# Setup Virtual Environment and Install Pyspark:

Here are the commands used to set up a virtual environment and install Pyspark.

```
python3 -m venv itvg-venv
source itvg-venv/scripts/activate
pip install pyspark
```

Here is how we can get started with local development of data engineering pipelines using Spark.

- Create Virtual Environment - `python3.7 -m venv itvg-venv`
- Activate virtual environment - `source itvg-venv/bin/activate`
- Install PySpark for local development - `pip install pyspark==2.4.*`
- Open using PyCharm and make sure appropriate virtual environment is used from the virtual environment which we have setup.
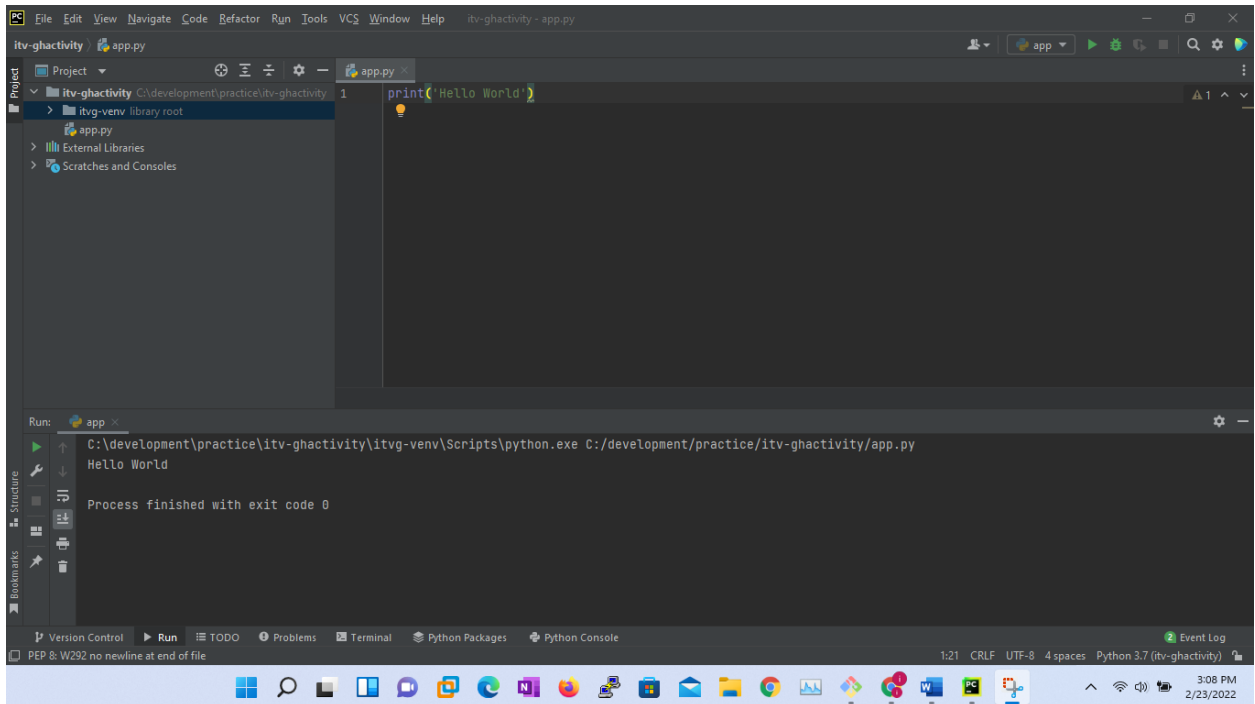


