

CS4234/5234 Project Report

Question 1:

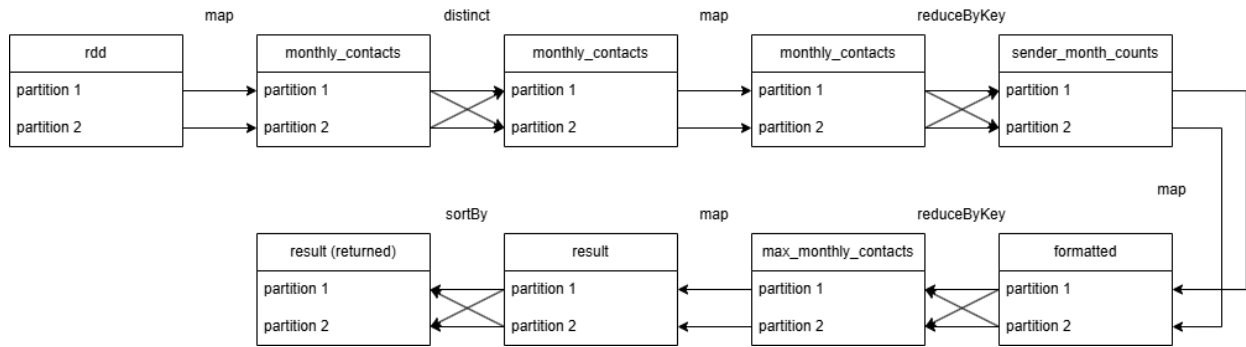
```
# Q5: replace pass with your code
def get_monthly_contacts(rdd):
    monthly_contacts = rdd.map(lambda x: (x[0], f"{x[2].month}/{x[2].year}", x[1]))
    sender_month_counts = monthly_contacts.distinct().map(lambda x: ((x[0], x[1]), 1))\
        .reduceByKey(lambda a, b: a + b)
    formatted = sender_month_counts.map(lambda x: (x[0][0], (x[0][1], x[1])))
    max_monthly_contacts = formatted.reduceByKey(lambda a, b: a if a[1] >= b[1] else b)
    result = max_monthly_contacts.map(lambda x: (x[0], x[1][0], x[1][1]))
    return result.sortBy(lambda x: (-x[2], x[0]))
```

The function `get_monthly_contacts()` takes an RDD as input and computes, for each sender, the month where they contacted the highest number of unique recipients, the function follows these steps.

1. Extract month and year: the timestamp of each email is converted to the MM/YYYY format
2. Ensure unique contacts: the function ensures that each recipient is counted only once per sender per month.
3. Count unique contacts per month: using `reduceByKey()`, it calculates the number of unique recipients each sender contacted per month
4. Identify the peak month: For each sender, the function determines the month with the highest number of unique contacts
5. Sorting: the results are sorted by the number of unique contacts (descending) and then by sender email (ascending)

Challenges:

- **Getting the correct date format:** unlike question 1 where it was easy to use the provided `date_to_dt()` function in this question I only need the month and year. If I wanted to use the provided function then I would have to cut it up, instead I decided to get the date different and instead just use the map function instead, and it would create cleaner code.
- **Getting the unique recipient counts per sender per month:** this is something that stumped me for a while, at first, I tried to do it all in one trying to think of different way of extracting the information I wanted. In the end what caused this to be so much of a challenge was not thinking to use `.distinct()`, which I eventually did get after going back to my notes.



