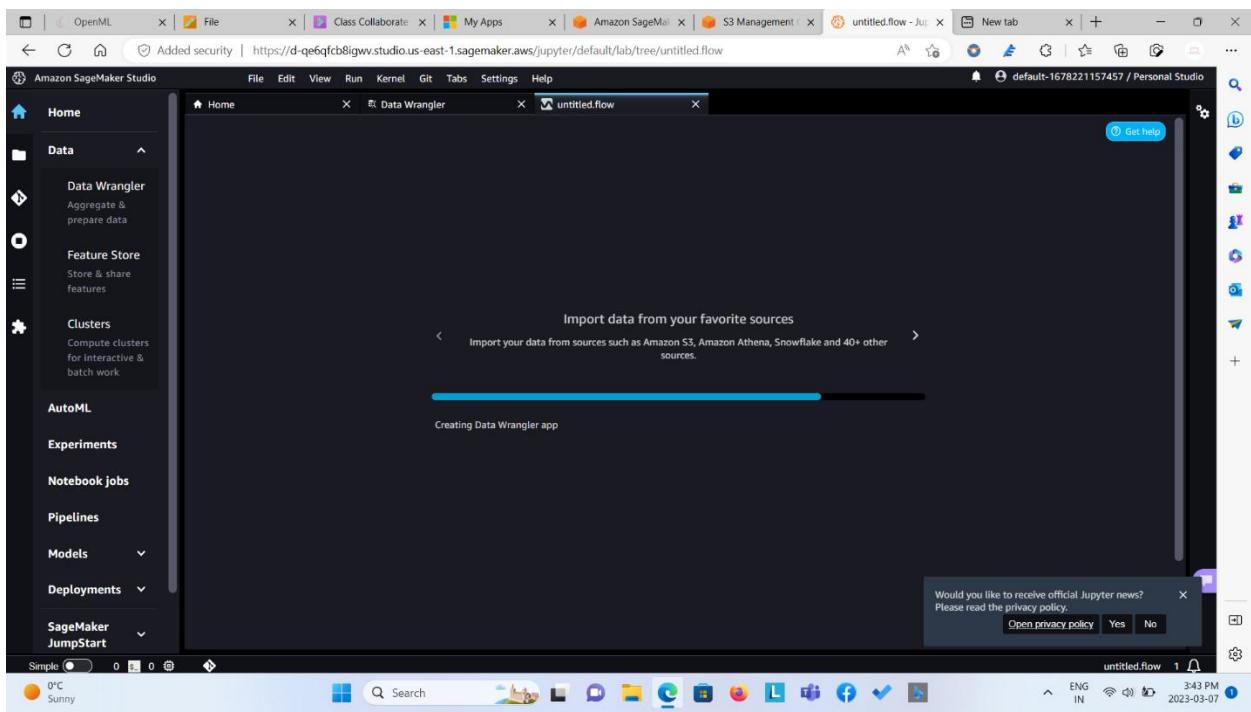
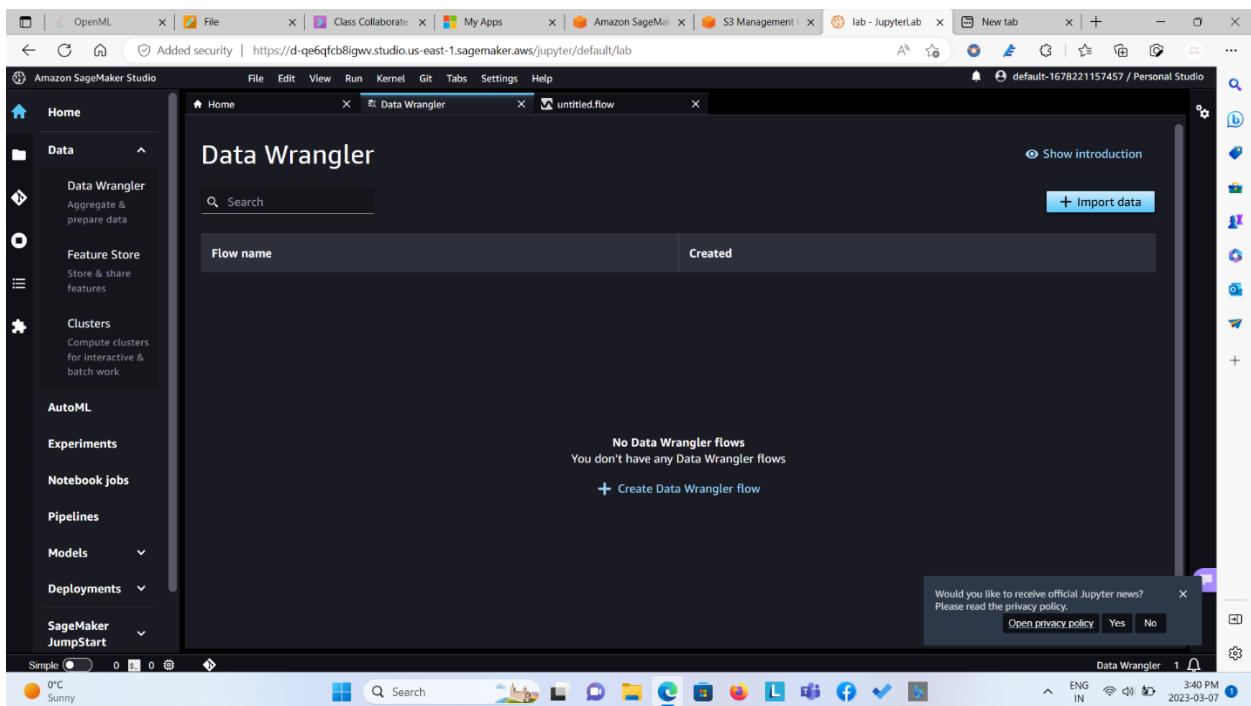
The AWS navigation bar includes: Services, Search, [Alt+S], Global, fast-ai-academic-12-Student-Azure/n01586439@humber.ca @ 2246...

