In [24]: ▶
```python
import os
import json
from pathlib import Path
import zipfile
import email
from email.policy import default
from email.parser import Parser
from datetime import timezone
from collections import namedtuple

import pandas as pd
import s3fs
from bs4 import BeautifulSoup
from dateutil.parser import parse
from chardet.universaldetector import UniversalDetector

from pyspark.ml import Pipeline
from pyspark.ml.feature import CountVectorizer
from pyspark.ml.feature import HashingTF, Tokenizer
from pyspark.sql import SparkSession
from pyspark.sql.functions import col
from pyspark.ml.pipeline import Transformer
from pyspark.sql.functions import udf
from pyspark.sql.types import StructType, StringType

import pandas as pd
import shutil
```

In [26]: ▶
```python
from pyspark.sql import functions
```

In [2]: ⏭

```python
1  current_dir = Path(os.getcwd()).absolute()
2  results_dir = current_dir.joinpath('results')
3  results_dir.mkdir(parents=True, exist_ok=True)
4  data_dir = current_dir.joinpath('data')
5  data_dir.mkdir(parents=True, exist_ok=True)
6  enron_data_dir = data_dir.joinpath('enron')
7
8  output_columns = [
9          'payload',
10         'text',
11         'Message_ID',
12         'Date',
13         'From',
14         'To',
15         'Subject',
16         'Mime-Version',
17         'Content-Type',
18         'Content-Transfer-Encoding',
19         'X-From',
20         'X-To',
21         'X-cc',
22         'X-bcc',
23         'X-Folder',
24         'X-Origin',
25         'X-FileName',
26         'Cc',
27         'Bcc'
28 ]
29
30 columns = [column.replace('-', '_') for column in output_columns]
31
32 ParsedEmail = namedtuple('ParsedEmail', columns)
33
34 spark = SparkSession\
35     .builder\
36     .appName("Assignment04")\
37     .getOrCreate()
```

The following code loads data to your local JupyterHub instance. You only need to run this once.

In [3]: ⏭

```python
1  def move_data_to_local():
2      dst_data_path = data_dir.joinpath('enron.zip')
3      src_data_path = "C:/Users/bibek/Documents/GitHub/dsc650/data/external
4      shutil.copy(src_data_path, dst_data_path)
5
6      with zipfile.ZipFile(src_data_path) as f_zip:
7          f_zip.extractall(path=enron_data_dir)
8
9  move_data_to_local()
```

This code reads emails and creates a Spark dataframe with three columns.

# Assignment 4.1

In [4]:
```python
def read_raw_email(email_path):
    detector = UniversalDetector()

    try:
        with open(email_path) as f:
            original_msg = f.read()
    except UnicodeDecodeError:
        detector.reset()
        with open(email_path, 'rb') as f:
            for line in f.readlines():
                detector.feed(line)
                if detector.done:
                    break
        detector.close()
        encoding = detector.result['encoding']
        with open(email_path, encoding=encoding) as f:
            original_msg = f.read()

    return original_msg

def make_spark_df():
    records = []
    for root, dirs, files in os.walk(enron_data_dir):
        for file_path in files:
            ## Current path is now the file path to the current email.
            current_path = Path(root).joinpath(file_path)

            # get original message
            original_msg = read_raw_email(current_path)

            # get username
            path = os.path.normpath(root)
            x = path.split(os.sep)
            username = x[11]

            # get relative path
            relative_path = os.path.relpath(current_path, enron_data_dir)
            id = relative_path

            records.append((id,username,original_msg))

    # PySpark dataframe
    df = spark.createDataFrame(records).toDF("id","username","original_ms
    return df
```

In [5]:
```python
df = make_spark_df()
```

In [6]:  ▶| 1 df.printSchema()

```
root
 |-- id: string (nullable = true)
 |-- username: string (nullable = true)
 |-- original_msg: string (nullable = true)
```

