

# Day 4: Text Analysis





# Kenneth Lay



- lay-k/
  - \_sent/
  - all\_documents/
  - business/
  - deleted\_items/
  - family/
  - inbox/
  - sec\_panel/
  - sent/

- lay-k/

\_sent/

all\_documents/

business/



Email1

deleted\_items/

Email2

family/

....

inbox/

sec\_panel/

sent/

2

- Relevant terms

- Summaries
- Searching



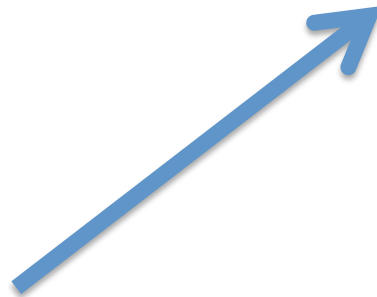
“Conference”  
“Call”  
“Saturday”  
“Noon”



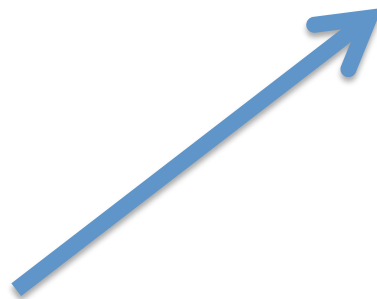
- Relevant terms
  - Summaries
  - Searching
- Similarity



- Relevant terms
  - Summaries
  - Searching
- Similarity



- Relevant terms
  - Summaries
  - Searching
- Similarity



# Relevant Terms

using TF-IDF

# Intuition

- Count # times each word is used

TF

We are in the process of trying to arrange a conference call with you on either Tuesday or Wednesday of next week to discuss the paper which is attached

We will be doing this by conference call and once we set a time to talk with you, will give you the number to call.

# Intuition

- Count # times each word is used

TF

We are in the process of trying to arrange a **conference** call with you on either Tuesday or Wednesday of next week to discuss the paper which is attached.

We will be doing this by **conference** call and once we set a time to talk with you, will give you the number to call.

**conference: 2**

# Intuition

- Count # times each word is used

TF

We are in the process of trying **to** arrange a conference call with you on either Tuesday or Wednesday of next week **to** discuss the paper which is attached.

We will be doing this by conference call and once we set a time **to** talk with you, will give you the number **to** call.

**to: 4**

# Intuition

- Count # times each word is used **TF**
- Penalize if most documents use word **IDF**

We are in the process of trying **to** arrange a conference call with you on either Tuesday or Wednesday of next week **to** discuss the paper which is attached.

We will be doing this by conference call and once we set a time **to** talk with you, will give you the number **to** call.

