

Code Highlights

Parse all sub-folders and files from a given folder:

- Universal Detector: Detects the encoding type of a file and returns the file text properly decoded. This is so helpful!
- `os.walk`

for root, dirs, files in `os.walk(enron_data_dir)`:

```
    for file_path in files:
```

- Parse email using the python email package: <https://docs.python.org/3/library/email.examples.html> (<https://docs.python.org/3/library/email.examples.html>)
- Create PySpark Dataframe and Schema:
 - <https://www.geeksforgeeks.org/how-to-create-pyspark-dataframe-with-schema/> (<https://www.geeksforgeeks.org/how-to-create-pyspark-dataframe-with-schema/>)
 - <https://sparkbyexamples.com/pyspark/pyspark-structtype-and-structfield/> (<https://sparkbyexamples.com/pyspark/pyspark-structtype-and-structfield/>)
- Proper Installation Notes to get Pyspark running in Jupyter.
 - <https://changhsinlee.com/install-pyspark-windows-jupyter/> (<https://changhsinlee.com/install-pyspark-windows-jupyter/>)

```
In [31]: 1 #!/*** Run this code may or may not ensure pyspark is happy
          2 """
          3 import findspark
          4 findspark.init()
          5
          6 import pyspark # only run after findspark.init()
          7 from pyspark.sql import SparkSession
          8 spark = SparkSession.builder.getOrCreate()
          9
         10 df = spark.sql(''select 'spark' as hello '')
         11 df.show()
         12 """
         13 print()
         14
         15
```

```
In [32]: 1 import os
          2 import json
          3 from pathlib import Path
          4 import zipfile
          5 import email
          6 from email.policy import default
          7 from email.parser import Parser
          8 from datetime import timezone
          9 import datetime
```

```
10
11 from collections import namedtuple
12
13 import pandas as pd
14 #import s3fs
15 from bs4 import BeautifulSoup
16 from dateutil.parser import parse
17 from chardet.universaldetector import UniversalDetector
18
19 #!/** Must Run before pyspark
20 import findspark
21 findspark.init()
22
23 from pyspark.ml import Pipeline
24 from pyspark.ml.feature import CountVectorizer
25 from pyspark.ml.feature import HashingTF, Tokenizer
26 from pyspark.sql import SparkSession
27 from pyspark.sql.functions import col
28 from pyspark.ml.pipeline import Transformer
29 from pyspark.sql.functions import udf
30 from pyspark.sql.types import *
31
32 import pandas as pd
33
34 current_dir = Path(os.getcwd()).absolute()
35 results_dir = current_dir.joinpath('results')
36 results_dir.mkdir(parents=True, exist_ok=True)
37 data_dir = current_dir.joinpath('data')
38 data_dir.mkdir(parents=True, exist_ok=True)
39 enron_data_dir = data_dir.joinpath('enron')
40
41 output_columns = [
42     'username',
43     'original_msg',
44     'payload',
45     'Message-ID',
46     'Date',
47     'From',
48     'To',
49     'Subject',
50     'Mime-Version',
51     'Content-Type',
52     'Content-Transfer-Encoding',
53     'X-From',
54     'X-To',
55     'X-cc',
56     'X-bcc',
57     'X-Folder',
58     'X-Origin',
59     'X-FileName',
60     'Cc',
61     'Bcc'
62 ]
63
64 columns = [column.replace('-', '_') for column in output_columns]
65
```

```

66 ParsedEmail = namedtuple('ParsedEmail', columns)
67
68 spark = SparkSession\
69     .builder\
70     .appName("Assignment04")\
71     .getOrCreate()
72
73

```

In []:

In []:

In [33]:

```

1 print("Current Dir: ", current_dir)
2 print("Results Dir: ", results_dir)
3 print("Data Dir: ", data_dir)
4 print("Enron Data Dir: ", enron_data_dir)
5 print(ParsedEmail.username)
6
7

```

```

Current Dir:  C:\Users\family\DSCProjects\DSC\DSC650\assignment04
Results Dir:  C:\Users\family\DSCProjects\DSC\DSC650\assignment04\results
Data Dir:     C:\Users\family\DSCProjects\DSC\DSC650\assignment04\data
Enron Data Dir:  C:\Users\family\DSCProjects\DSC\DSC650\assignment04\data\enron
<property object at 0x00000168A2252D18>