In [7]:  ▶| 1 df.show()

```
+--------------------+--------+--------------------+
|                  id|username|        original_msg|
+--------------------+--------+--------------------+
|   davis-d\2_trash\1_|  davis-d|Message-ID: <1774...|
|   davis-d\2_trash\2_|  davis-d|Message-ID: <2467...|
|   davis-d\2_trash\3_|  davis-d|Message-ID: <2833...|
|   davis-d\2_trash\4_|  davis-d|Message-ID: <1972...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <1964...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <7345...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <5686...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <7218...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <3016...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <1233...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <2215...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <1365...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <2251...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <8556...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <1807...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <2705...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <2977...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <3065...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <2798...|
|davis-d\2_trash\c...|  davis-d|Message-ID: <3108...|
+--------------------+--------+--------------------+
only showing top 20 rows
```

In [7]:  ▶|
```python
1 from pyspark.sql import SparkSession
2
3 spark = SparkSession.builder.master("local").getOrCreate()
4 print(spark.sparkContext.version)
```

```
2.4.5
```

# Assignment 4.2

Use `plain_msg_example` and `html_msg_example` to create a function that parses an email message.

In [8]:

```
1   plain_msg_example = """
2   Message-ID: <6742786.1075845426893.JavaMail.evans@thyme>
3   Date: Thu, 7 Jun 2001 11:05:33 -0700 (PDT)
4   From: jeffrey.hammad@enron.com
5   To: andy.zipper@enron.com
6   Subject: Thanks for the interview
7   Mime-Version: 1.0
8   Content-Type: text/plain; charset=us-ascii
9   Content-Transfer-Encoding: 7bit
10  X-From: Hammad, Jeffrey </O=ENRON/OU=NA/CN=RECIPIENTS/CN=NOTESADDR/CN=CB
11  X-To: Zipper, Andy </O=ENRON/OU=NA/CN=RECIPIENTS/CN=AZIPPER>
12  X-cc:
13  X-bcc:
14  X-Folder: \Zipper, Andy\Zipper, Andy\Inbox
15  X-Origin: ZIPPER-A
16  X-FileName: Zipper, Andy.pst
17
18  Andy,
19
20  Thanks for giving me the opportunity to meet with you about the Analyst/
21
22  Thanks and Best Regards,
23
24  Jeff Hammad
25  """
26
27  html_msg_example = """
28  Message-ID: <21013632.1075862392611.JavaMail.evans@thyme>
29  Date: Mon, 19 Nov 2001 12:15:44 -0800 (PST)
30  From: insynconline.6jy5ympb.d@insync-palm.com
31  To: tstaab@enron.com
32  Subject: Last chance for special offer on Palm OS Upgrade!
33  Mime-Version: 1.0
34  Content-Type: text/plain; charset=us-ascii
35  Content-Transfer-Encoding: 7bit
36  X-From: InSync Online <InSyncOnline.6jy5ympb.d@insync-palm.com>
37  X-To: THERESA STAAB <tstaab@enron.com>
38  X-cc:
39  X-bcc:
40  X-Folder: \TSTAAB (Non-Privileged)\Staab, Theresa\Deleted Items
41  X-Origin: Staab-T
42  X-FileName: TSTAAB (Non-Privileged).pst
43
44  <html>
45
46  <html>
47  <head>
48  <title>Paprika</title>
49  <meta http-equiv="Content-Type" content="text/html;">
50  </head>
51  <body bgcolor="#FFFFFF" TEXT="#333333" LINK="#336699" VLINK="#6699cc" AL
52  <table border="0" cellpadding="0" cellspacing="0" width="582">
53  <tr valign="top">
54    <td width="582" colspan="9"><nobr><a href="http://insync-online.p04.co
55  </tr>
56  <tr valign="top">
```

```
57    <td width="4" bgcolor="#CCCCCC"><img src="http://images4.postdirect.co
58    <td width="20"><img src="http://images4.postdirect.com/master-images/4
59    <td width="165"><br><a href="http://insync-online.p04.com/u.d?LkReaQA5
60    <td width="20"><img src="http://images4.postdirect.com/master-images/4
61    <td width="165"><br><a href="http://insync-online.p04.com/u.d?BkReaQA5
62    <td width="20"><img src="http://images4.postdirect.com/master-images/4
63    <td width="165"><br><a href="http://insync-online.p04.com/u.d?JkReaQA5
64    <td width="19"><img src="http://images4.postdirect.com/master-images/4
65    <td width="4" bgcolor="#CCCCCC"><img src="http://images4.postdirect.co
66  </tr>
67  </table>
68  <table border="0" cellpadding="0" cellspacing="0" width="582">
69  <tr valign="top">
70    <td width="4" bgcolor="#CCCCCC"><img src="http://images4.postdirect.co
71    <td width="574"><br>
72      <table border="0" cellpadding="0" cellspacing="0" width="574" bgcolo
73      <tr>
74        <td width="50"><img src="http://images4.postdirect.com/master-imag
75        <td width="474"><font face="verdana, arial" size="-2"color="#00000
76          <br>
77          Dear THERESA,