*Figure 1Created new environment with pyspark and python and opened in pycharm*
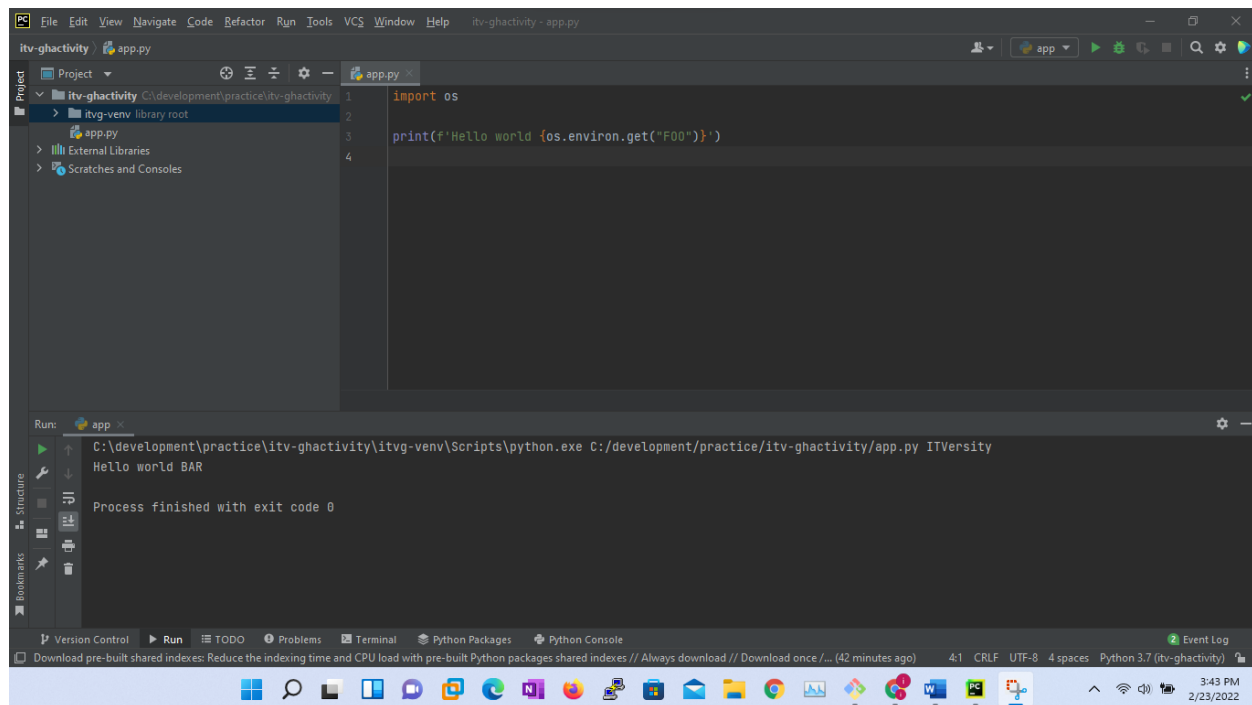
*Figure 2 Passed the environment variable from operating system into the program*

## Read data from files:

Let us develop the code to read the data from files into Spark Dataframes.

Create directory for data and copy some files into it.

mkdir -p data/itv-github/landing/ghactivity

cd data/itv-github/landing/ghactivity

wget https://data.gharchive.org/2021-01-13-0.json.gz

wget https://data.gharchive.org/2021-01-14-0.json.gz

wget https://data.gharchive.org/2021-01-15-0.json.gz

Create a Python program by name read.py. We will create a function by name from_files. It reads the data from files into Dataframe and returns it.

```
def from_files(spark, data_dir, file_pattern, file_format):
    df = spark. \
        read. \
        format(file_format). \
        load(f'{data_dir}/{file_pattern}')
```