```

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

In [34]:

```

1 #!/*** Copied Files manually, avoided Amazon S3 due to ongoing S3 iss
2 """
3 def copy_data_to_local():
4     dst_data_path = data_dir.joinpath('enron.zip')
5     endpoint_url='https://storage.budsc.midwest-datascience.com'
6     enron_data_path = 'data/external/enron.zip'
7
8     s3 = s3fs.S3FileSystem(
9         anon=True,
10        client_kwargs={
11            'endpoint_url': endpoint_url
12        }
13    )
14
15
16    s3.get(enron_data_path, str(dst_data_path))
17
18    with zipfile.ZipFile(dst_data_path) as f_zip:
19        f_zip.extractall(path=data_dir)
20
21 copy_data_to_local()
22 """
23 print()
24

```

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

Assignment 4.1

```
In [35]: 1 #!/*** Use Universal Detector to ascertain the message encoding type.
2 #!/*** Returns a text based on the detected encoding type
3 def read_raw_email(email_path):
4     detector = UniversalDetector()
5
6     try:
7         with open(email_path) as f:
8             original_msg = f.read()
9     except UnicodeDecodeError:
10        detector.reset()
11        with open(email_path, 'rb') as f:
12            for line in f.readlines():
13                detector.feed(line)
14                if detector.done:
15                    break
16        detector.close()
17        encoding = detector.result['encoding']
18        with open(email_path, encoding=encoding) as f:
19            original_msg = f.read()
20
21    return original_msg
22
23 def make_spark_df():
24
25     #!/*** All Fields are Stringtype except Date which is Timestamp t
26     #!/*** PySpark will accept a Datetime Object for timestamp type
27     schema = StructType([
28         StructField("id", StringType(), True),
29         StructField("username", StringType(), True),
30         StructField("original_msg", StringType(), True)
31     ])
32     records = []
33     sc = spark.sparkContext
34
35     for root, dirs, files in os.walk(enron_data_dir):
36         for file_path in files:
37             ## Current path is now the file path to the current email
38             ## Use this path to read the following information
39             ## original_msg
40             ## username (Hint: It is the root folder)
41             ## id (The relative path of the email message)
42             current_path = Path(root).joinpath(file_path)
43
44             #!/*** Get raw Email Text Message from File
45             raw_email = read_raw_email(current_path)
46
47             row = []
48
49
```

```

50         #!/*** Append ID
51         id_path = str(current_path).replace(os.getcwd(), "").replace(
52         row.append(id_path)
53
54         #!/*****
55         #!/*** Find username
56         #!/*****
57
58
59         limit = 10
60         tgt = 'enron'
61         val = ""
62         i = -1
63         n = 0
64         #!/*** Username will be the folder name after enron
65         #!/*** Search through the Path Parent names to find the i
66         #!/*** Username will be parent[i-1].name
67         while val != tgt:
68             i += 1
69             val = str(Path(root).parents[i].name)
70
71             #!/*** Prevent Infinite Loops with a maximum Loop Lin
72             n += 1
73             if n > limit:
74                 #!/*** Limit Reached Set i to 2
75                 i=1
76                 break
77             #!/*** i can't be negative, reset to 0
78             if i == 0:
79                 i = 1
80
81         try:
82             username = str(Path(root).parents[i-1].name)
83         except:
84             print(i)
85         #!/*** Append Username
86         row.append(username)
87
88
89
90         #!/*** Add Original Message
91         row.append(raw_email)
92
93         #!/*** Add Row as Record
94         records.append(row)
95
96         #!/*** Print a Sample, every 200 records
97         if len(records) % 200 == 0:
98             print(f"username: {username} Path: {id_path} Msg Len:
99
100         ## TODO: Complete the code to code to create the Spark dataframe
101         return spark.createDataFrame(records,schema)
102
103 df = make_spark_df()
104
105 username: davis-d Path: \davis-d\all_documents\244_ Msg Len: 716