78          <br><br>
79          Due to overwhelming demand for the Palm OS&#174; v4.1 Upgrade wi
80          extending the special offer of 25% off through November 30, 2001
81          increase the functionality of your Palm&#153; III, IIIx, IIIxe,
82          new Palm OS v4.1 through this extended special offer. You'll rec
83          <b>for just $29.95 when you use Promo Code <font color="#FF0000"
84          <b>$10 savings</b> off the list price.
85          <br><br>
86          <a href="http://insync-online.p04.com/u.d?NkReaQA5eczXRh=51">Cli
87          <br><br>
88          <a href="http://insync-online.p04.com/u.d?MkReaQA5eczXRm=61"><im
89          <br><br>
90          You can do a lot more with your Palm&#153; handheld when you upg
91          favorite features just got even better and there are some terrif
92          <br><br>
93          <LI> Handwrite notes and even draw pictures right on your Palm&#
94          <LI> Tap letters with your stylus and use Graffiti&#174; at the
95          <LI> Improved Date Book functionality lets you view, snooze or c
96          <LI> You can easily change time-zone settings</LI>
97
98          <br><br>
99          <a href="http://insync-online.p04.com/u.d?WkReaQA5eczXRb=71"><im
100         <br><br>
101         <LI> <nobr>Mask/unmask</nobr> private records or hide/unhide dir
102         <LI> Lock your device automatically at a designated time using t
103         <LI> Always remember your password with our new Hint feature*</L
104
105         <br><br>
106         <a href="http://insync-online.p04.com/u.d?VEReaQA5eczXRQ=81"><im
107         <br><br>
108         <LI> Use your GSM compatible mobile phone or modem to get online
109         <LI> Stay connected with email, instant messaging and text messa
110         <LI> Send applications or records through your cell phone to sch
111             important information to others</LI>
112
113         <br><br>
```

```
114          All this comes in a new operating system that can be yours for j
115          upgrade to the new Palm&#153; OS v4.1</a> and you'll also get th
116          <nobr>1-800-881-7256</nobr> to order via phone.
117          <br><br>
118          Sincerely,<br>
119          The Palm Team
120          <br><br>
121          P.S. Remember, this extended offer opportunity of 25% savings ab
122          and is only available through the Palm Store when you use Promo
123          <br><br>
124          <img src="http://images4.postdirect.com/master-images/404707/bot
125          <br><img src="http://images4.postdirect.com/master-images/404707
126          </font></td>
127        <td width="50"><img src="http://images4.postdirect.com/master-imag
128      </tr>
129      </table></td>
130      <td width="4" bgcolor="#CCCCCC"><img src="http://images4.postdirect.
131    </tr>
132    <tr>
133    <td colspan="3"><img src="http://images4.postdirect.com/master-images/
134    </tr>
135 </table>
136 <table border="0" cellpadding="0" cellspacing="0" width="582">
137   <tr>
138     <td width="54"><img src="http://images4.postdirect.com/master-images
139     <td width="474"><font face="arial, verdana" size="-2" color="#000000
140     * This feature is available on the Palm&#153; IIIx, Palm&#153; IIIxe
141     ** Note: To use the MIK functionality, you need either a Palm OS&#17
142     with  <nobr>built-in</nobr> modem or data capability that has either
143     are using a phone, you must have data services from your mobile serv
144     a list of tested and supported phones that you can use with the MIK.
145     <br><br>
146     ------------------<br>
147     To modify your profile or unsubscribe from Palm newsletters, <a href
148     Or, unsubscribe by replying to this message, with "unsubscribe" as t
149     <br><br>
150     ------------------<br>
151     Copyright&#169; 2001 Palm, Inc. Palm OS, Palm Computing, HandFAX, Ha
152     HotSync, iMessenger, MultiMail, Palm.Net, PalmConnect, PalmGlove, Pa
153     and the Palm Platform Compatible Logo are registered trademarks of P
154     AnyDay, EventClub, HandMAIL, the HotSync Logo, PalmGear, PalmGlove,
155     trade dress, PalmSource, Smartcode, and Simply Palm are trademarks o
156     product names may be trademarks or registered trademarks of their re
157     <img src="http://images4.postdirect.com/master-images/404707/clear.g
158     <td width="54"><img src="http://images4.postdirect.com/master-images
159   </tr>
160 </table><br><br><br><br>
161 <!-- The following image is included for message detection -->
162 <img src="http://p04.com/1x1.dyn" border="0" alt="" width="1" height="1"
163 <img src="http://p04.com/1x1.dyn?0vEGou8Ig30ba2L2bLn" width=1 height=1><
164 </html>
165
166 </html>
167 """
168 plain_msg_example = plain_msg_example.strip()
169 html_msg_example = html_msg_example.strip()
```