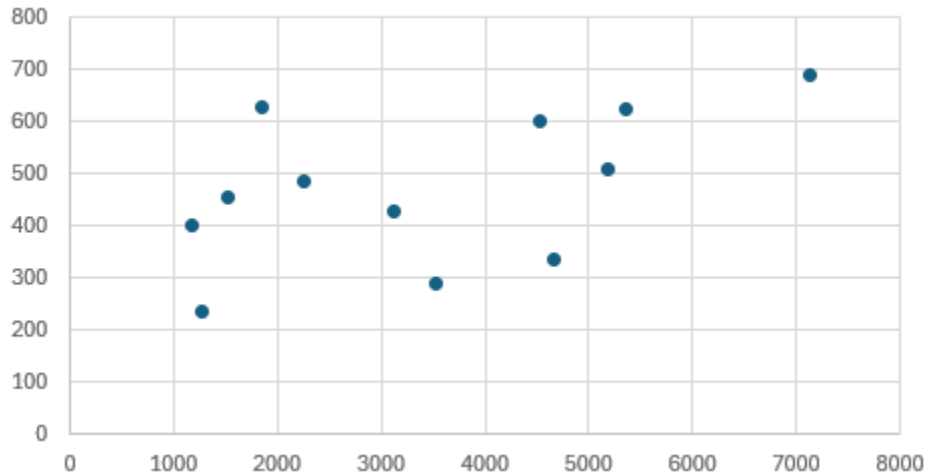
Narrow Dependencies: Map and Distinct

Wide Dependencies: Distinct, reduceByKey, and sortBy

Question 2:

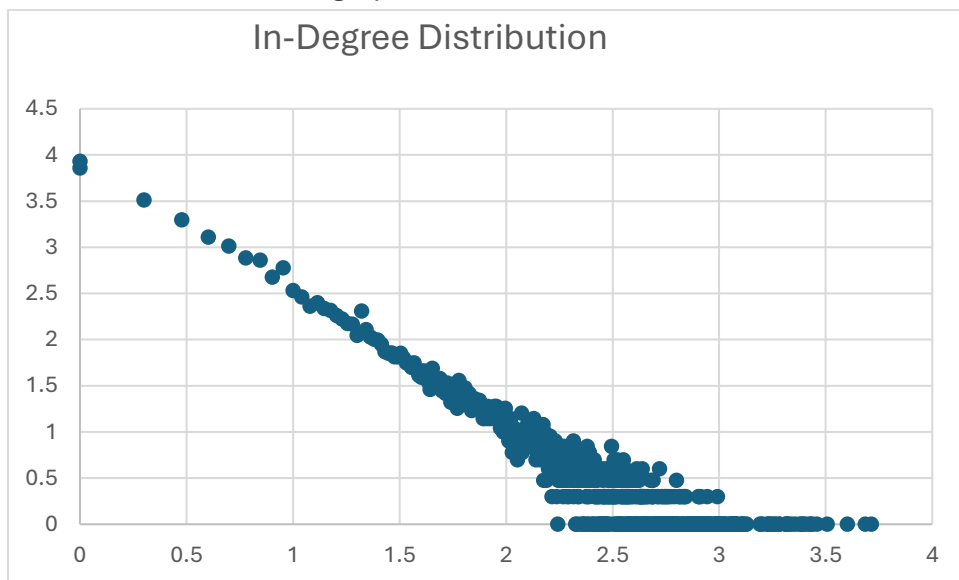
Note The code used to compute these answers is at the bottom of the report