```

```
username: davis-d Path: \davis-d\all_documents\425_ Msg Len: 938
username: davis-d Path: \davis-d\deleted_items\101_ Msg Len: 1436
username: davis-d Path: \davis-d\deleted_items\296_ Msg Len: 2400
username: davis-d Path: \davis-d\discussion_threads\129_ Msg Len: 1000
username: davis-d Path: \davis-d\discussion_threads\30_ Msg Len: 698
username: davis-d Path: \davis-d\finanial_operations\14_ Msg Len: 1816
username: davis-d Path: \davis-d\inbox\family\13_ Msg Len: 1187
username: davis-d Path: \davis-d\sap\18_ Msg Len: 2219
username: davis-d Path: \davis-d\sent\64_ Msg Len: 1122
username: davis-d Path: \davis-d\_sent_mail\55_ Msg Len: 2223
username: gay-r Path: \gay-r\all_documents\238_ Msg Len: 3570
username: gay-r Path: \gay-r\all_documents\41_ Msg Len: 678
username: gay-r Path: \gay-r\discussion_threads\138_ Msg Len: 3172
username: gay-r Path: \gay-r\inbox\18_ Msg Len: 1792
username: gay-r Path: \gay-r\sent\210_ Msg Len: 2076
username: gay-r Path: \gay-r\sent\393_ Msg Len: 3247
username: gay-r Path: \gay-r\_sent_mail\40_ Msg Len: 3233
username: may-l Path: \may-l\all_documents\75_ Msg Len: 1647
username: may-l Path: \may-l\inbox\1009_ Msg Len: 1354
username: may-l Path: \may-l\inbox\196_ Msg Len: 2635
username: may-l Path: \may-l\inbox\380_ Msg Len: 36915
username: may-l Path: \may-l\inbox\565_ Msg Len: 5189
username: may-l Path: \may-l\inbox\746_ Msg Len: 6925
username: may-l Path: \may-l\inbox\929_ Msg Len: 1937
username: may-l Path: \may-l\sent_items\1_ Msg Len: 713
username: meyers-a Path: \meyers-a\deleted_items\1123_ Msg Len: 2359
username: meyers-a Path: \meyers-a\deleted_items\287_ Msg Len: 1716
username: meyers-a Path: \meyers-a\deleted_items\467_ Msg Len: 2675
username: meyers-a Path: \meyers-a\deleted_items\695_ Msg Len: 1605
username: meyers-a Path: \meyers-a\deleted_items\883_ Msg Len: 1955
username: mims-thurston-p Path: \mims-thurston-p\all_documents\133_ Ms
g Len: 717
username: mims-thurston-p Path: \mims-thurston-p\deleted_items\111_ Ms
g Len: 4454
username: mims-thurston-p Path: \mims-thurston-p\deleted_items\297_ Ms
g Len: 3746
username: mims-thurston-p Path: \mims-thurston-p\deleted_items\490_ Ms
g Len: 1207
username: mims-thurston-p Path: \mims-thurston-p\deleted_items\673_ Ms
g Len: 1436
username: mims-thurston-p Path: \mims-thurston-p\inbox\11_ Msg Len: 54
65
username: mims-thurston-p Path: \mims-thurston-p\inbox\320_ Msg Len: 2
230
username: mims-thurston-p Path: \mims-thurston-p\sent\167_ Msg Len: 50
4
username: mims-thurston-p Path: \mims-thurston-p\sent_items\159_ Msg L
en: 1793
username: mims-thurston-p Path: \mims-thurston-p\_sent_mail\138_ Msg L
en: 1283
username: mims-thurston-p Path: \mims-thurston-p\_sent_mail\99_ Msg Le
n: 677
username: quenet-j Path: \quenet-j\deleted_items\8_ Msg Len: 2588
username: reitmeyer-j Path: \reitmeyer-j\deleted_items\103_ Msg Len: 2
343
username: reitmeyer-j Path: \reitmeyer-j\inbox\207_ Msg Len: 889
```

```
username: reitmeyer-j Path: \reitmeyer-j\inbox\67_ Msg Len: 1577
username: sanchez-m Path: \sanchez-m\inbox\41_ Msg Len: 2489
username: shively-h Path: \shively-h\all_documents\144_ Msg Len: 458
username: shively-h Path: \shively-h\all_documents\55_ Msg Len: 542
username: shively-h Path: \shively-h\deleted_items\283_ Msg Len: 533
username: shively-h Path: \shively-h\deleted_items\507_ Msg Len: 991
username: shively-h Path: \shively-h\deleted_items\700_ Msg Len: 11895
username: shively-h Path: \shively-h\discussion_threads\2_ Msg Len: 10
23
username: shively-h Path: \shively-h\inbox\81_ Msg Len: 4439
username: shively-h Path: \shively-h\sent\238_ Msg Len: 934
username: shively-h Path: \shively-h\sent_items\53_ Msg Len: 534
username: shively-h Path: \shively-h\_sent_mail\231_ Msg Len: 794
username: staab-t Path: \staab-t\crestone\29_ Msg Len: 657
username: staab-t Path: \staab-t\deleted_items\275_ Msg Len: 1527
username: staab-t Path: \staab-t\mark_whitt\6_ Msg Len: 1194
username: zipper-a Path: \zipper-a\all_documents\45_ Msg Len: 694
username: zipper-a Path: \zipper-a\deleted_items\161_ Msg Len: 924
username: zipper-a Path: \zipper-a\deleted_items\352_ Msg Len: 1603
username: zipper-a Path: \zipper-a\deleted_items\540_ Msg Len: 1295
username: zipper-a Path: \zipper-a\inbox\114_ Msg Len: 803
username: zipper-a Path: \zipper-a\inbox\92_ Msg Len: 575
username: zipper-a Path: \zipper-a\sent_items\222_ Msg Len: 609
username: zipper-a Path: \zipper-a\sent items\89_ Msg Len: 998
```