In [9]:
```python
def parse_html_payload(payload):
    """
    This function uses Beautiful Soup to read HTML data
    and return the text.  If the payload is plain text, then
    Beautiful Soup will return the original content
    """
    soup = BeautifulSoup(payload, 'html.parser')
    return str(soup.get_text()).encode('utf-8').decode('utf-8')

def parse_email(original_msg):

    if bool(BeautifulSoup(original_msg, "html.parser").find()):
        # convert HTML to plain email message
        original_msg = parse_html_payload(original_msg)

    result = {}
    email = Parser(policy=default).parsestr(original_msg) # email parser
    # get values for each field
    result['text']= email.get_content()
    result['payload'] = email.get_payload()
    for key, value in email.items():
        result[key] =  email[key]

    # Convert the dictionary to a tuple and return the tuple
    tuple_result = tuple([str(result.get(column, None)) for column in co]
    return ParsedEmail(*tuple_result)
```

In [10]:
```python
parsed_msg = parse_email(plain_msg_example)
print(parsed_msg.text)
```

Andy,

Thanks for giving me the opportunity to meet with you about the Analyst/ As
sociate program.  I enjoyed talking to you, and look forward to contributin
g to the success that the program has enjoyed.

Thanks and Best Regards,

Jeff Hammad

```
In [11]:    1  parsed_html_msg = parse_email(html_msg_example)
            2  print(parsed_html_msg.text)
```

        You can do a lot more with your Palm\u2122 handheld when you upgr
    ade to the Palm OS v4.1. All your
            favorite features just got even better and there are some terrifi
    c new additions:

     Handwrite notes and even draw pictures right on your Palm\u2122 handheld
     Tap letters with your stylus and use Graffiti◆ at the same time with th
    e enhanced onscreen keyboard
     Improved Date Book functionality lets you view, snooze or clear multiple
    alarms all with a single tap
     You can easily change time-zone settings

     Mask/unmask private records or hide/unhide directly within the applicati
    on
     Lock your device automatically at a designated time using the new Autolo
    cking feature

**Functions to read all example emails from path:**
**'.../dsc650/assignments/assignment04/examples'**

```
In [12]:    1  def get_email_paths(directory):
            2      pathlist = []
            3      for root, dirs, files in os.walk(directory):
            4          for file_path in files:
            5              ## Current path is now the file path to the current email.
            6              current_path = Path(root).joinpath(file_path)
            7              pathlist.append(str(current_path))
            8      return pathlist
```

```
In [13]:    1  def get_email_text(i):
            2      original_msg = read_raw_email(i)
            3      original_msg = original_msg.strip()
            4      parsed_msg = parse_email(original_msg)
            5      text = parsed_msg.text
            6      return text
```

```
In [14]:    1  examples_dir = current_dir.joinpath('examples')
            2  examples_dir.mkdir(parents=True, exist_ok=True)
```

