# PREDICT 420

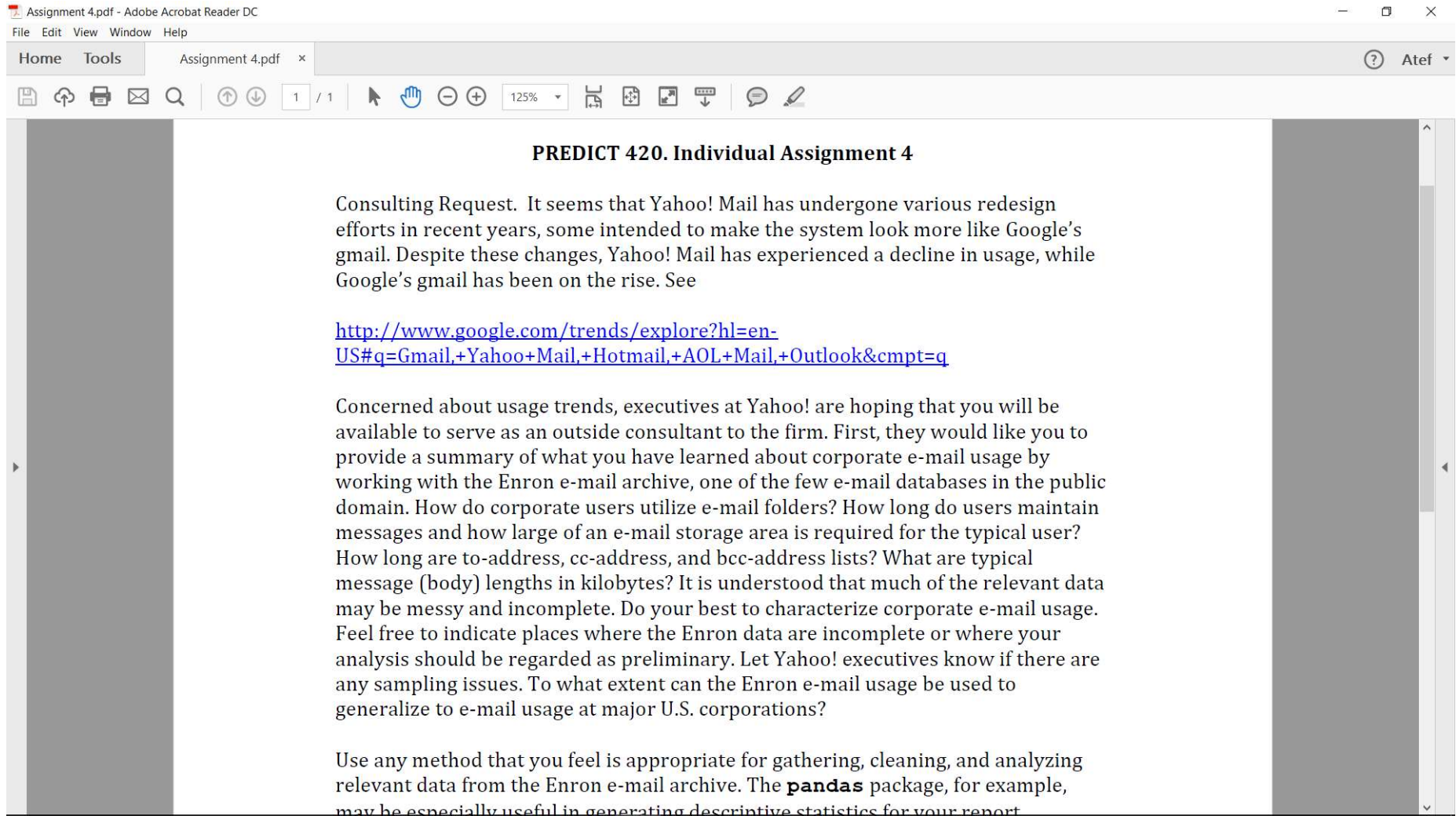# Atef Bader, PhD

# Agenda

- Assignment #4 - Walkthrough
    - Due at the end of this week (week #6)

- Exercise #6 - Walkthrough
    - Due at the end of this week (week #6)

