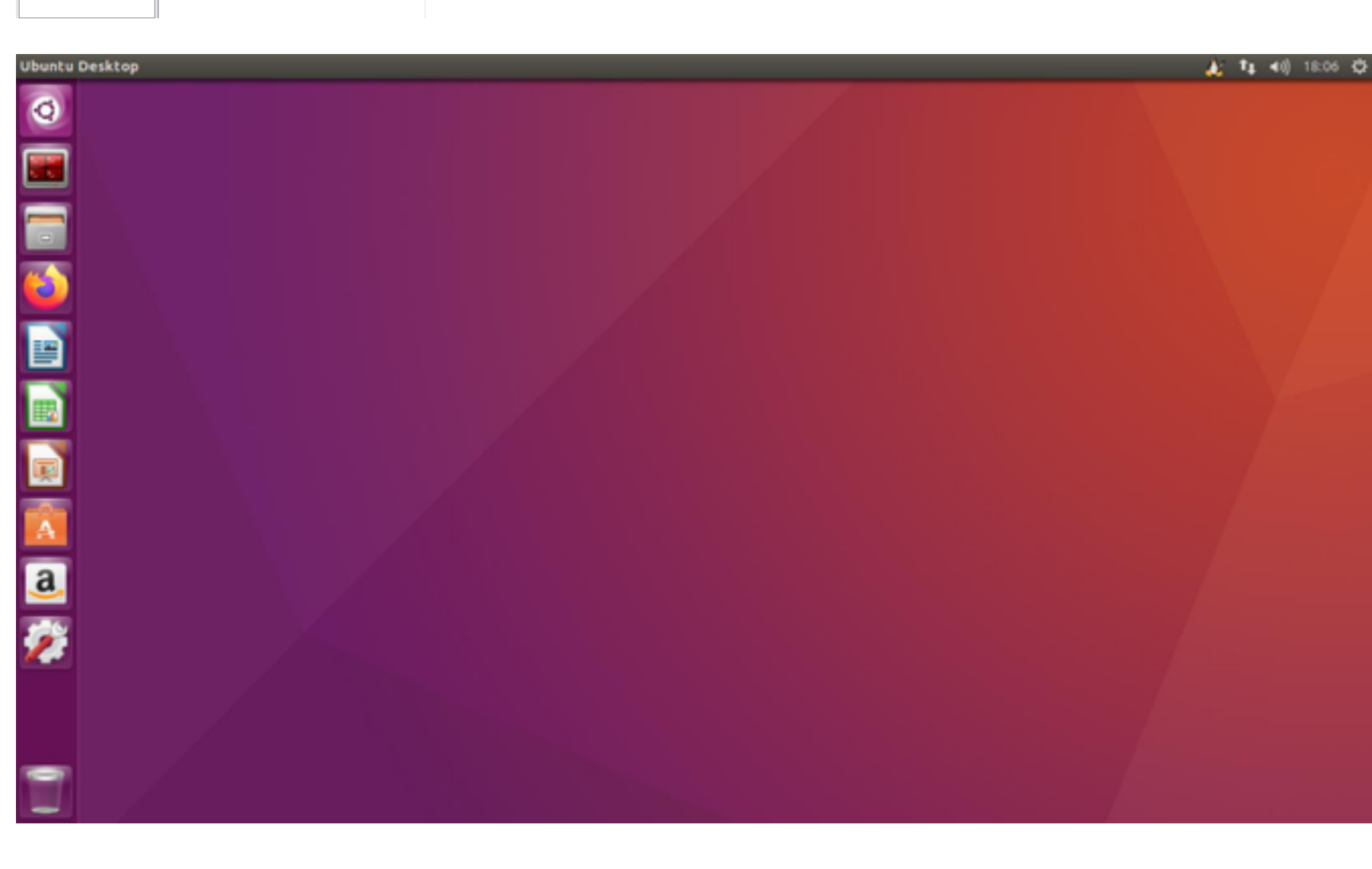


- 309554001 劉雨恩
- 309706013 黃柏元
- 0613144 葉之晴
- 309707007 黃如君

Step1. 虛擬機安裝

The screenshot shows the Ubuntu Software Center interface. The 'LibreOffice 64-bit' package is highlighted in the search results. The package description indicates it is a free office suite that can open and save files in the same format as Microsoft Office. The version number 4.2.4.2 is visible.



