# databricks™ Explore sample datasets

## List sample datasets available in Databricks CE

```
display(dbutils.fs.ls("/databricks-datasets"))
```

| path |
| --- |
| dbfs:/databricks-datasets/ |
| dbfs:/databricks-datasets/README.md |
| dbfs:/databricks-datasets/Rdatasets/ |
| dbfs:/databricks-datasets/SPARK_README.md |
| dbfs:/databricks-datasets/adult/ |
| dbfs:/databricks-datasets/airlines/ |
| dbfs:/databricks-datasets/amazon/ |
| dbfs:/databricks-datasets/asa/ |
| dbfs:/databricks-datasets/atlas_higgs/ |
| dbfs:/databricks-datasets/bikeSharing/ |
| dbfs:/databricks-datasets/cctvVideos/ |

```
selected_dataset = "sample_logs"
selected_dataset_path = "/databricks-datasets/" + selected_dataset

print(selected_dataset_path)
```

```
/databricks-datasets/sample_logs
```

## List the files in a dataset

```
filelist = dbutils.fs.ls(selected_dataset_path)
display(filelist)
```

| path |
| --- |
| dbfs:/databricks-datasets/sample_logs/_SUCCESS |
| dbfs:/databricks-datasets/sample_logs/part-00000 |
| dbfs:/databricks-datasets/sample_logs/part-00001 |

dbfs:/databricks-datasets/sample_logs/part-00002

dbfs:/databricks-datasets/sample_logs/part-00003

dbfs:/databricks-datasets/sample_logs/part-00004

dbfs:/databricks-datasets/sample_logs/part-00005

dbfs:/databricks-datasets/sample_logs/part-00006

dbfs:/databricks-datasets/sample_logs/part-00007

⬇

```python
from pyspark.sql import *

filelistDF = sqlContext.createDataFrame(filelist)


just_readme = filelistDF.filter((filelistDF.name == "README.md") | (filelistDF.name ==
"readme.md"))
display(just_readme)
```

OK

```python
filelistDF.createOrReplaceTempView("filelist_view")
```

```sql
%sql
select `path` from filelist_view where `name` like 'README%'
```

OK

```python
#print(filelist)
print(filelist[0][1])

with open("/dbfs/" + selected_dataset_path + "/part-00004") as f:
  x = ''.join(f.readlines())

print(x)
```

```
_SUCCESS
2.2.2.2 - - [21/Mar/2014:10:00:00 -0400] "GET /endpoint_875 HTTP/1.1" 500 396
127.0.0.1 - - [21/Jun/2014:10:00:00 -0700] "GET /endpoint_314 HTTP/1.1" 200 220
```

```python
import re
from pyspark.sql import Row

apache_access_log_format = '^(\S+) (\S+) (\S+) \[([\w:/]+\s[+\-]\d{4})\] "(\S+) (\S+) (\S+)" (\d{3}) (\d+)'

#Return a dictionary containing the parts of the Apache access log
def parse_apache_log_line(line):
  match = re.search(apache_access_log_format, line)
  if match is None:
    # Might want to just ignore, since it'a common to have dirty logs
    raise Error("Invalid log line format: %s" % line)
  return Row(
    ipAddress = match.group(1),
    clientIdentd = match.group(2),
    userId = match.group(3),
    dateTime = match.group(4),
    method = match.group(5),
    endpoint = match.group(6),
    protocol = match.group(7),
    responseCode = int(match.group(8)),
    contentSize = long(match.group(9)))
```

```
AttributeError: 'RemoteContext' object has no attribute 'textfile'
```