In [36]:

```
1 print(df)
2 print(df.printSchema())
3 df.show(n=10,truncate = True)
4 ids = df.select("id").collect()
5 print(ids[0].id)
6 print(len(ids))
7
8
```

```
DataFrame[id: string, username: string, original_msg: string]  
root
```

Assignment 4.2

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

```
In [37]: 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=NOTESADMS/OU=ENRON/OU=NA/CN=RECIPIENTS/CN=AZIPPER>  
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 Analy  
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"
52 <table border="0" cellpadding="0" cellspacing="0" width="582">
53 <tr valign="top">
54 <td width="582" colspan="9"><noabr><a href="http://insync-online.p04
55 </tr>
56 <tr valign="top">
57 <td width="4" bgcolor="#CCCCCC"><br><a href="http://insync-online.p04.com/u.d?LkRea
60 <td width="20"><br><a href="http://insync-online.p04.com/u.d?BkRea
62 <td width="20"><br><a href="http://insync-online.p04.com/u.d?JkRea
64 <td width="19">
69 <tr valign="top">
70 <td width="4" bgcolor="#CCCCCC"><br>
72 <table border="0" cellpadding="0" cellspacing="0" width="574" bgc
73 <tr>
74 <td width="50"><font face="verdana, arial" size="-2"color="#00
76 <br>
77 Dear THERESA,
78 <br><br>
79 Due to overwhelming demand for the Palm OS&#174; v4.1 Upgrade
80 extending the special offer of 25% off through November 30, 2
81 increase the functionality of your Palm&#153; III, IIIx, IIIx
82 new Palm OS v4.1 through this extended special offer. You'll
83 <b>for just $29.95 when you use Promo Code <font color="#FF00
84 <b>$10 savings</b> off the list price.
85 <br><br>
86 <a href="http://insync-online.p04.com/u.d?NkReaQA5eczXRh=51">
87 <br><br>
88 <a href="http://insync-online.p04.com/u.d?MkReaQA5eczXRm=61">
89 <br><br>
90 You can do a lot more with your Palm&#153; handheld when you
91 favorite features just got even better and there are some tex
92 <br><br>
93 <LI> Handwrite notes and even draw pictures right on your Pal
94 <LI> Tap letters with your stylus and use Graffiti&#174; at t
95 <LI> Improved Date Book functionality lets you view, snooze 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">
100 <br><br>
101 <LI> <noabr>Mask/unmask</noabr> private records or hide/unhide
102 <LI> Lock your device automatically at a designated time usin
103 <LI> Always remember your password with our new Hint feature*