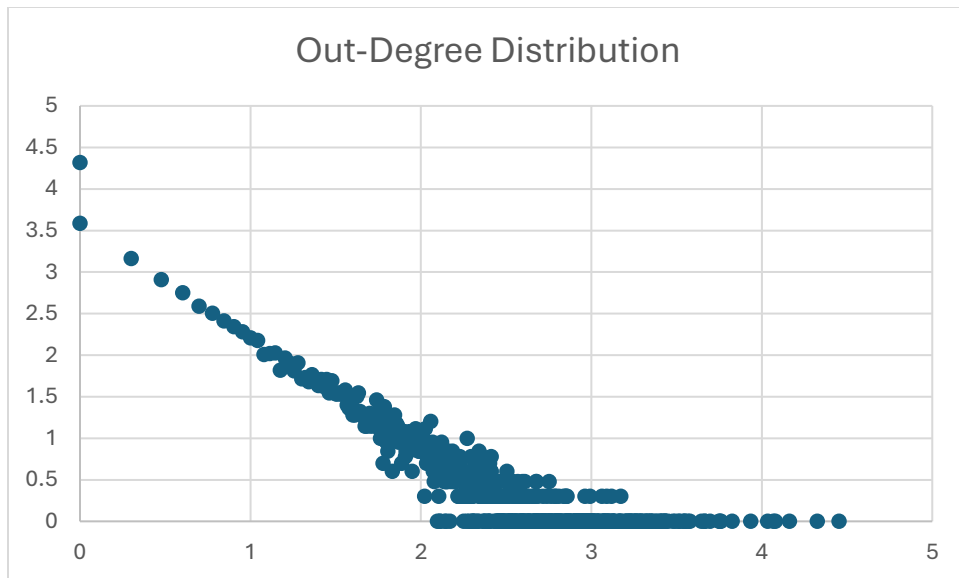
1. To Evaluate the 80/20 rule, I extracted the weighted network for the selected period using `conver_to_weighted_network()`. Then I computed the weighted in-degree and out-degree of each email using my answers from `project.py`. My calculations indicate that these high-degree nodes account for 99.08% of outgoing communications and 92.64% of incoming communications, confirming the presence of a small number of nodes that make up a large amount of the network traffic.
2. To examine the “the rich-get-richer” effect, I tracked the maximum node degree (k_{max}) over times, progressively adding monthly data slices and computing the corresponding number of nodes and edges. Using the functions `get_in_degrees()` and `get_out_degrees()` from `project.py`, I observed that k_{max} increases as more months were included, the results where: (3532, 290), (4660, 335), (5184, 510), (7137, 690), (4523, 600), (1515, 454), (1851, 628), (1274, 236), (1167, 402), (5362, 626), (2266, 487), (3130, 430). When I plotted this in Excel:



The relationship between k_{max} and the number of nodes appears approximately linear, supporting the idea that as the network expands, the most connected nodes continue accumulating an increasing share of traffic.

- To test the hypothesis, I generated in-degree and out-degree distributions and then printed them into .csv files, these results were then plotted using Excel. I then plotted these on a log-log graph, where a linear trend would indicate power-law behavior. These are the graphs:





The results show a clear linear trend line, this confirms that the Enron email network follows a power-law degree distribution. This finding validates the hypothesis that the Enron email network exhibits the self-organizing properties characteristic of scale-free structures.