**to: 4**  Want to penalize



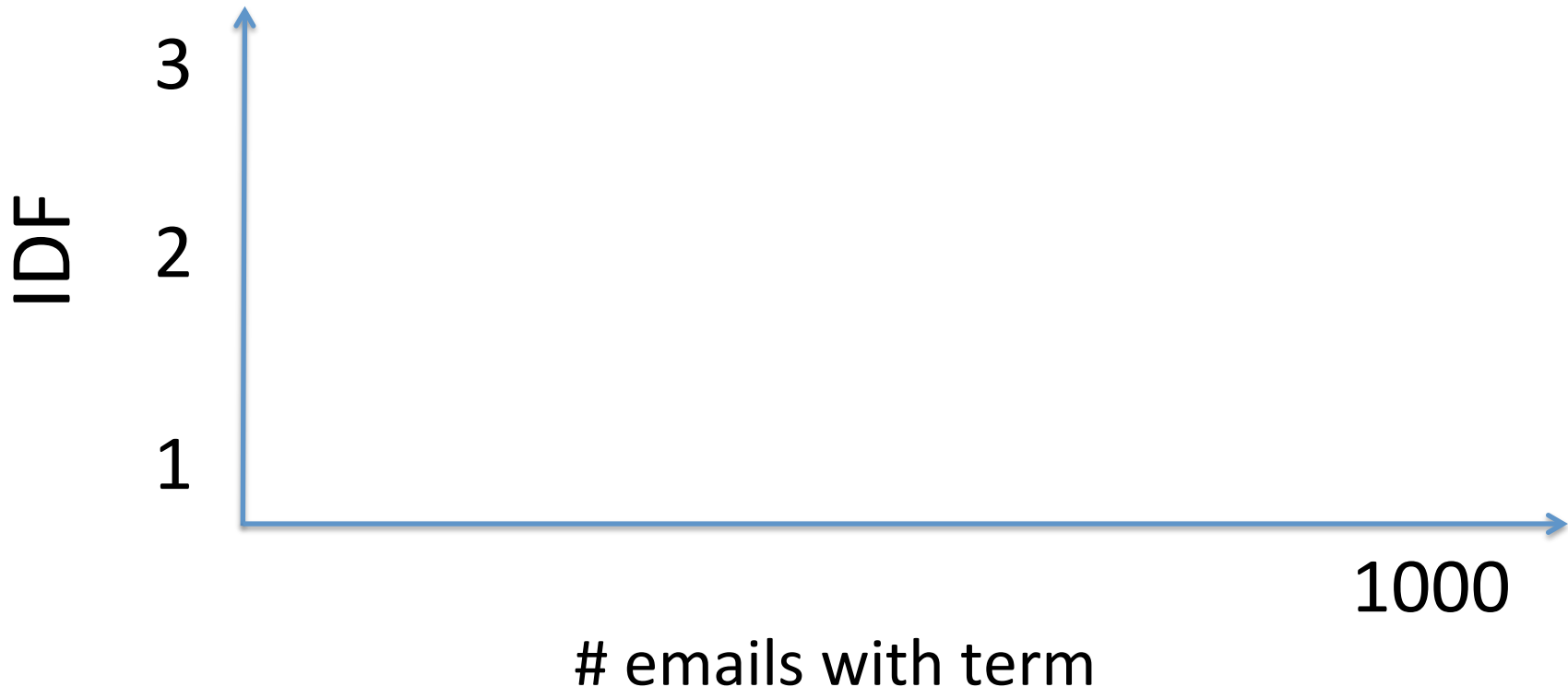
# IDF. How do they work?

