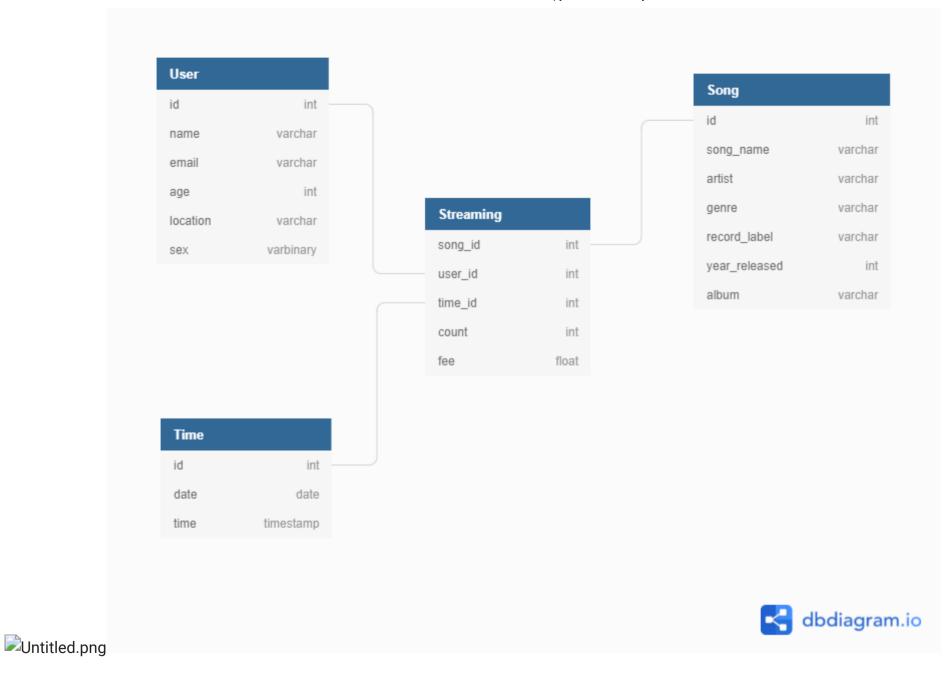


```
#!pip install cubes
#!pip install -Iv sqlalchemy==1.3.9

from sqlalchemy import create_engine
from cubes.tutorial.sql import create_table_from_csv
```

# ▼ ECS766 Coursework 2 - Elliot Linsey

Q.1A



# Q.1B

We need to perform a dicing operation as you are acting on 2 dimensions. This involves slicing by the given time period you wanted (October 2021), you may need to drill down to get to this specific month. Then slice by song dimension to the given song. From here you would use an aggregate function to produce the total fee from within that month.

# Q1.C

Using the formula  $\Pi_{i=1}^n(L_i+1)$  equates to 5x2x2 = 20. Therefore, this cube will contain 20 cuboids.

# Q2

#### **Base Table**

RID	Brand	Branch	
R1	Audi	TH	
R2	Audi	N	
R3	Audi	Н	
R4	Ford	TH	
R5	Ford	N	
R6	Ford	Н	
R7	Mini	TH	
R8	Mini	N	
R9	Mini	Н	

#### **Bitmap Index Table on Brand**

RID	Audi	Ford	Mini
R1	1	0	0
R2	1	0	0

RID	Audi	Ford	Mini
R3	1	0	0
R4	0	1	0
R5	0	1	0
R6	0	1	0
R7	0	0	1
R8	0	0	1
R9	0	0	1

# **-** Q3

```
engine = create_engine('sqlite:///data.sqlite')
create table from csv(engine,
                      "IBRD Balance Sheet FY2010.csv",
                      table name="ibrd balance",
                      fields=[
                          ("category", "string"),
                          ("category_label", "string"),
                          ("subcategory", "string"),
                          ("subcategory_label", "string"),
                          ("line_item", "string"),
                          ("year", "integer"),
                          ("amount", "integer")],
                      create id=True
from cubes import Workspace
workspace = Workspace()
workspace.register_default_store("sql", url="sqlite:///data.sqlite")
workspace.import_model("tutorial_model.json")
```

```
cube = workspace.cube("ibrd_balance")
browser = workspace.browser(cube)
result = browser.aggregate()
result.summary["record_count"]
```

