


```
Applications Places System cloudera@quickstart:~
Access documents, folders and network places cloudera@quickstart:~
File Edit View Search Terminal Help
In [11]: #remove whitespaces from headers

In [12]: buyclicksDataFrame = pd.read_csv('./buy-clicks.csv')

In [13]:

In [13]: buyclicksDataFrame = buyclicksDataFrame.rename(columns=lambda x: x.strip())

In [14]:

In [14]: buyclicksDataFrame.head(n=2)
Out[14]:
   timestamp txId  userSessionId  team  userId  buyId  price
0  2016-05-26 15:36:54  6004         5820    9   1300     2    3.0
1  2016-05-26 15:36:54  6005         5775   35    868     4   10.0

In [15]:

In [15]: #select 'userId' and 'price' and drops all others columns

In [16]: PurchasesDataFrame = buyclicksDataFrame[['userId', 'price']]

In [17]: PurchasesDataFrame.head(n=2)
Out[17]:
   userId  price
0    1300     3.0
1     868    10.0

In [18]:

cloudera@quickstart:~ script-cp.txt (~) - gedit Cloudera Live: Welco... cloudera@quickstart:~
```

```
Applications Places System cloudera@quickstart:~
Browse and run installed applications cloudera@quickstart:~
File Edit View Search Terminal Help

In [18]: ##select 'userId' and 'adCount' and drops all others columns

In [19]: useradClicksDataFrame = adclicksDataFrame[['userId', 'adCount']]

In [20]: useradClicksDataFrame.head(n=2)
Out[20]:
   userId  adCount
0     611         1
1    1874         1

In [21]:

In [21]: #creates new file by adding each adCount per userId

In [22]: PerUserDataFrame = useradClicksDataFrame.groupby('userId').sum()

In [23]: PerUserDataFrame = PerUserDataFrame.reset_index()

In [24]: #rename the columns

In [25]: PerUserDataFrame.columns = ['userId', 'totalAdClicks']

In [26]: PerUserDataFrame.head(n=2)
Out[26]:
   userId  totalAdClicks
0         1             44
1         8             10

In [27]:

cloudera@quickstart:~ script-cp.txt (~) - gedit Cloudera Live: Welco... cloudera@quickstart:~
```