In [15]: ►|
```python
1  pathlist = get_email_paths(examples_dir)
2  for item in pathlist:
3      text = get_email_text(item)
4      print(text)
```

```
4 ROUND-TRIP AIR TICKETS - ABSOLUTELY FREE
 to exciting destinations across the U.S., Caribbean,
 Hawaii or Mexico. Take the family or take two trips!
 And except for your hotel stay, there's NO OBLIGATION
 and NO PURCHASE REQUIRED for your trip. Imagine
 flying your family FREE to Florida, California,
 Las Vegas, Jamaica, Aruba, Bahamas, or other exciting
 vacation paradise! Plus, you get a 30-day FREE
 trial in Preferred Traveller, where you'll save big
 whenever and wherever you travel,
 dine or entertain!
Act now! This FREE offer is good only for a limited
 time!




 The 4 FREE AIRLINE TICKETS are yours to keep,
 just for registering! Plus, your FREE TICKETS
```

## Assignment 4.3

In [16]: ►|
```python
1  ## This creates a schema for the email data
2  email_struct = StructType()
3
4  for column in columns:
5      email_struct.add(column, StringType(), True)
```

In [35]:

```python
## This creates a user-defined function which can be used in Spark
parse_email_func = udf(lambda z: parse_email(z), email_struct)

def parse_emails(input_df):
    new_df = input_df.select(
        'username', 'id', 'original_msg', parse_email_func('original_msg
    )
    for column in columns:
        new_df = new_df.withColumn(column, new_df.parsed_email[column])

    new_df = new_df.drop('parsed_email')
    return new_df

class ParseEmailsTransformer(Transformer):
    def _transform(self, dataset):
        """
        Transforms the input dataset.

        :param dataset: input dataset, which is an instance of :py:class
        :returns: transformed dataset
        """
        return dataset.transform(parse_emails)

## Use the custom ParseEmailsTransformer, Tokenizer, and CountVectorizer
## to create a spark pipeline
email_pipeline = Pipeline(stages=[ParseEmailsTransformer(),
                                  Tokenizer(inputCol='text',
                                            outputCol='words'),
                                  CountVectorizer(inputCol='words',
                                                  outputCol='features',
                                                  vocabSize=3)])
```

In [37]:

```python
model = email_pipeline.fit(df)
```

In [38]:

```python
result = model.transform(df)
```

In [78]:  ▶|  | 1 | result.select('id', 'words', 'features').show()

```
+--------------------+--------------------+--------------------+
|                  id|               words|            features|
+--------------------+--------------------+--------------------+
|   davis-d\2_trash\1_|[, >, , , , , >, ...|(3,[0,1,2],[697.0...|
|   davis-d\2_trash\2_|[fyi..., thanks.,...|(3,[0,1,2],[57.0,...|
|   davis-d\2_trash\3_|[----------------...|(3,[0,1,2],[118.0...|
|   davis-d\2_trash\4_|[-----original, m...|(3,[0,2],[17.0,1.0])|
|davis-d\2_trash\c...|[hi, mommy!, , ye...|(3,[0,1,2],[4.0,1...|
|davis-d\2_trash\c...|[hey, sweetie,, ,...| (3,[0,1],[8.0,1.0])|
|davis-d\2_trash\c...|[----------------...|       (3,[0],[41.0])|
|davis-d\2_trash\c...|[----------------...|(3,[0,1,2],[65.0,...|
|davis-d\2_trash\c...|[----------------...|(3,[0,2],[58.0,1.0])|
|davis-d\2_trash\c...|[----------------...|(3,[0,1,2],[29.0,...|
|davis-d\2_trash\c...|[----------------...|(3,[0,1,2],[18.0,...|
|davis-d\2_trash\c...|[----------------...|       (3,[0],[43.0])|
|davis-d\2_trash\c...|[, , , , ---------...|       (3,[0],[11.0])|
|davis-d\2_trash\c...|[----------------...|       (3,[0],[14.0])|
|davis-d\2_trash\c...|[are, you, on, th...| (3,[0,1],[1.0,1.0])|
|davis-d\2_trash\c...|[listen, girly!, ...| (3,[0,2],[9.0,1.0])|
|davis-d\2_trash\c...|[candis, all, you...|(3,[0,1,2],[16.0,...|
|davis-d\2_trash\c...|[what, is, your, ...| (3,[0,2],[1.0,1.0])|
|davis-d\2_trash\c...|[candis, -, , why...|        (3,[0],[3.0])|
|davis-d\2_trash\c...|[----------------...|       (3,[0],[16.0])|
+--------------------+--------------------+--------------------+
only showing top 20 rows
```

In [ ]:  ▶|  | 1 |