The screenshot shows the Windows 7 desktop with the Windows Update window open. The window title is 'Windows Update - Windows 7.0.7600.17514'. The main area shows 'In operation' with a progress bar. Below the progress bar, there is a table showing the Windows Update history for 'Windows 7.0.7600.17514'. The table has columns: Name, My Address, Last Check, Report, Capacity, Size, and Status. The status is '233'. At the bottom, there is a 'Showing 1 of 1 entries' message and buttons for 'Previous' and 'Next'.

Name	My Address	Last Check	Report	Capacity	Size	Status
Windows 7.0.7600.17514	Windows 7.0.7600.17514	26	26	87.2 GB	26 GB	233

Showing 1 of 1 entries

Previous Next

```
huser@master:~$ source ~/.bashrc
huser@master:~$ scala
Welcome to Scala version 2.11.6 (OpenJDK 64-Bit Server VM, Java 1.8.0_275).
Type in expressions to have them evaluated.
Type :help for more information.

scala>
```

```

root@ip-10-0-1-10:~# source ~/.bashrc
hadoopmaster~# pyspark
Python 2.7.12 (default, Jul 21 2016, 15:19:50)
Type "help", "copyright", "credits" or "license()" for more information.
In [0]: %configure
Spark's default log4j profile: org.apache.spark/log4j-default.properties
Use "setLogLevel" to adjust logging level or setLogLevel=(name,value).
Example: "setLogLevel('org.apache.spark.executor', 'WARN')".
Note: If you are unable to load native-hadoop library for
your platform... using built-in Java classes where applicable
Welcome to
      Spark
      version 2.0.0

Using Python version 2.7.12 (default, Jul 21 2016 15:19:50)
SparkSession available as 'spark'.
>>>

```

```
(gedit6046): Gtk-WARNING **: Calling inhibit failed: GDBus.Error::freedesktop:GDBus.Error::org.gnome.SessionManager: The name org.gnome.SessionManager was not provided by any service files

** (gedit6046): WARNING **: Set document metadata failed: 不支援指定屬性 metadata
(gedit6046): libxml2-ERROR **: Failed to enable module: /usr/share/xml/libxml/catalog/catalog.xml.gz
** (gedit6046): WARNING **: Set document metadata failed: 不支援指定屬性 metadata
to:gedit-encoding

** (gedit6046): WARNING **: Set document metadata failed: 不支援指定屬性 metadata
to:gedit-poolfont
libusermaster= source -v .bashrc
libusermaster= python -version
python 2.7.11 :: daemona 2.0 (64-bit)
libusermaster=:
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 878049 entries, 0 to 878048
Data columns (total 6 columns):
  Dates      878049 non-null object
  Category   878049 non-null object
  DayOfWeek  878049 non-null object
  PDistrict  878049 non-null object
  X          878049 non-null float64
  Y          878049 non-null float64
dtypes: float64(2), object(4)
```