The main interface shows the "Upload" step of the wizard. The "Destination" field is set to `s3://basketballbucket`. The "Files and folders" section shows one item: `dataset_2200_basketball.csv` (text/csv, 3.2 KB). The "Upload" button is visible at the bottom.



Step 5: Performing Data Cleaning/Data Pre-processing/Data Transformation/Data Analysis via. Data Wrangler from Studio by linking Dataset from S3 Bucket.





The screenshot shows the Amazon SageMaker Studio interface with the 'Data Wrangler' tab selected. A modal window titled 'Create connection' is open, prompting the user to select a data source to import a dataset. The 'Available' section lists six options: Amazon S3, Amazon Athena, Amazon Redshift, Snowflake, Amazon EMR, and Databricks. A search bar is available at the top of the modal. On the right side of the modal, there is a link to 'Use sample data'. A small notification at the bottom right of the modal asks if the user wants to receive official Jupyter news, with 'Open privacy policy' and 'Yes' or 'No' buttons.

The screenshot shows the 'Import a dataset from S3' step in the Data Wrangler interface. The left panel displays the 'Import data' section, which includes an 'Advanced configuration' button and a 'S3 URI path' input field containing 'S3 / baskballbucket / dataset_2200_baskball.csv'. Below this, a table shows the object details: 'dataset_2200_baskball.csv', '3.21KB', and 'Last modified 2023-03-07 20:38:29+00:00'. The right panel is the 'DETAILS' panel, which contains configuration settings for the import: 'Name' set to 'dataset_2200_baskball.csv', 'File type' set to 'csv', 'First row is header' checked, 'Delimiter' set to 'COMMA', 'Sampling' set to 'First K', 'Sample size' set to '50000', and 'Filename as separate column' unchecked. A preview of the first 100 rows of the CSV file is shown below the configuration panel.

Object name	Size	Last modified
dataset_2200_baskball.csv	3.21KB	2023-03-07 20:38:29+00:00

PREVIEW - dataset_2200_baskball.csv (First 100 rows shown. The preview doesn't reflect your sampling configuration.)

sex	assists_per_minute	height	time_played	age	points_per_minute
Male	0.0888	201	36.02	28	0.5885
Male	0.1399	198	39.52	30	0.8291
Female	0.0747	198	38.8	26	0.4974
Female	0.0983	191	40.71	30	0.5772
Female	0.1276	196	38.4	28	0.5703

Screenshot of Amazon SageMaker Studio Data Wrangler interface showing a dataset preview and step details.