Below is my JSON Model file, the two added functions are named 'maximum' and 'minimum' in the 'aggregates' section.

```
# {
      "dimensions": [
#
           "name":"item",
           "levels": [
                       "name": "category",
                       "label": "Category",
                       "attributes": ["category", "category_label"]
                  },
                       "name": "subcategory",
                       "label": "Sub-category",
                       "attributes": ["subcategory", "subcategory_label"]
                  },
                       "name": "line item",
                       "label": "Line Item",
                       "attributes": ["line item"]
          {"name":"year", "role": "time"}
      ],
      "cubes": [
#
#
               "name": "ibrd_balance",
```

```
"dimensions": ["item", "year"],
#
              "measures": [{"name":"amount", "label":"Amount"}],
              "aggregates": [
                      {
                          "name": "amount sum",
                          "function": "sum",
                          "measure": "amount"
                      },
                          "name": "record count",
                          "function": "count"
                      },
                           "name": "maximum",
                          "function": "max",
                          "measure": "amount"
                      },
                           "name": "minimum",
                          "function": "min",
                          "measure": "amount"
                  ],
              "mappings": {
                            "item.line item": "line item",
                            "item.subcategory": "subcategory",
                             "item.subcategory label": "subcategory label",
                            "item.category": "category",
                            "item.category_label": "category_label"
                           },
              "info": {
                  "min_date": "2010-01-01",
                  "max date": "2010-12-31"
```

### **-** Q4

```
engine = create engine('sqlite:///data.sqlite')
create_table_from_csv(engine,
                      "country-income.csv",
                      table name="country income",
                      fields=[
                          ("region", "string"),
                          ("age", "integer"),
                          ("income", "integer"),
                          ("online shopper", "string")
                      create id=True
workspace = Workspace()
workspace.register default store("sql", url="sqlite:///data.sqlite")
workspace.import model("country income model.json")
cube = workspace.cube("country income")
browser = workspace.browser(cube)
result = browser.aggregate()
result.summary["record count"]
```

Below is my JSON file representation for the model of this data cube.

```
# {
      "dimensions": [
#
             "name": "region",
            "levels":[{"name":"region",
                       "label": "Region",
                       "attributes": ["region"]}]
          },
            "name": "age",
            "levels":[{"name":"age",
                       "label": "Age",
                       "attributes": ["age"]}]
          },
            "name": "online shopper",
            "levels":[{"name":"online shopper",
                       "label": "Online Shopper",
                       "attributes": ["online shopper"]}]
      "cubes": [
               "name": "country_income",
              "dimensions": ["region", "age", "online_shopper"],
               "measures": [{"name":"income", "label":"Income"}],
               "aggregates": [
                           "name": "amount_sum",
                           "function": "sum",
                           "measure": "income"
                       },
```

```
11/11/21, 8:26 PM
   #
                              "name": "record_count",
                              "function": "count"
                          },
                              "name": "maximum",
                              "function": "max",
                              "measure": "income"
                          },
                              "name": "minimum",
                              "function": "min",
                              "measure": "income"
                          },
                              "name": "average",
                              "function": "avg",
                              "measure": "income"
                          }
                      ],
                  "mappings": {
                                 "region.region": "region",
                                "age.age": "age",
                                 "online shopper.online shopper": "online shopper"
   # }
   result = browser.aggregate()
   print('Full Summary:')
   print(result.summary)
   print('Region Summary:')
   result = browser.aggregate(drilldown=["region"])
   for record in result:
       print(record)
```

```
print('Online Shopper Summary:')
result = browser.aggregate(drilldown=["online shopper"])
for record in result:
    print(record)
     Full Summary:
     {'amount sum': 768200, 'record count': 10, 'maximum': 99600, 'minimum': 57600, 'average': 76820.0}
     Region Summary:
     {'region': 'Brazil', 'amount sum': 193200, 'record count': 3, 'maximum': 73200, 'minimum': 57600, 'average': 64400.0}
     {'region': 'India', 'amount sum': 331200, 'record count': 4, 'maximum': 94800, 'minimum': 69600, 'average': 82800.0}
     {'region': 'USA', 'amount sum': 243800, 'record count': 3, 'maximum': 99600, 'minimum': 64800, 'average': 81266.6666666667}
     Online Shopper Summary:
     {'online shopper': 'No', 'amount sum': 386400, 'record count': 5, 'maximum': 99600, 'minimum': 62400, 'average': 77280.0}
     {'online shopper': 'Yes', 'amount sum': 381800, 'record count': 5, 'maximum': 94800, 'minimum': 57600, 'average': 76360.0}
import cubes as cubes
cuts = [cubes.RangeCut("age",[40],[50])]
cell = cubes.Cell(cube, cuts)
result = browser.aggregate(cell, drilldown=["age"])
print('Ages between 40 and 50 Summary:')
for record in result:
    print(record)
     Ages between 40 and 50 Summary:
     {'age': 40, 'amount sum': 69600, 'record count': 1, 'maximum': 69600, 'minimum': 69600, 'average': 69600.0}
     {'age': 42, 'amount sum': 80400, 'record count': 1, 'maximum': 80400, 'minimum': 80400, 'average': 80400.0}
     {'age': 43, 'amount sum': 73200, 'record count': 1, 'maximum': 73200, 'minimum': 73200, 'average': 73200.0}
     {'age': 45, 'amount sum': 79400, 'record count': 1, 'maximum': 79400, 'minimum': 79400, 'average': 79400.0}
     {'age': 46, 'amount_sum': 62400, 'record_count': 1, 'maximum': 62400, 'minimum': 62400, 'average': 62400.0}
     {'age': 49, 'amount sum': 86400, 'record count': 1, 'maximum': 86400, 'minimum': 86400, 'average': 86400.0}
```