Question 2 Code:

```
1 from pyspark import SparkContext, SparkConf
2 from pyspark.sql import SparkSession
3 from datetime import datetime, timezone
4 import numpy as np
5 from project import extract_email_network
6 import csv
7
8 conf = SparkConf().setAppName("Enron Network Analysis")
9 sc = SparkContext(conf=conf)
10 spark = SparkSession(sc)
11
12 def utf8_decode_and_filter(rdd):
13     def utf_decode(s):
14         try:
15             return str(s, 'utf-8')
16         except:
17             pass
18     return rdd.map(lambda x: utf_decode(x[1])).filter(lambda x: x != None)
19
20
21 def convert_to_weighted_network(rdd, drange=None):
22     def filter_by_date(record):
23         sender, recipient, timestamp = record
24         if drange is None:
25             return True
26         return drange[0] <= timestamp <= drange[1]
27
28     return (rdd.filter(filter_by_date)
29             .map(lambda x: ((x[0], x[1]), 1))
30             .reduceByKey(lambda a, b: a + b)
31             .map(lambda x: (x[0][0], x[0][1], x[1])))
32
33
34 def get_out_degrees(rdd):
35     print("Q3.1")
36     all_nodes = rdd.flatMap(lambda x: [x[0], x[1]]).distinct()
37     out_degrees = rdd.map(lambda x: (x[0], x[2])).reduceByKey(lambda a, b: a + b)
38     final_rdd = all_nodes.map(lambda n: (n, 0)).leftOuterJoin(out_degrees)\
39         .map(lambda x: (x[1][1] if x[1][1] else 0, x[0]))\
40         .sortBy(lambda x: (-x[0], x[1]))
41
42     return final_rdd
43
44 def get_in_degrees(rdd):
45     all_nodes = rdd.flatMap(lambda x: [x[0], x[1]]).distinct()
46     in_degrees = rdd.map(lambda x: (x[1], x[2])).reduceByKey(lambda a, b: a + b)
47     final_rdd = all_nodes.map(lambda n: (n, 0)).leftOuterJoin(in_degrees)\
48         .map(lambda x: (x[1][1] if x[1][1] else 0, x[0]))\
49         .sortBy(lambda x: (-x[0], x[1]))
50
51     return final_rdd
52
53
54 rdd = utf8_decode_and_filter(sc.sequenceFile("/user/ufac001/project/enron-full"))
55
56 email_rdd = extract_email_network(rdd)
57
58 date_range = (datetime(2001, 1, 1, tzinfo=timezone.utc), datetime(2001, 12, 31, tzinfo=timezone.utc))
59
60 weighted_rdd = convert_to_weighted_network(email_rdd, date_range)
61
62 out_degrees = get_out_degrees(weighted_rdd).collect()
63 in_degrees = get_in_degrees(weighted_rdd).collect()
64
65 total_out_degree = sum(d for d, _ in out_degrees)
66 total_in_degree = sum(d for d, _ in in_degrees)
67
68 top_20_percent = int(len(out_degrees) * 0.2) or 1
69 top_out_degree = sum(d for d, _ in sorted(out_degrees, reverse=True)[:top_20_percent])
70 top_in_degree = sum(d for d, _ in sorted(in_degrees, reverse=True)[:top_20_percent])
```

```

70 k_max_values = []
71 node_counts = []
72 for n in range(1, 13):
73     month_range = (datetime(2001, n, 1, tzinfo=timezone.utc), datetime(2001, n, 28, tzinfo=timezone.utc))
74     monthly_rdd = convert_to_weighted_network(email_rdd, month_range)
75     max_out = max(get_out_degrees(monthly_rdd).collect(), key=lambda x: int(x[0]) if isinstance(x[0], str) else x[0], default=(0, ""))[0]
76     max_in = max(get_in_degrees(monthly_rdd).collect(), key=lambda x: int(x[0]) if isinstance(x[0], str) else x[0], default=(0, ""))[0]
77     k_max_values.append((max_out, max_in))
78     node_counts.append(monthly_rdd.count())
79
80 print(f"top_out_degree = {top_out_degree}")
81 print(f"total out degree = {total_out_degree}")
82 print(f"top in degree = {top_in_degree}")
83 print(f"total in degree = {total_in_degree}")
84 print("80/20 Rule Compliance (Out-degree):", top_out_degree / total_out_degree if total_out_degree else 0)
85 print("80/20 Rule Compliance (In-degree):", top_in_degree / total_in_degree if total_in_degree else 0)
86 print("Max Degree (k_max) Over Time:", k_max_values)
87
88 def get_out_degree_dist(rdd):
89     out_degrees = get_out_degrees(rdd)
90     degree_dist = out_degrees.map(lambda x: (x[0], 1)).reduceByKey(lambda a, b: a + b)
91     return degree_dist.sortBy(lambda x: x[0])
92
93 def get_in_degree_dist(rdd):
94     in_degrees = get_in_degrees(rdd)
95     degree_dist = in_degrees.map(lambda x: (x[0], 1)).reduceByKey(lambda a, b: a + b)
96     return degree_dist.sortBy(lambda x: x[0])
97
98 out_degrees = get_out_degrees(weighted_rdd)
99 in_degrees = get_in_degrees(weighted_rdd)
100
101 out_degree_dist = get_out_degree_dist(weighted_rdd).collect()
102 in_degree_dist = get_in_degree_dist(weighted_rdd).collect()
103
104 with open("out_degree_results.csv", mode='a', newline='') as write_file:
105     writer = csv.writer(write_file)
106     for entry in out_degree_dist:
107         writer.writerow(entry)
108
109 with open("in_degree_results.csv", mode='a', newline='') as write_file:
110     writer = csv.writer(write_file)
111     for entry in in_degree_dist:
112         writer.writerow(entry)
113
114
115 print(in_degree_dist)
116
117 sc.stop()

```

Note: I am happy to provide the original file, or data files for part 3