Data types - Transform: dataset_2200_basketball.csv

Step 2. Data types

sex (string)	assists_per_minute (float)	height (long)	time_played (float)	age (long)	points
Male	0.0888	201	36.02	28	0.58
Male	0.1399	198	39.32	30	0.82
Female	0.0747	198	38.8	26	0.49
Female	0.0983	191	40.71	30	0.57
Female	0.1276	196	38.4	28	0.57
Female	0.1671	201	34.1	31	0.58
Male	0.1906	193	36.2	30	0.52
Female	0.1061	191	36.75	27	0.55
Female	0.2446	185	38.43	29	0.40
Male	0.167	203	33.54	24	0.47
Male	0.2485	188	35.01	27	0.45
Male	0.1227	198	36.67	29	0.49
Female	0.124	185	33.88	24	0.56
Female	0.1461	191	35.59	30	0.51
Female	0.2315	191	38.01	28	0.37
Female	0.0494	193	32.38	32	0.55
Female	0.1107	196	35.22	25	0.47
Male	0.2521	183	31.73	29	0.57

ALL STEPS

- + Add step
- 1. 1.53 Source
- 2. Data types

Screenshot of Amazon SageMaker Studio Data Wrangler interface showing a data flow diagram and a context menu.

Data flow

Choose the plus sign to add a step to the flow. Select a step to edit it.

Get insights from your data

Use built-in analyses to better understand your data. You can use the information to help you with processing.

Validation complete 0 errors

Source - Sampled → **Data types**

Transform: dataset_2200_bask...

Context Menu (right-clicked on the flow)

- Add transform
- Add analysis
- Train model NEW
- Get data insights
- Add destination > Amazon S3
- Export to > SageMaker Feature Store NEW
- Join
- Concatenate
- Edit

The screenshot displays the Amazon SageMaker Studio interface with several open tabs and panels.

Top Bar: Includes standard browser controls (Back, Forward, Stop, Refresh), a search bar, and a tab bar with multiple entries: OpenML, File, Class Collaborate, My Apps, Amazon SageMaker, S3 Management, untitled.flow - lab, and default-1678221157457 / Personal Studio.

Main Area:

- Data Wrangler Tab:** Shows a data table for "dataset_2200_basketball.csv". The table has columns: assists_per_minute, height, time_played, age, and points_per_minute. A tooltip for "Create analysis" is visible over the table.
- Feature Store Tab:** Shows a sidebar with sections: Feature Store, Data Wrangler, Clusters, AutoML, Experiments, Notebook jobs, Pipelines, Models, Deployments, and SageMaker JumpStart.
- ML Flow Tab:** Shows a flow named "untitled.flow" with one step. The step details show a "Data types" transform for "dataset_2200_basketball.csv". The target column is "points_per_minute". Problem type is set to "Classification".
- Bottom Panel:** Displays various icons for file operations, a search bar, and system status information (0°C, Sunny, 3:56 PM, 2023-03-07).

The screenshot shows the Amazon SageMaker Studio interface with the 'Data Wrangler' tab selected. On the left, a sidebar lists various services: Home, Data (selected), Feature Store, Clusters, AutoML, Experiments, Notebook Jobs, Pipelines, Models, Deployments, and SageMaker JumpStart. The main area displays a 'Table Summary: Basketball Summary' report. It includes a table of summary statistics:

summary	sex	assists_per_minute	height	time_played
count	96	95	95	95
mean	None	0.16127368421052637	189.89473684210526	25.947265157894703
stddev	None	0.06011012271908513	6.994638804432945	8.666799384948705
min	Female	0.0494	160	10.08
max	Male	0.3437	203	40.71

Below this is a 'Data table' showing individual data points:

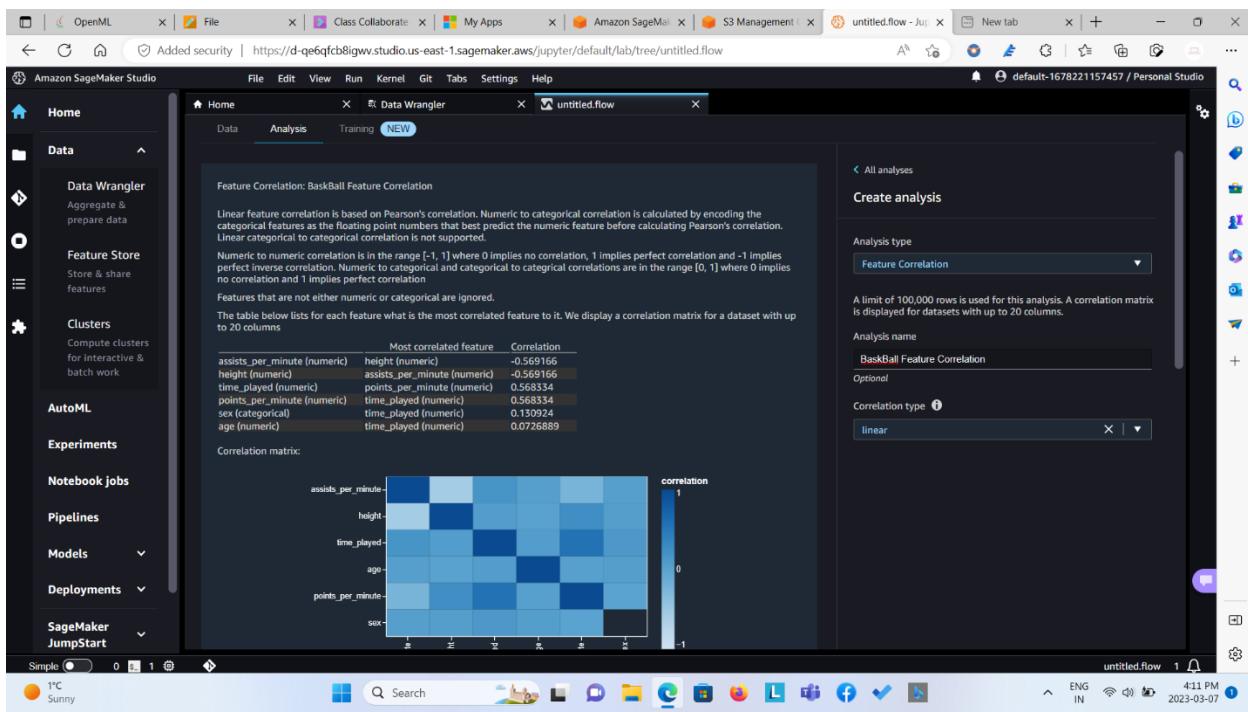
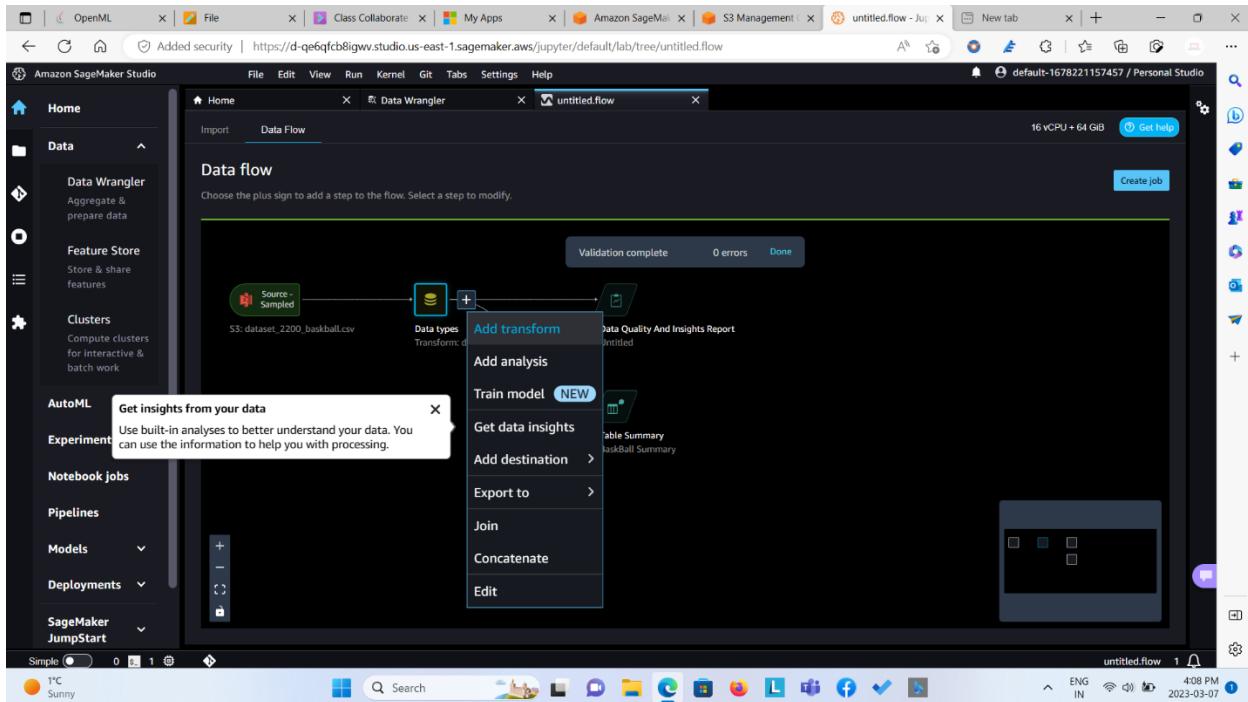
sex	assists_per_minute	height	time_played	age
Male	0.0888	201	36.02	28
Male	0.1399	198	39.32	30
Female	0.0747	198	38.8	26
Female	0.0983	191	40.71	30
Female	0.1276	196	38.4	28
Female	0.1671	201	34.1	31
Male	0.1906	193	36.2	30
Female	0.1061	191	36.75	27
Female	0.2446	185	38.43	29