return df



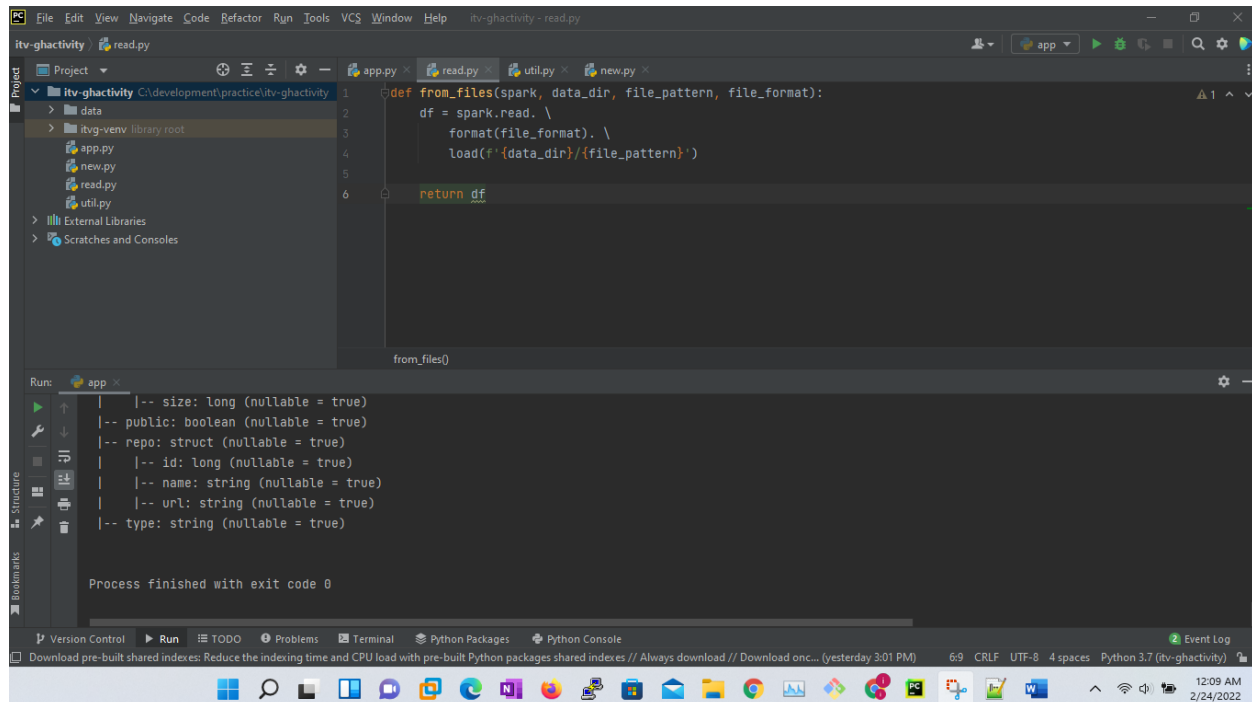*Figure 3read.py*

Call the program from app.py. For now review schema and data.

```python
import os
from util import get_spark_session
from read import from_files
def main():
    env = os.environ.get('ENVIRON')
    src_dir = os.environ.get('SRC_DIR')
    file_pattern = f"{os.environ.get('SRC_FILE_PATTERN')}-*"
    src_file_format = os.environ.get('SRC_FILE_FORMAT')
    spark = get_spark_session(env, 'GitHub Activity - Reading Data')
    df = from_files(spark, src_dir, file_pattern, src_file_format)
    df.printSchema()
    df.select('repo.*').show()
if __name__ == '__main__':
    main()
```
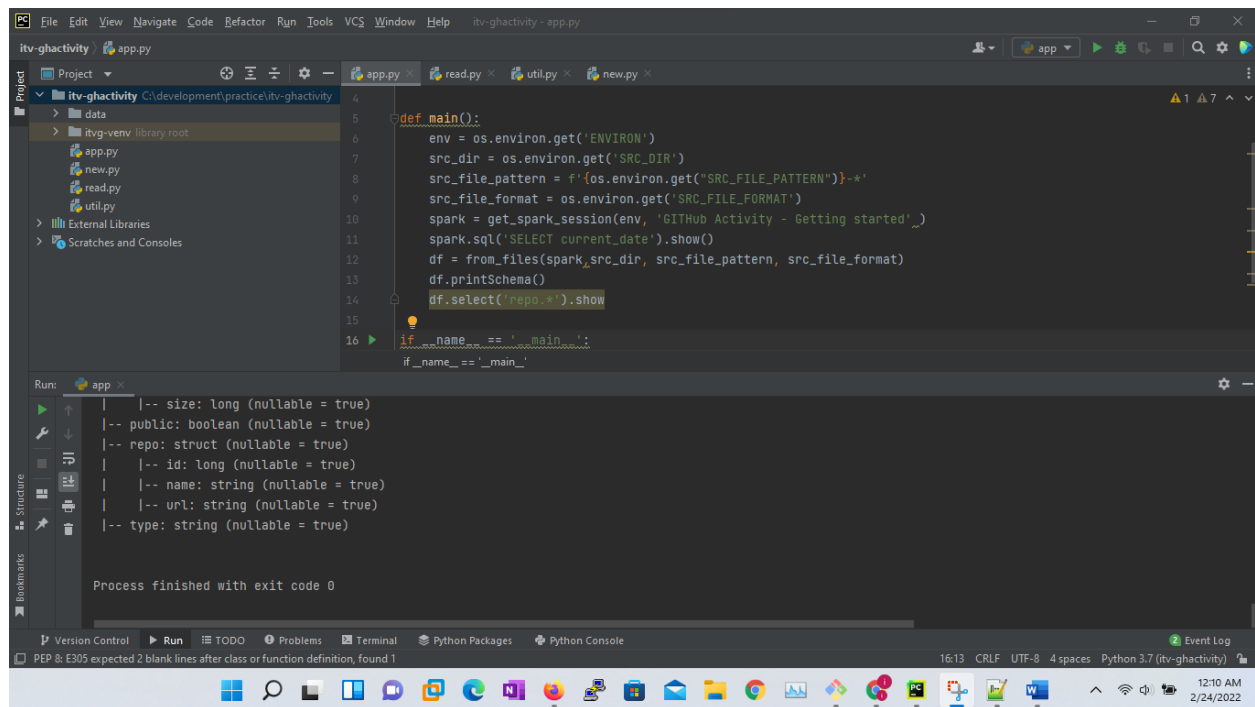
## Process data using Spark APIs:

Create a Python program by name process.py. We will create a function by name df_transform. It partitions the Dataframe using specified field.

```python
from pyspark.sql.functions import year, \
    month, dayofmonth

def transform(df):
    return df.withColumn('year', year('created_at')). \
        withColumn('month', month('created_at')). \
        withColumn('day', dayofmonth('created_at'))
```

Call the program from app.py. For now review schema and data.

```python
import os
from util import get_spark_session
from read import from_files
from process import transform
```

```python
def main():

    env = os.environ.get('ENVIRON')

    src_dir = os.environ.get('SRC_DIR')

    file_pattern = f"{os.environ.get('SRC_FILE_PATTERN')}-*"

    src_file_format = os.environ.get('SRC_FILE_FORMAT')

    spark = get_spark_session(env, 'GitHub Activity - Partitioning Data')

    df = from_files(spark, src_dir, file_pattern, src_file_format)

    df_transformed = transform(df)

    df_transformed.printSchema()

    df_transformed.select('year', 'month', 'day').show()
```