$\log(\text{total \# emails} / \text{\# emails with word})$

# IDF. How do they work?

$\log(\text{total \# emails} / \text{\# emails with word})$

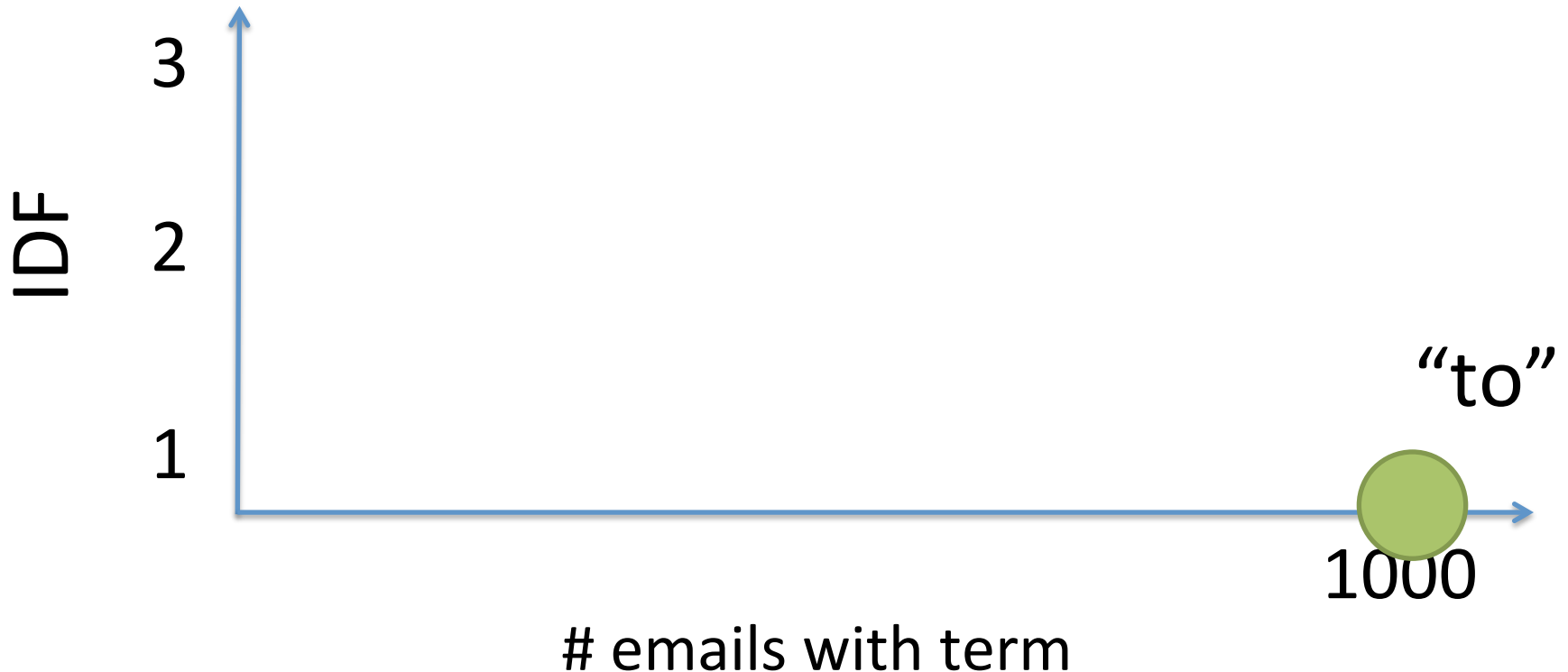
1000 Total Emails



# IDF. How do they work?

$\log(\text{total \# emails} / \text{\# emails with word})$

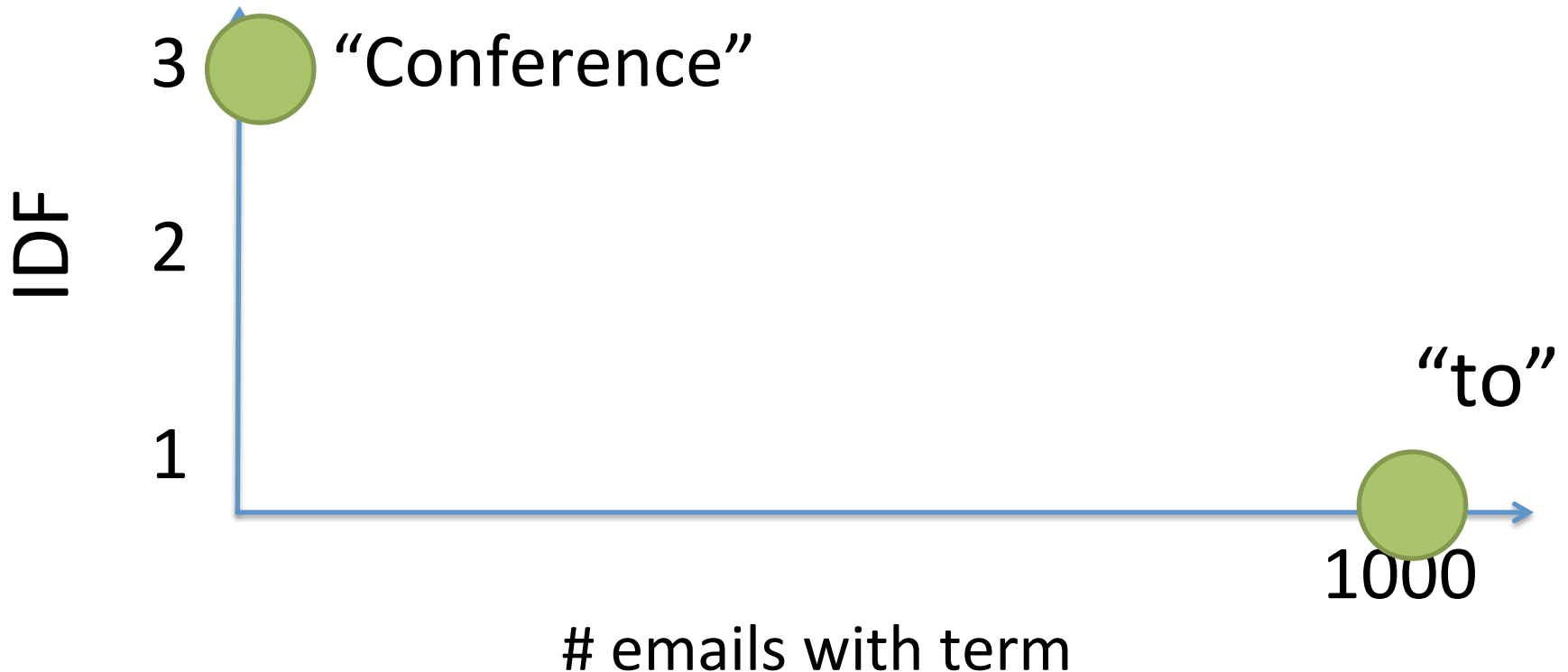
1000 Total Emails