# Assignment #4 - Deliverable

**PREDICT 420. Individual Assignment 4**

Consulting Request. It seems that Yahoo! Mail has undergone various redesign efforts in recent years, some intended to make the system look more like Google's gmail. Despite these changes, Yahoo! Mail has experienced a decline in usage, while Google's gmail has been on the rise. See

http://www.google.com/trends/explore?hl=en-
US#q=Gmail,+Yahoo+Mail,+Hotmail,+AOL+Mail,+Outlook&cmpt=q

Concerned about usage trends, executives at Yahoo! are hoping that you will be available to serve as an outside consultant to the firm. First, they would like you to provide a summary of what you have learned about corporate e-mail usage by working with the Enron e-mail archive, one of the few e-mail databases in the public domain. How do corporate users utilize e-mail folders? How long do users maintain messages and how large of an e-mail storage area is required for the typical user? How long are to-address, cc-address, and bcc-address lists? What are typical message (body) lengths in kilobytes? It is understood that much of the relevant data may be messy and incomplete. Do your best to characterize corporate e-mail usage. Feel free to indicate places where the Enron data are incomplete or where your analysis should be regarded as preliminary. Let Yahoo! executives know if there are any sampling issues. To what extent can the Enron e-mail usage be used to generalize to e-mail usage at major U.S. corporations?

Use any method that you feel is appropriate for gathering, cleaning, and analyzing relevant data from the Enron e-mail archive. The **pandas** package, for example, may be especially useful in generating descriptive statistics for your report.

# Assignment #4 - Deliverable

# Assignment #4 - Deliverable

# Assignment #4 - Deliverable

# Enron Email Dataset

This dataset was collected and prepared by the CALO Project (A Cognitive Assistant that Learns and Organizes). It contains data from about 150 users, mostly senior management of Enron, organized into folders. The corpus contains a total of about 0.5M messages. This data was originally made public, and posted to the web, by the Federal Energy Regulatory Commission during its investigation.

The email dataset was later purchased by Leslie Kaelbling at MIT, and turned out to have a number of integrity problems. A number of folks at SRI, notably Melinda Gervasio, worked hard to correct these problems, and it is thanks to them (not me) that the dataset is available. The dataset here does not include attachments, and some messages have been deleted "as part of a redaction effort due to requests from affected employees". Invalid email addresses were converted to something of the form user@enron.com whenever possible (i.e., recipient is specified in some parse-able format like "Doe, John" or "Mary K. Smith") and to no_address@enron.com when no recipient was specified.

I get a number of questions about this corpus each week, which I am unable to answer, mostly because they deal with preparation issues and such that I just don't know about. If you ask me a question and I don't answer, please don't feel slighted.

I am distributing this dataset as a resource for researchers who are interested in improving current email tools, or understanding how email is currently used. This data is valuable; to my knowledge it is the only substantial collection of