On the right, there is a 'Create analysis' panel with fields for 'Analysis type' (set to 'Table Summary'), 'Analysis name' (set to 'BasketBall Summary'), and an optional note. Buttons for 'Preview' and 'Save' are at the bottom.

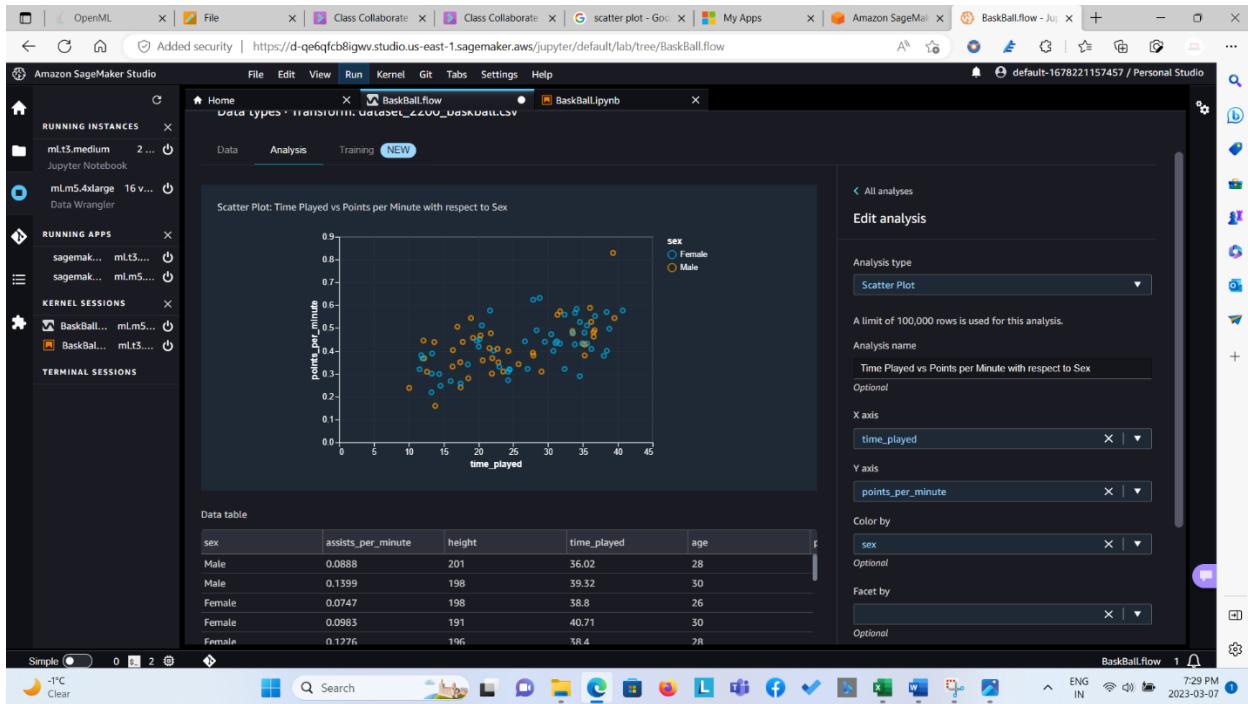
From the data analysis report I came to know about various things –