# IDF. How do they work?

$\log(\text{total \# emails} / \text{\# emails with word})$

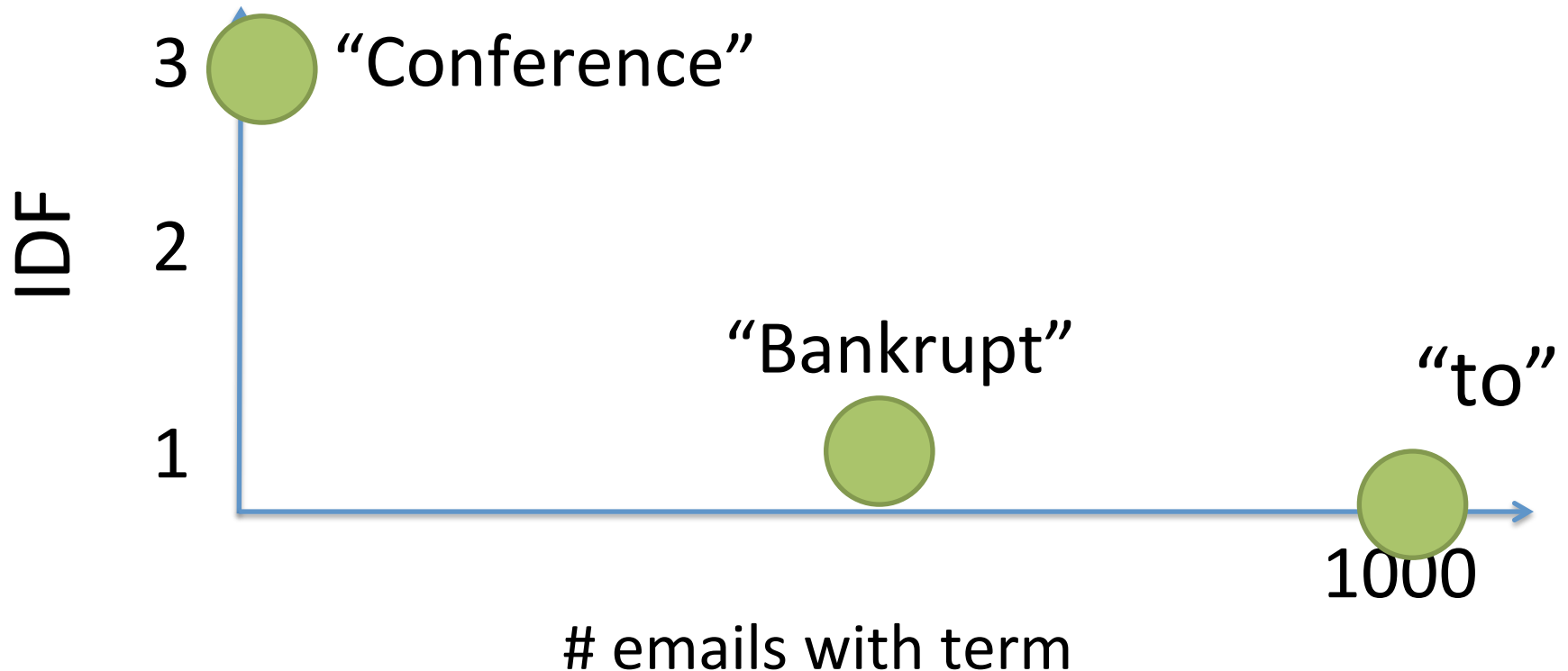
1000 Total Emails



# IDF. How do they work?

$\log(\text{total \# emails} / \text{\# emails with word})$

1000 Total Emails



# Relevant Terms

Frequent  
Words in email

But not in all  
emails

# Relevant Terms

Frequent  