# Assignment #4 - Deliverable

# Assignment #4 - Deliverable

File  Edit  View  History  Bookmarks  Tools  Help

Enron Email Dataset

https://www.cs.cmu.edu/~./enron/

Search

## Enron Email Dataset

This dataset was collected and prepared by the CALO Project (A Cognitive Assistant that Learns and Organizes). It contains data from about 150 users, mostly senior management of Enron, organized into folders. The corpus contains a total of about 0.5M messages. This data was originally made public, and posted to the web, by the Federal Energy Regulatory Commission during its investigation.
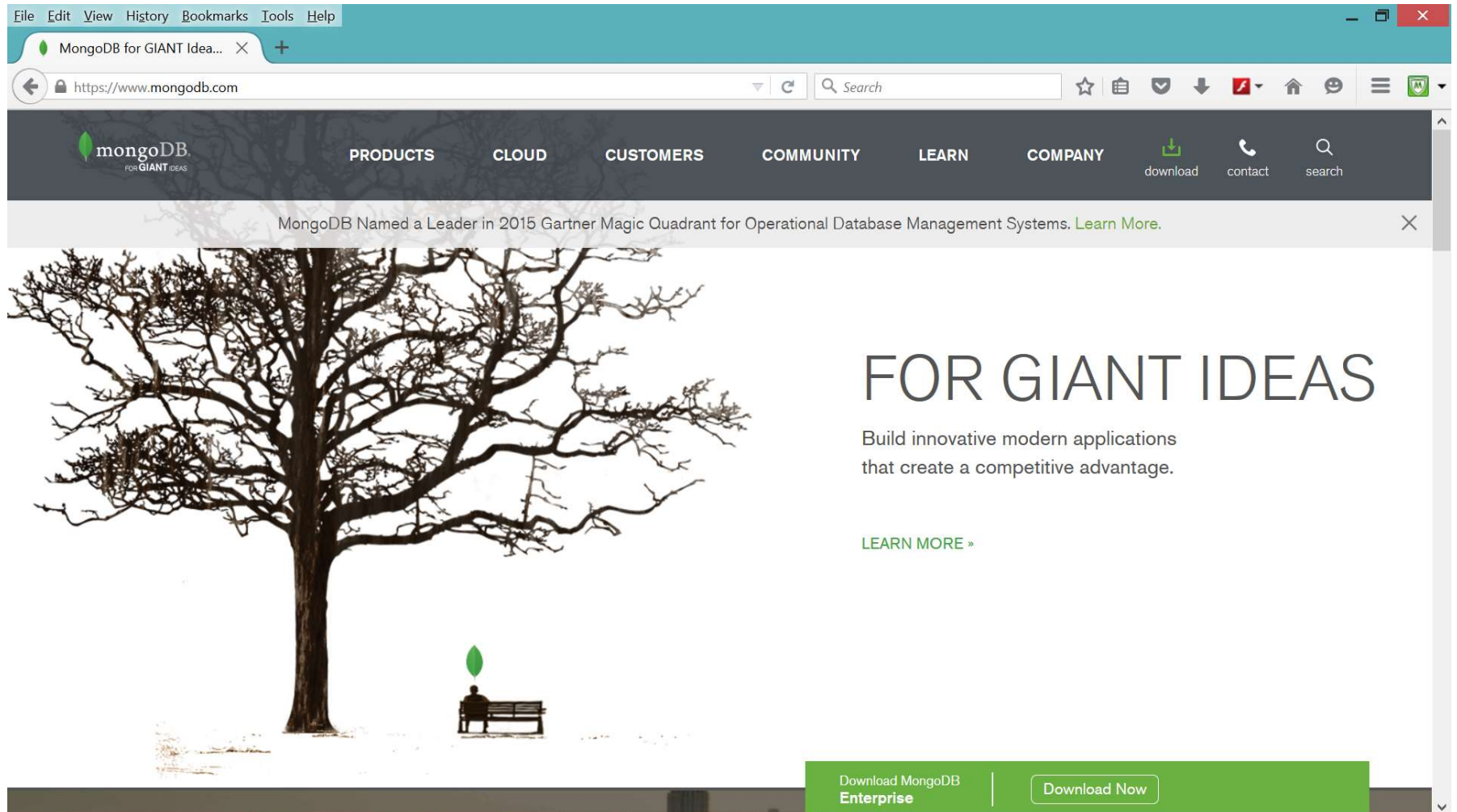
The email dataset was later purchased by Leslie Kaelbling at MIT, and turned out to have a number of integrity problems. A number of folks at SRI, notably Melinda Gervasio, worked hard to correct these problems, and it is thanks to them (not me) that the dataset is available. The dataset here does not include attachments, and some messages have been deleted "as part of a redaction effort due to requests from affected employees". Invalid email addresses were converted to something of the form user@enron.com whenever possible (i.e., recipient is specified in some parse-able format like "Doe, John" or "Mary K. Smith") and to no_address@enron.com when no recipient was specified.