```

```

104
105         <br><br>
106         <a href="http://insync-online.p04.com/u.d?VEReaQA5eczXRQ=81">
107         <br><br>
108         <LI> Use your GSM compatible mobile phone or modem to get onl
109         <LI> Stay connected with email, instant messaging and text me
110         <LI> Send applications or records through your cell phone to
111             important information to others</LI>
112
113         <br><br>
114         All this comes in a new operating system that can be yours fo
115         upgrade to the new Palm&#153; OS v4.1</a> and you'll also get
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
122         and is only available through the Palm Store when you use Pro
123         <br><br>
124         </td>
127         <td width="50">
137         <tr>
138             <td width="54"><font face="arial, verdana" size="-2" color="#000
140             * This feature is available on the Palm&#153; IIx, Palm&#153; II
141             ** Note: To use the MIK functionality, you need either a Palm OS&
142             with <nobr>built-in</nobr> modem or data capability that has eit
143             are using a phone, you must have data services from your mobile s
144             a list of tested and supported phones that you can use with the M
145             <br><br>
146             -----<br>
147             To modify your profile or unsubscribe from Palm newsletters, <a h
148             Or, unsubscribe by replying to this message, with "unsubscribe" a
149             <br><br>
150             -----<br>
151             Copyright&#169; 2001 Palm, Inc. Palm OS, Palm Computing, HandFAX,
152             HotSync, iMessenger, MultiMail, Palm.Net, PalmConnect, PalmGlove,
153             and the Palm Platform Compatible Logo are registered trademarks o
154             AnyDay, EventClub, HandMAIL, the HotSync Logo, PalmGear, PalmGlov
155             trade dress, PalmSource, Smartcode, and Simply Palm are trademark
156             product names may be trademarks or registered trademarks of their
157             
165
166 </html>
167 """
168 plain_msg_example = plain_msg_example.strip()
169 html_msg_example = html_msg_example.strip()
170
171

```

```

In [38]: 1 def parse_html_payload(payload):
2         """
3         This function uses BeautifulSoup to read HTML data
4         and return the text. If the payload is plain text, then
5         BeautifulSoup will return the original content
6         """
7         soup = BeautifulSoup(payload, 'html.parser')
8         return str(soup.get_text()).encode('utf-8').decode('utf-8')
9
10 def parse_email(original_msg):
11     result = {}
12     msg = Parser(policy=default).parsestr(original_msg)
13     ## TODO: Use Python's email library to read the payload and the h
14     ## https://docs.python.org/3/library/email.examples.html
15
16     for column in output_columns:
17
18         #!/*** Username is the root parent folder name
19         if column == 'username':
20             continue
21
22         #!/*** Append the whole unprocessed text file
23         if column == 'original_msg':
24             result[column] = original_msg
25             continue
26
27         #!/*** Append the Payload, which is the message body
28         #!/*** Send raw Text to parse_html_payload. This will use bea
29         #!/*** html and return plain text
30         if column == 'payload':
31             #print('payload: ')
32             result[column] = parse_html_payload(msg.get_content())
33             continue
34
35         #!/*** Convert Date to Datetime
36         if column == 'Date':
37             #!/*** Convert Text to Datetime Object
38             # dt = datetime.datetime.strptime(msg[column], "%a, %d %b %
39
40             # result[column] = dt
41             # continue
42         #!/*** All other headers are pass-thru strings
43         if column in msg.keys():
44             #print(column, )

```

```

45         result[column] = msg[column]
46     else:
47
48         #!/*** Empty Field: return Zero Length String
49         result[column] = ""
50
51
52     tuple_result = tuple([str(result.get(column, None)) for column in
53     return ParsedEmail(*tuple_result)
54
55

```

```

In [39]: 1 msg = Parser(policy=default).parsestr(plain_msg_example)
         2 msg.keys()
         3
         4

```

```

Out[39]: ['Message-ID',
          'Date',
          'From',
          'To',
          'Subject',
          'Mime-Version',
          'Content-Type',
          'Content-Transfer-Encoding',
          'X-From',
          'X-To',
          'X-cc',
          'X-bcc',
          'X-Folder',
          'X-Origin',
          'X-FileName']

```

```

In [40]: 1 parsed_msg = parse_email(plain_msg_example)
         2 print(parsed_msg.payload)
         3
         4