memory usage: 40.2+ MB							
	Dates	Category	DayOfWeek	PdDistrict	X	Y	
0	2015-05-12 23:53:00	WARRANTS	Wednesday	NORTHERN	-122.4528982	37.774499	
1	2015-05-12 23:53:00	OTHER OFFENSES	Wednesday	NORTHERN	-122.4528982	37.774499	
2	2015-05-12 23:33:00	OTHER OFFENSES	Wednesday	NORTHERN	-122.4243683	37.800413	
3	2015-05-12 23:33:00	LARCENY/HI/FT	Wednesday	NORTHERN	-122.4243683	37.800413	
4	2015-05-12 23:33:00	LARCENY/HI/FT	Wednesday	PARK	-122.4387878	37.771151	
878044	2013-01-06 00:15:00	ROBBERY	Monday	TARPAVAL	-122.4590303	37.714548	
878045	2013-01-06 00:15:00	LARCENY/HI/FT	Monday	INGLSIDE	-122.447364	37.714548	
878046	2013-01-06 00:01:00	LARCENY/HI/FT	Monday	SOUTHERN	-122.4503390	37.780206	
878047	2013-01-06 00:01:00	VANDALISM	Monday	SOUTHERN	-122.390551	37.780206	
878048	2013-01-06 00:01:00	FORGERY/COUNTERFEITING	Monday	BAVIEW	-122.3846828	37.738212	

878049 rows x 6 columns

```

1 # Change 'DayofWeek' to dummy features
2 all_dum = pd.get_dummies(pd.Series(data['DayofWeek']))
3 new_data = pd.concat([all_dum, axis=1], drop(columns=['DayofWeek']))
4
5 # Change 'PdDistrict' to dummy features
6 all_dum = pd.get_dummies(pd.Series(data['PdDistrict']))
7 new_data = pd.concat([new_data, all_dum], axis=1) \
8                     .drop(columns='PdDistrict')
9 new_data.info()

```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 878049 entries, 0 to 878048
Data columns (total 21 columns):
Dates                878049 non-null object
Category             878049 non-null object
X                    878049 non-null float64
Y                    878049 non-null float64
Friday               878049 non-null uint8
Monday               878049 non-null uint8
Saturday             878049 non-null uint8
Sunday               878049 non-null uint8
Thursday             878049 non-null uint8
```

Tuesday	878049	non-null	uint8
Wednesday	878049	non-null	uint8
DAYVIEW	878049	non-null	uint8
CENTRAL	878049	non-null	uint8
INGLESIDE	878049	non-null	uint8
MISSION	878049	non-null	uint8
NORTHERN	878049	non-null	uint8
PARK	878049	non-null	uint8
RICHMOND	878049	non-null	uint8
SOUTHERN	878049	non-null	uint8
TARAVAL	878049	non-null	uint8
TENDERLOIN	878049	non-null	uint8

```
2.3 Take hour from 'Dates' & change to dummy features
```

```
1 new_data['Dates'] = new_data['Dates'].str[11:-6]
2 all_dum = pd.get_dummies(pd.Series(new_data['Dates']))
3 new_data = pd.concat([new_data, all_dum], axis=1).drop(columns=['Dates'])
4 new_data.info()
5
6 # output preprocessed data csv
7 new_data.to_csv('data/train.csv')
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 878049 entries, 0 to 878048
Data columns (total 44 columns):
Category      878049 non-null object
X              878049 non-null float64
Y              878049 non-null float64
Friday         878049 non-null uint8
Monday         878049 non-null uint8
Saturday       878049 non-null uint8
```

Sunday	878049	non-null	uint8
Thursday	878049	non-null	uint8
Tuesday	878049	non-null	uint8
Wednesday	878049	non-null	uint8
RAYVIEW	878049	non-null	uint8
CENTRAL	878049	non-null	uint8
INGLESIDE	878049	non-null	uint8
MISSION	878049	non-null	uint8
NORTHERN	878049	non-null	uint8
PARK	878049	non-null	uint8
RICHMOND	878049	non-null	uint8
SOUTHERN	878049	non-null	uint8
TARAVAL	878049	non-null	uint8
TENDERLOIN	878049	non-null	uint8
00	878049	non-null	uint8
01	878049	non-null	uint8
02	878049	non-null	uint8
03	878049	non-null	uint8
04	878049	non-null	uint8
05	878049	non-null	uint8
06	878049	non-null	uint8
07	878049	non-null	uint8
08	878049	non-null	uint8
09	878049	non-null	uint8
10	878049	non-null	uint8
11	878049	non-null	uint8
12	878049	non-null	uint8
13	878049	non-null	uint8
14	878049	non-null	uint8
15	878049	non-null	uint8
16	878049	non-null	uint8
17	878049	non-null	uint8
18	878049	non-null	uint8
19	878049	non-null	uint8
20	878049	non-null	uint8
21	878049	non-null	uint8
22	878049	non-null	uint8
23	878049	non-null	uint8

Spark

- 訓練模型 `DecisionTree.trainClassifier`，並計