1. The overview of the dataset in term for number of instances, number of features, missing data percentage, duplicate rows, data types in short all about the statistic of the dataset.
2. The descriptive summary of all the features in terms of min, max, mean, median, standard deviation, mode.
3. Summary and Distribution of the target variable as well as other features..
4. Quick overview about the regression model in terms of Validation scores and Test scores having R2, MSE, RMSE, MAE, Max Error, and Median Absolutely error as its metrics.
5. Machine Learning Terminologies

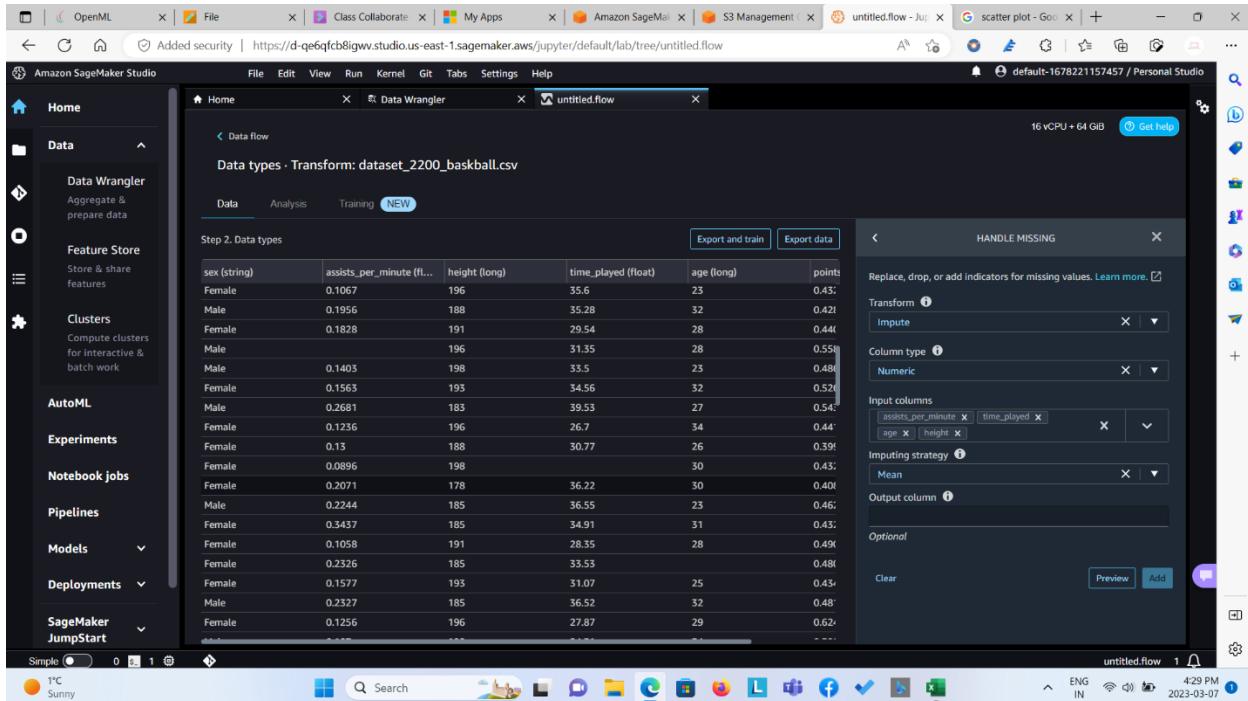
This was report was necessary to generate as from it we could know about the dataset which we will be working upon as well as we would get clear idea about the features and domain we are working for.