Words in email

But not in all  
emails

TF

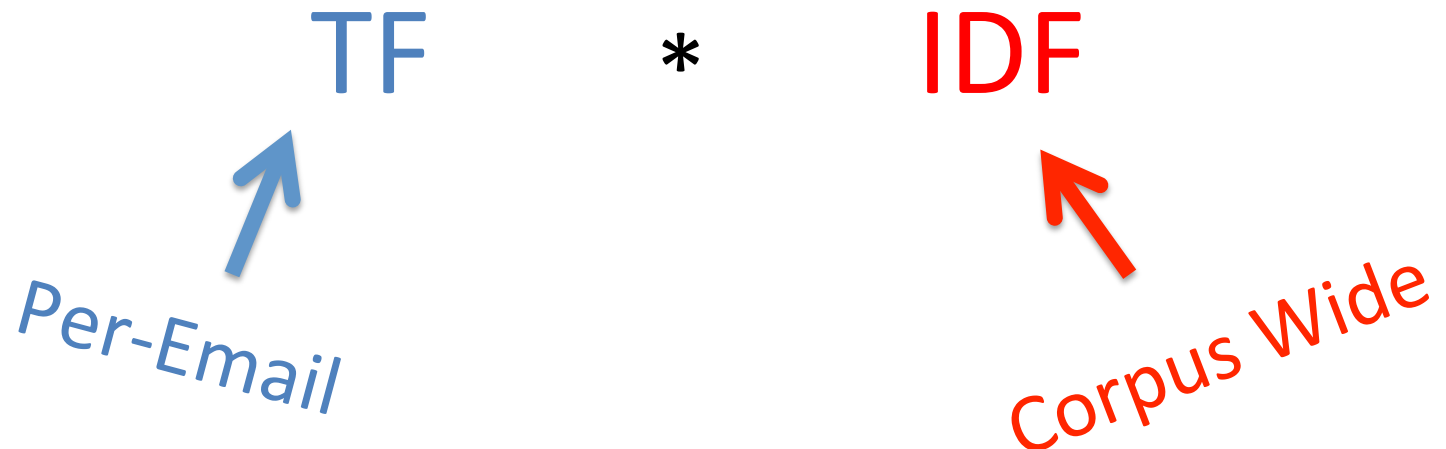
\*

IDF

# Relevant Terms

Frequent  
Words in email

But not in all  
emails





# Relevant Terms

We are in the process of trying to arrange a **conference** call with you on either Tuesday or Wednesday of next week to discuss the paper which is attached.

