Attribution 4.0 International (CC BY 4.0)

# for a given day, get the popularity of a drive model: # cat data_Q3_2018.zip_folder/2018-07-27.csv | sed '1d' | cut -d',' -f3 | sort | uniq -c | sort -g -k1,1 # for every CSV file, get the date # find . | grep csv | while read fullpath; do

fullpath | sed 's/\//_/g' | sed 's/\.csv//g' | sed 's/zip_folder/_/g' | sed 's/data_//g'; done

# create a file per day containing the popularity of each model
https://stackoverflow.com/questions/17017732/changing-delimiter-of-the-uniq-command
https://pandas.pydata.org/pandas-docs/stable/reference/api/pandas.read_fwf.html

date; find . | grep csv | while read fullpath; do
  filename=`echo

echo fullpath | sed 's/zip_folder/_/g' | sed 's/data_//g' | sed 's/\.csv//g' | sed 's/\//_/g' | sed 's/\.//g' | sed 's/^_//g'`; cat

fullpath | sed '1d' | cut -d',' -f3 | sort | uniq -c | sort -g -k1,1 | sed 's/^ *//;s/ /,/' > count_of_models_on_{filename}.dat; done; date

```
In [1]:  import pandas
         print('pandas',pandas.__version__)
         import glob
         import pickle
         import numpy
         import seaborn
         import time
         import datetime
         import matplotlib.pyplot as plt

         pandas 0.23.4
```

```
In [2]:  list_of_dat = glob.glob('data_synthesized_from_csvs/count_of_models_per_
         day/count_of_models_on_*.dat')
         print(len(list_of_dat))

         2088
```

```
In [3]:  list_of_df=[]
         for path_to_dat in list_of_dat:
             date_str = path_to_dat[:-len('.dat')].split('_')[-1]
             date_as_dt = datetime.datetime.strptime(date_str, '%Y-%m-%d')
         #    print(path_to_dat)
             df = pandas.read_csv(path_to_dat,header=None)
             df.columns=[date_as_dt,'model']
             df=df.set_index('model')
             list_of_df.append(df)
```

list_of_models=[] for df in list_of_df: for model_name in df.index: list_of_models.append(model_name) list_of_models = list(set(list_of_models)) print(len(list_of_models))

```
In [4]:  df = pandas.concat(list_of_df,sort=False,axis=1) # join all the datafram
         es into a single df
         df = df.reindex(sorted(df.columns), axis=1) # order columns by calendar
          date
```

```
In [5]: df.shape

Out[5]: (113, 2088)
```

```
In [21]: sorted_df = df.loc[df.sum(axis=1).sort_values(ascending=False).index]
```

```
In [27]: seaborn.set(rc={'figure.figsize':(12,10)})
         seaborn.heatmap(sorted_df);
         plt.title('Backblaze drives by model over time',fontsize=14);
```



Backblaze drives by model over time