```

Andy,

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

Thanks and Best Regards,

Jeff Hammad

```

In [41]: 1 parsed_html_msg = parse_email(html_msg_example)
         2 print(parsed_html_msg.payload)

```

Paprika

Dear THERESA,

Due to overwhelming demand for the Palm OS® v4.1 Upgrade with Mobile Connectivity, we are extending the special offer of 25% off through November 30, 2001. So there's still time to significantly increase the functionality of your Palm™ III, IIIx, IIIxe, IIIc, V or Vx handheld. Step up to the new Palm OS v4.1 through this extended special offer. You'll receive the brand new Palm OS v4.1 for just \$29.95 when you use Promo Code OS41WAVE. That's a \$10 savings off the list price.

[Click here to view a full product demo now.](#)

You can do a lot more with your Palm™ handheld when you upgrade to the Palm OS v4.1. All your favorite features just got even better and there are some terrific new additions:

Handwrite notes and even draw pictures right on your Palm™ handheld
Tap letters with your stylus and use Graffiti® at the same time with

the 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 application

Lock your device automatically at a designated time using the new Auto locking feature

Always remember your password with our new Hint feature*

Use your GSM compatible mobile phone or modem to get online and access the web

Stay connected with email, instant messaging and text messaging to GSM mobile phones

Send applications or records through your cell phone to schedule meetings and even "beam"

important information to others

All this comes in a new operating system that can be yours for just \$29.95! Click here to

upgrade to the new Palm™ OS v4.1 and you'll also get the latest Palm desktop software. Or call

1-800-881-7256 to order via phone.

Sincerely,

The Palm Team

P.S. Remember, this extended offer opportunity of 25% savings absolutely ends on November 30, 2001

and is only available through the Palm Store when you use Promo Code OS41WAVE.

* This feature is available on the Palm™ IIIx, Palm™ IIIxe, and Palm™ Vx.

** Note: To use the MIK functionality, you need either a Palm OS®

compatible modem or a phone
 with built-in modem or data capability that has either an infrared port or cable exits. If you
 are using a phone, you must have data services from your mobile service provider. Click here for
 a list of tested and supported phones that you can use with the MI K. Cable not provided.

 To modify your profile or unsubscribe from Palm newsletters, click here.

Or, unsubscribe by replying to this message, with "unsubscribe" as the subject line of the message.

 Copyright© 2001 Palm, Inc. Palm OS, Palm Computing, HandFAX, HandsTAMP, HandWEB, Graffiti,
 HotSync, iMessenger, MultiMail, Palm.Net, PalmConnect, PalmGlove, PalmModem, PalmPoint, PalmPrint,
 and the Palm Platform Compatible Logo are registered trademarks of Palm, Inc. Palm, the Palm logo,
 AnyDay, EventClub, HandMAIL, the HotSync Logo, PalmGear, PalmGlove, PalmPix, Palm Powered, the Palm
 trade dress, PalmSource, Smartcode, and Simply Palm are trademarks of Palm, Inc. All other brands and
 product names may be trademarks or registered trademarks of their respective owners.

Assignment 4.3

Notes:

<https://www.analyticsvidhya.com/blog/2019/11/build-machine-learning-pipelines-pyspark/>
[\(https://www.analyticsvidhya.com/blog/2019/11/build-machine-learning-pipelines-pyspark/\)](https://www.analyticsvidhya.com/blog/2019/11/build-machine-learning-pipelines-pyspark/)

<https://towardsdatascience.com/data-transformation-in-pyspark-6a88a6193d92>
[\(https://towardsdatascience.com/data-transformation-in-pyspark-6a88a6193d92\)](https://towardsdatascience.com/data-transformation-in-pyspark-6a88a6193d92)

Pyspark UDF - User Defined Functions:

<https://sparkbyexamples.com/pyspark/pyspark-udf-user-defined-function/>
[\(https://sparkbyexamples.com/pyspark/pyspark-udf-user-defined-function/\)](https://sparkbyexamples.com/pyspark/pyspark-udf-user-defined-function/)

```
In [42]: 1  ## This creates a schema for the email data
          2  email_struct = StructType()
          3
          4  #!/*** All Columns are StringType except the Date Column which is Tim
          5  for column in columns:
          6      #if column == "Date":
          7      #     email_struct.add(column, TimestampType(), True)
          8      #else:
          9      email_struct.add(column, StringType(), True)
```

```
In [43]: 1  ## This creates a user-defined function which can be used in Spark
```

```

2  ## It transforms the custom function parse_emails into a function use
3  parse_email_func = udf(lambda z: parse_email(z), email_struct)
4
5  #!/*** Takes the existing df, passes each column into parsed_email()
6  #!/*** This step parses the email header into separate columns
7  def parse_emails(input_df):
8      #!/*** Selects the first three columns of the input_df
9      #!/*** The fourth column applies the transformed parse_email_func
10     #!/*** and outputs the results in parsed_email column
11     new_df = input_df.select(
12         'username', 'id', 'original_msg', parse_email_func('original_
13     )
14
15     #!/*** Extracts each (sub)column in parsed_email into a df column
16     for column in columns:
17         new_df = new_df.withColumn(column, new_df.parsed_email[column
18
19     #!/*** Removed the parsed_email column filled with sub columns
20     new_df = new_df.drop('parsed_email')
21
22     return new_df
23
24 #!/*** Transformer function used to apply parse_emails, which in turn
25 class ParseEmailsTransformer(Transformer):
26     def transform(self, dataset):
27         """
28         Transforms the input dataset.
29
30         :param dataset: input dataset, which is an instance of :py:cl
31         :returns: transformed dataset
32         """
33         return dataset.transform(parse_emails)
34
35
36 ## Use the custom ParseEmailsTransformer, Tokenizer, and CountVectori
37 ## to create a spark pipeline
38 email_pipeline = Pipeline(
39     ## TODO: Complete code
40     stages=[
41         ParseEmailsTransformer(),
42         Tokenizer(inputCol='payload', outputCol="words"),
43         CountVectorizer(inputCol='words', outputCol='features')
44     ]
45 )
46
47 model = email_pipeline.fit(df)
48 result = model.transform(df)
49
50

```

```
In [44]: 1 result.select('id', 'words', 'features').show(n=20, truncate=True)
```

```
2
3
```



```

+-----+-----+-----+
|          id|          words|          features|
+-----+-----+-----+
| \davis-d\2_trash\1_|[, >, , , , , >, ...|(99771,[0,1,2,3,4...|
| \davis-d\2_trash\2_|[fyi..., thanks.,...|(99771,[0,1,2,3,5...|
| \davis-d\2_trash\3_|[-----...|(99771,[0,1,2,6,7...|
| \davis-d\2_trash\4_|[-----original, m...|(99771,[0,2,6,7,9...|
| \davis-d\2_trash\...|[hi, mommy!, , ye...|(99771,[0,1,2,6,7...|
| \davis-d\2_trash\...|[hey, sweetie,, ,...|(99771,[0,1,7,10,...|
| \davis-d\2_trash\...|[-----...|(99771,[0,10,25,2...|
| \davis-d\2_trash\...|[-----...|(99771,[0,1,2,3,5...|
| \davis-d\2_trash\...|[-----...|(99771,[0,2,3,6,7...|
| \davis-d\2_trash\...|[-----...|(99771,[0,1,2,3,7...|
| \davis-d\2_trash\...|[-----...|(99771,[0,1,2,3,4...|
| \davis-d\2_trash\...|[-----...|(99771,[0,10,25,2...|
| \davis-d\2_trash\...|[, , , , -----...|(99771,[0,7,10,12...|
| \davis-d\2_trash\...|[-----...|(99771,[0,10,14,2...|
| \davis-d\2_trash\...|[are, you, on, th...|(99771,[0,1,7,10,...|
| \davis-d\2_trash\...|[listen, girly!, ...|(99771,[0,2,3,7,1...|
| \davis-d\2_trash\...|[candis, all, you...|(99771,[0,1,2,3,7...|
| \davis-d\2_trash\...|[what, is, your, ...|(99771,[0,2,7,11,...|
| \davis-d\2_trash\...|[candis, -, , why...|(99771,[0,7,12,14...|
| \davis-d\2_trash\...|[-----...|(99771,[0,4,10,13...|

```