I get a number of questions about this corpus each week, which I am unable to answer, mostly because they deal with preparation issues and such that I just don't know about. If you ask me a question and I don't answer, please don't feel slighted.
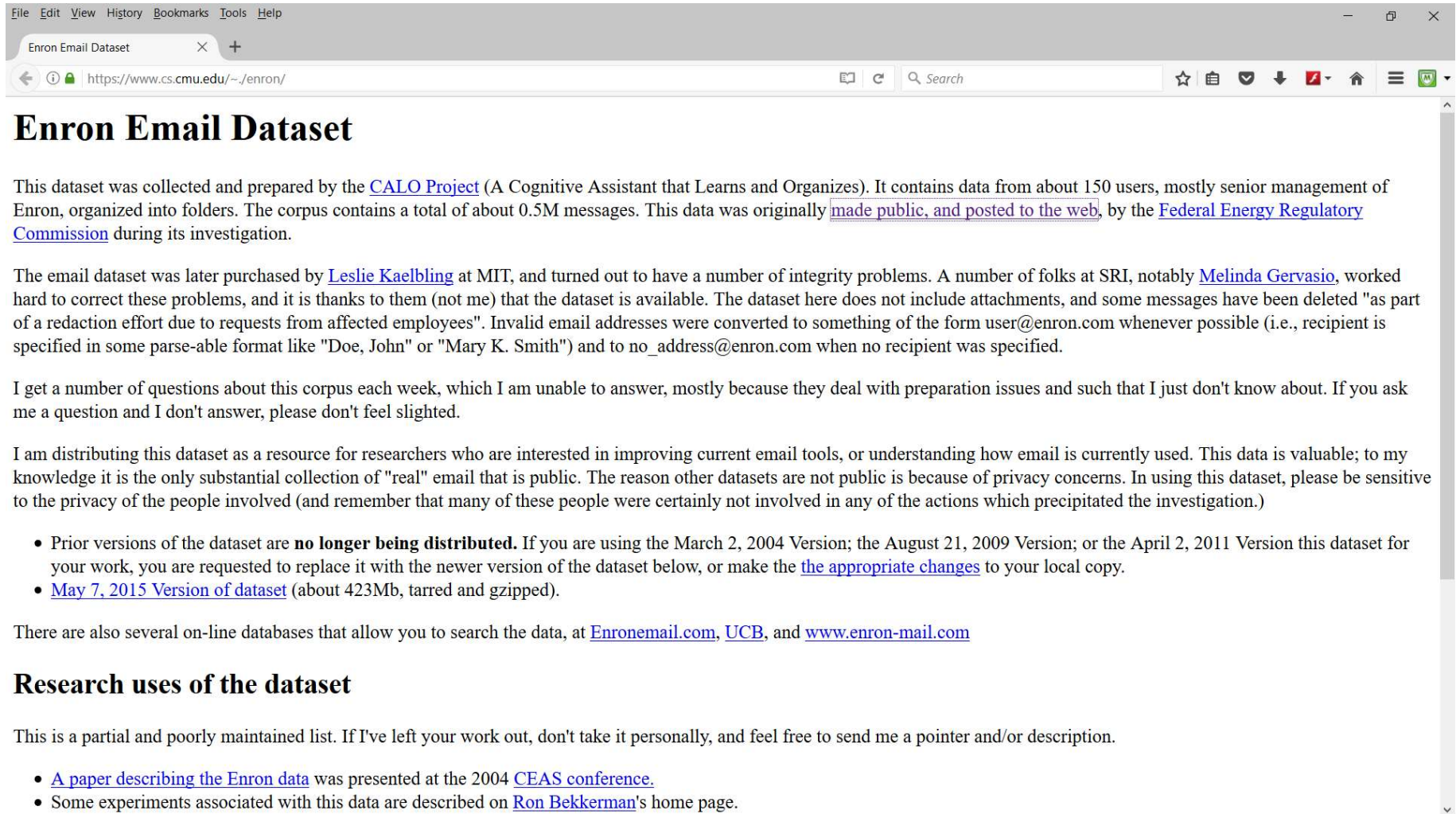
I am distributing this dataset as a resource for researchers who are interested in improving current email tools, or understanding how email is currently used. This data is valuable; to my knowledge it is the only substantial collection of "real" email that is public. The reason other datasets are not public is because of privacy concerns. In using this dataset, please be sensitive to the privacy of the people involved (and remember that many of these people were certainly not involved in any of the actions which precipitated the investigation.)

- Prior versions of the dataset are **no longer being distributed.** If you are using the March 2, 2004 Version; the August 21, 2009 Version; or the April 2, 2011 Version this dataset for your work, you are requested to replace it with the newer version of the dataset below, or make the the appropriate changes to your local copy.
- May 7, 2015 Version of dataset (about 423Mb, tarred and gzipped).

There are also several on-line databases that allow you to search the data, at Enronemail.com, UCB, and www.enron-mail.com

## Research uses of the dataset

This is a partial and poorly maintained list. If I've left your work out, don't take it personally, and feel free to send me a pointer and/or description.

- A paper describing the Enron data was presented at the 2004 CEAS conference.
- Some experiments associated with this data are described on Ron Bekkerman's home page.

# Assignment #4 - Deliverable

- Submit your log (script) file that has the commands you executed and the output you got

# Exercises #6 - Deliverable

- ❑     Submit your Ipython Notebook Script
- ❑     Show the Python code snippet you wrote