#### Part 2

#### **-** 01

The formula for Euclidean distance is  $d(p,q) = \sqrt{\sum_{i=1}^n (q_i - p_i)^2}$ 

With these x values below we find the distances between x1 and x2 to x3.

It would classify x3 as class y2 as the euclidean distance is closest to point x2 at 2.24 compared to x1 at 5.0.

### Q2

When K is odd in a KNN algorithm we do not need a tie-breaking policy because one class will always have a greater number of points closer to the new point than the other class. For example if K=5, the closest you can get to a tie will be 3 points of class y1 and 2 points of class y2 being the nearest neighbours, in this situation and with all odd K numbers the point will always be designated as the class with the highest number of neighbours, therefore it would be classed as label y1.

### **Q3**

A classifier that has an accuracy of 99.9% can be terrible for some datasets if 0.1% of the data is another class. This algorithm may label every data point as the same class and therefore miss the 0.1% that is different. While this appears to give a very high accuracy, it is actually not truly classifying anything rather than just assigning the same class to all data points. This can be extremely bad if misidentifying that 0.1% of data

results in a far greater cost than identifying 99.9% of the other data. In essence, you have an accuracy of 99.9% but a TPR (if the 0.1% of data is the positive label) of 0%.

## 04

A precision of 1.0 and low recall of 0.1 means that out of the data you predicted positive, all were correctly identified as true positive. However, the low recall means that you actually missed a lot of the true positives and mislabelled them as false negatives. Therefore, when it detects a positive it can be trusted but it is far less trustworthy when detecting negatives as there is a high chance they could be positive. In terms of this classifier, it should not be trusted if it detects that point to not belong to class y and it should be trusted if it detects it to belong to class y.

# Q5

The K=1 classifier struggled the most with classes 4 and 9, correctly identifying 40 4s but misidentifying 4 of them (10%) as 9s.

### **-** Q6

```
import gzip
import pickle

import numpy as np
import pandas as pd
import matplotlib.pyplot as plt
import seaborn as sns
sns.set_style('darkgrid')

# Selecting the training data from the original dataset
f = gzip.open('mnist.pkl.gz', 'rb')
X, y = pickle.load(f, encoding='latin1')[0]
f.close()
```

```
# Subsampling
sample_size = 2000
X, y = X[:sample_size], y[:sample_size]

from sklearn.model_selection import train_test_split

X_train, X_test, y_train, y_test = train_test_split(X, y, test_size=0.2, random_state=0)

from sklearn import svm
clf = svm.SVC()
clf.fit(X_train,y_train)
print('Test dataset accuracy: ' + str(clf.score(X_test,y_test)))

Test dataset accuracy: 0.9275
```

#### ▼ Q7

```
from sklearn.model_selection import GridSearchCV
from sklearn.ensemble import RandomForestClassifier

parameters = {'n_estimators': [50,100,200], 'max_features': [0.1,0.25]}

rfc = RandomForestClassifier()
    rfc_cv = GridSearchCV(rfc, parameters, cv=5)
    rfc_cv.fit(X_train, y_train)
    print('Best hyperparameter settings: {0}.'.format(rfc_cv.best_params_))
    print('Average accuracy across folds of best hyperparameter setting: {0}.'.format(rfc_cv.best_score_))
    print('Test dataset accuracy of best hyperparameter setting: {0}.'.format(rfc_cv.score(X_test, y_test)))

    Best hyperparameter settings: {'max_features': 0.1, 'n_estimators': 100}.
    Average accuracy across folds of best hyperparameter setting: 0.911875.
    Test dataset accuracy of best hyperparameter setting: 0.9025.
```