From the above heatmap, it is inferred that there is partial correlation between features which says that the features have partial significance over the target variable in order to predict the “points scored in a minute”.



From the above scatter diagram, it is said the there is a relationship between time played and points per minute. It shows positive correlation and it plays vital role in predicting the points scored in a minute.



The screenshot shows the Amazon SageMaker Studio Data Wrangler interface. On the left, a sidebar lists various services like Home, Data Wrangler, Feature Store, Clusters, AutoML, Experiments, Notebook jobs, Pipelines, Models, Deployments, and SageMaker JumpStart. The main area is titled "Data types · Transform: dataset_2200_basketball.csv". It has tabs for Data, Analysis, Training, and NEW. The Data tab is selected, showing a preview of the dataset with columns: sex (string), assists_per_minute (float), height (float), time_played (float), age (float), and points (float). Below the preview, there are buttons for "Export and train" and "Export data". To the right, a panel titled "HANDLE MISSING" is open, showing the configuration for handling missing values. It includes sections for Transform (Impute), Column type (Numeric), Input columns (assists_per_minute, time_played), Imputing strategy (Mean), and Output column. There are also optional settings and a preview button.

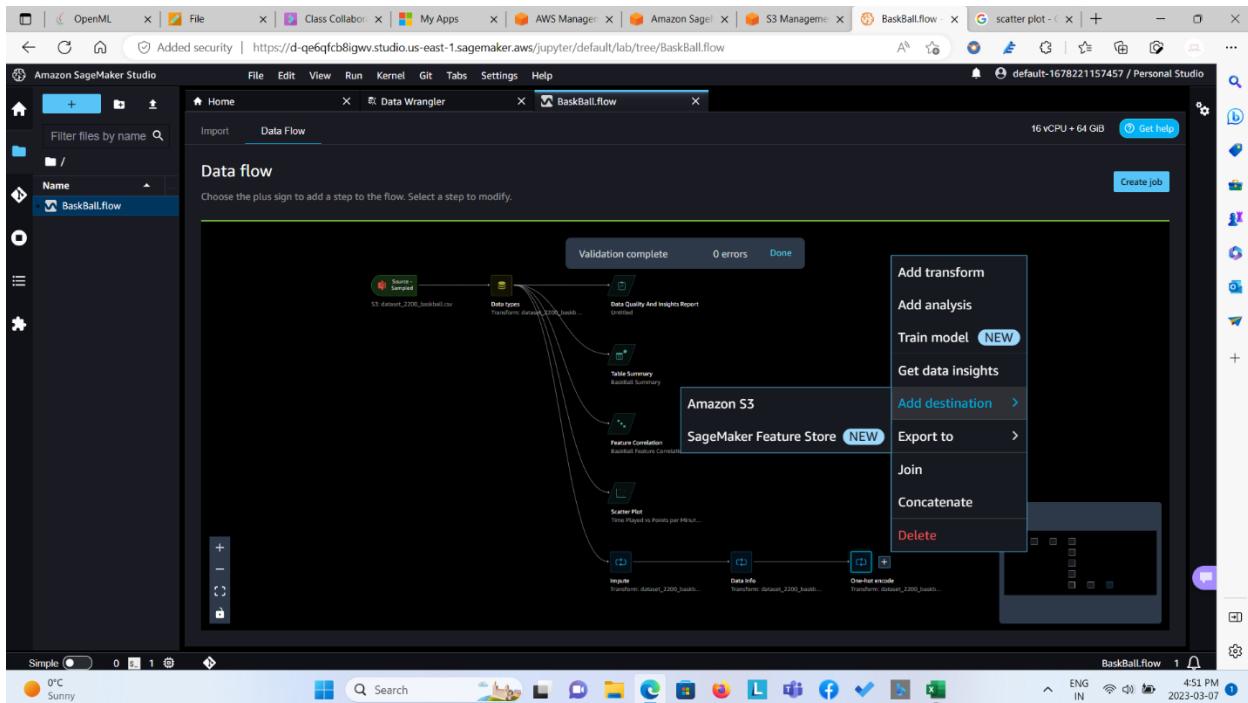
From the above data pre-processing step, we have replaced the missing values with the mean of the particular column. We did this so that we have a complete structured dataset in order to get the desire results in terms of accuracy and precision.