We will be doing this by **conference** call and once we set a time to talk with you, will give you the number to call.

# Relevant Terms

We are in the process of trying to arrange a **conference** call with you on either Tuesday or Wednesday of next week to discuss the paper which is attached.

We will be doing this by **conference** call and once we set a time to talk with you, will give you the number to call.



	TFIDF
conference	4.1
call	3
...	
to	0.03
a	0.0012

How similar is email 1 to email 2?

Cosine Similarity

Email1 = conference, enron, donuts,...

Email2 = enron, call, appointment,...

Email1 = conference, **enron**, donuts,...

Email2 = **enron**, call, appointment,...

Email1  $\cap$  Email2

Email1 = conference, **enron**, donuts,...

Email2 = **enron**, call, appointment,...

Email1  $\cap$  Email2

---

# words in both emails

Email1 • Email2

---

||Email1|| \* ||Email2||

E1 = { 'conference' : 4, 'enron' : 3 }  
E2 = { 'enron' : 1, 'call' : 2 }



E1 = { 'conference':4, 'enron':3 }

E2 = { 'enron': 1, 'call':2 }

E1[ 'enron' ] \* E2[ 'enron' ] + ...

E1 = { 'conference':4, 'enron':3 }

E2 = { 'enron': 1, 'call':2 }

$$\frac{E1['enron'] * E2['enron'] + \dots}{\underbrace{\text{sqrt}(4^2+3^2)}_{E1} * \underbrace{\text{sqrt}(1^2+2^2)}_{E2}}$$

# Data Cleaning

ok. 10:40 to be safe:P

On Fri, Jan 6, 2012 at 5:15 PM, Eugene Wu <sirrice@gmail.com> wrote:

> i have a feeling that breakfast foods stop at 11 at clover because thats

> how it works at the truck.

> so maybe 1045 or something is better.

>

>

> On Fri, Jan 6, 2012 at 5:09 PM, Adam Marcus <marcua@csail.mit.edu> wrote:

>

>> see you there at 11!

>>

>>

>> On Fri, Jan 6, 2012 at 5:05 PM, Eugene Wu <sirrice@gmail.com> wrote:

>>

>>> lets have brunch at 11. That way we skip the rush as well.

>>>

```
<html><body><table width="100%" cellpadding="0" cellspacing="0" border="0"><tr><td bgcolor="#f3f7fe"><table cellpadding="10" bgcolor="f3f7fe" border="0" cellspacing="0"><tr><td align="left"><p style="margin-left: 20px;"><font face="arial" size="1" color="#666666"><a href="http://Link.p0.com/t.d?T4Go9L_LLgypcb=@HTML_2PREVIEW_2LINK_0a=5CRicuFlRqyK5Ptqj2VYenr&msgVersion=mobile"><font color="#0038a5">View mobile version</font></a></font></p><table border="0" cellspacing="0" cellpadding="0" width="580"><tr><td width="1" height="
```

- Regular Expressions

“Enron’s”

~~“sirrice@gmail.com”~~

- Normalization

“Word” == “word”

- Stop Words

Ignore “a”, “to”, ...

So we implemented all of that...

Sucked



# Email Body vs Subject

Garbage in garbage out

- Forwarded message
- Unsubscribe
- Love
- Sincerely

# Subject Line

- Fwd: 11/5 [Dinner](#) @ Jewel Bako. 7:45PM
- [Dinner](#)
- What time for [dinner](#) tomorrow?

# Pretty Long

(optional)



# Download your IMAP Email!

[dataiap/resources/download\\_emails.py](https://dataiap/resources/download_emails.py)

# Presentations Wednesday!

- 1 minute
- 2 slides by Tuesday night
- Google presentations

Interesting dataset

Interesting analysis

Topic Extensions

<http://dataiap.github.com/dataiap/day4>

git pull

- Remind them about
  - Downloading their own email
  - Last day presentations.
- Proton pump inhibitor