```
1. 將處理好的資料 (data/train.csv) load 進來

1 | print(" Load Data...")
2 | rawDataWithHeader = csv.textFile(Path("data/train.csv"))
3 | header = rawDataWithHeader.first()
4 | rawData = rawDataWithHeader.filter(lambda x:x != header)
5 | rData = rawData.map(lambda x : x.replace(",",""))
6 | lines = rData.map(lambda x: x.split("\n"))
7 | print("共計:" + str(lines.count()) + "筆")
8 | lines.take(10)
```

共計：878049冊

```
lines.take(10)
```

- [0.25, 0.0, 0.1, 0.1224150077, 0.97.0077]

```
2. Construct RDD[LabeledPoint]

1 | labelpointRDD = lines.map(lambda r:LabeledPoint(extract_label(r), \
2 |                                     extract_features(r, len(r) - 1)))
3 | print ("labelpointRDD = ", labelpointRDD.first(), "\n")
```

- ```
('labelpointRDD = ', LabeledPoint(37.0, [23.0,6.0,4.0,-122.4258917,37.7745986]), '\n')
```

```
3. Randomly divide the data into 3 parts: train (80%), validation (10%), and test (10%)

1 (trainData, validationData, testData) = labelpointRDD.randomSplit([8, 1, 1])
2 print(" trainData : " + str(trainData.count()) +
3 " validationData : " + str(validationData.count()) +
```

```
4 | " testData : " + str(testData)
```

- 

```
trainData: 702525 validationData: 87520 testData: 88004
```

#### 4. 訓練模型

使用 `pyspark.mllib.tree.DecisionTree.trainClassifier`，並計算 Accuracy、Precision 和 Recall

```
1 # train model
2 from pyspark.mllib.tree import DecisionTree
3 from pyspark.mllib.evaluation import MulticlassMetrics
4 model = DecisionTree.trainClassifier(trainData, numClasses=39,
5 categoricalFeaturesInfo={},
6 impurity="entropy",
7 maxDepth=10, maxBins=10)
8
9 # evaluate model
10 def evaluateModel(model, validationData):
11 score = model.predict(validationData.map(lambda p: p.features))
12 scoreAndLabels = score.zip(validationData.map(lambda p: p.labels))
13 metrics = MulticlassMetrics(scoreAndLabels)
14 accuracy = metrics.accuracy
15 recall = metrics.recall()
16 precision = metrics.precision()
17 print "Accuracy = ", str(accuracy)
18 print "Recall = ", str(recall)
19 print "Precision = ", str(precision)
```

[illegible]

## One computer

1. Randomly split to training (75%) and test (25%) dataset

```
1 from sklearn.model_selection import train_test_split
2 x_train, x_test, y_train, y_test = train_test_split(
3 new_data.drop(columns=['Category']),
4 new_data['Category'], test_size=0.25,
```

- ```
5 random_state=42)
```

2. 訓練模型 sklearn.tree.DecisionTreeClassifier，並計算 Accuracy、Precision 和 Recall

```
1 from sklearn import tree
2
3 # train model & predict
4 clf = tree.DecisionTreeClassifier()
5 clf = clf.fit(x_train, y_train)
6 y_test_pred = clf.predict(x_test)
7
8 from sklearn.metrics import accuracy_score, precision_score, recall_score
```

```
13 recall = recall_score(y_test, y_test_pred, av
14 print('Accuracy: {}, Precision: {}, Recall: {}
```

Value	Spark	One computer
Accuracy	0.229124771481	0.24598543138675158