This screenshot shows the "Encode Categorical" transform step in the Amazon SageMaker Studio Data Wrangler. The main area displays the transformed dataset with columns: time_played (float), age (float), points_per_minute (float), sex_Female (float), and sex_Male (float). The "sex" column has been converted into two binary columns indicating whether the row represents a female or male player. The right-hand panel, titled "ENCODE CATEGORICAL", contains configuration options for this transformation, including "One-hot encode" as the method, "sex" as the input column, and "Keep" as the invalid handling strategy. Other options like "Input already ordinal encoded" and "Drop last" are also present.

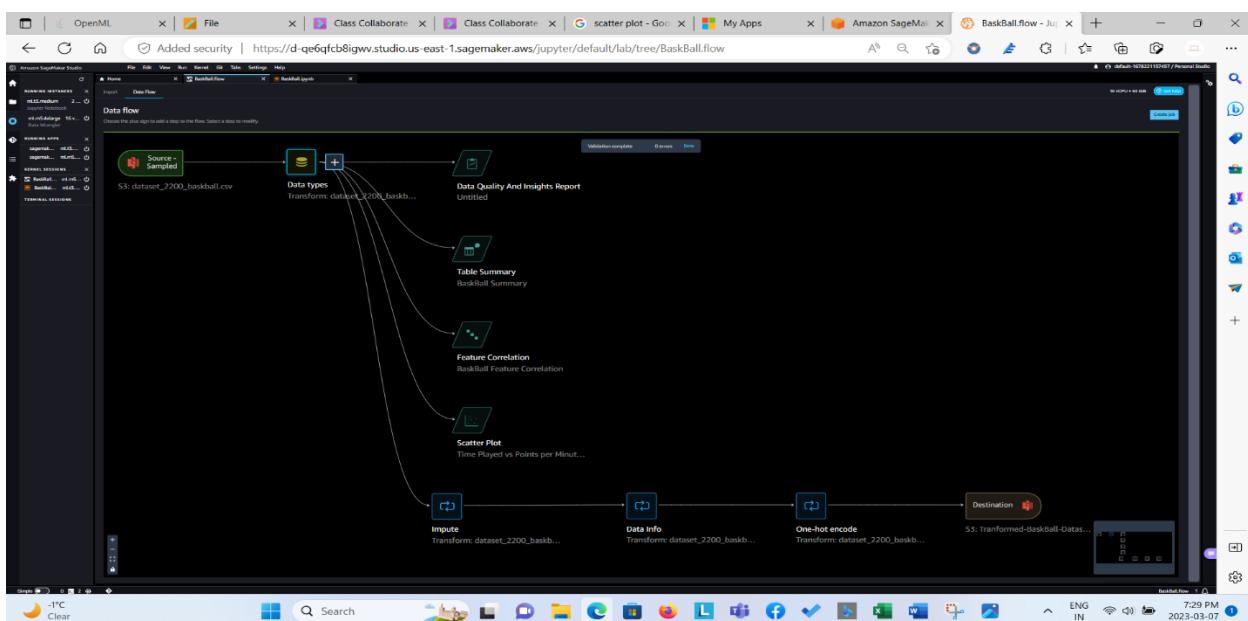
The step for one-hot-encoding is necessary for feature named “sex”, as we know that regression model can only be applied upon the features that are quantitative in nature. Here

sex is categorical variable, which also plays vital role in the prediction so to satisfy the regression condition we are converting it into numerical by encoding it as 0 and 1.

Step 6: Adding the Destination so that we can let the data flow to execute and the result to save into S3 bucket.



The Final Data-Flow



Later, created the job and the output was exported to the s3 bucket from data wrangler/data flow.