

Data Cleaning

February 12, 2019

Data Science CSCI 1951A

Brown University

Instructor: Ellie Pavlick

HTAs: Wennie Zhang, Maulik Dang, Gurnaaz Kaur

A message from your Health and Wellness Advocates!

- Health and Wellness Resource Fair - Feb 20 (Wed), 4-6PM, CIT 1st Floor Atrium
- Open Hours! wellness.advocates@lists.brown.edu

Announcements

- A message from your Health and Wellness Advocates!
- iClicker syncing—the saga continues
- Collab policy... ..
- Final project teaming—fill out the form please!

Announcements

COLLABORATION POLICY!

Announcements

COLLABORATION POLICY!

SERIOUSLY, PEOPLE.....

Today

- Problems with dirty data
- Cleaning and string matching heuristics
- Bash commands—greatest hits (for data scientists)

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Problems?

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Problems?

Inconsistent
Representations

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Problems?

Inconsistent
Representations

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providence	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Missing Values

Problems?

Inconsistent
Representations

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Typos

Missing Values

Duplicates

Problems?

Inconsistent
Representations

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Typos

Missing Values

Duplicates

Problems?

Inconsistent
Representations

ID	Name	Street	City	State	Zip	Hours
1	J Meltzer	123 University Ave	Providence	RI	98106	42
2	Erin Bugbee	245 3rd St	Pawtucket	RI	98052-1234	30
3	David Wang	345 Broadway	PVD	Rhode Island	98101	19
4	E Bugbe	245 Third Street	Pawtucket	NULL	98052	299
5	Dave Wang	345 Broadway St	Providnce	Rhode Island	98101	19
6	Jacob Meltzer	123 Univ Ave	PVD	Rhode Island	NULL	41
7	Haomo Ni	123 University Ave	Providence	Guyana	94305	NULL

...

Maybe Duplicates?

Typos

Missing Values

Dirty Data...

Dirty Data...

- Data is dirty on its own

Dirty Data...

- Data is dirty on its own
- Data sets are clean on their own but combining them introduces errors (e.g. duplicates, different naming conventions)

Dirty Data...

- Data is dirty on its own
- Data sets are clean on their own but combining them introduces errors (e.g. duplicates, different naming conventions)
- Data doesn't "age well" (inflation, restricting)

Dirty Data...

- Data is dirty on its own
- Data sets are clean on their own but combining them introduces errors (e.g. duplicates, different naming conventions)
- Data doesn't "age well" (inflation, restricting)
- Any combination of the above

Dirty Data...

Dirty Data...

- Parsing input data (e.g., separator issues)

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates
- Missing values and required fields (e.g., always use 0)

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates
- Missing values and required fields (e.g., always use 0)
- Different representations (2 vs Two)

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates
- Missing values and required fields (e.g., always use 0)
- Different representations (2 vs Two)
- Fields too long (get truncated)

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates
- Missing values and required fields (e.g., always use 0)
- Different representations (2 vs Two)
- Fields too long (get truncated)
- Primary key violations (from data merging)

Dirty Data...

- Parsing input data (e.g., separator issues)
- Naming conventions: NYC vs New York
- Formatting issues – esp. dates
- Missing values and required fields (e.g., always use 0)
- Different representations (2 vs Two)
- Fields too long (get truncated)
- Primary key violations (from data merging)
- Redundant Records (from data merging)

Clicker Lightning Round!





TAS

ID	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42
7	Haomo Ni	Warwick	RI	NULL

ID	Name	City	State	Hours
1	J Meltzer	Providence	Rhode Island	42
2	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38
7	Haomo Ni	Warwick	Rhode Island	0

Clicker Lightning Round!





TAS







 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI	 NULL

ID	Name	City	State	Hours
1	J Meltzer	Providence	Rhode Island	42
2	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38
7	Haomo Ni	Warwick	Rhode Island	0

Clicker Lightning Round!





TAS







 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI	 NULL

ID 	Name	City	State	Hours
 1 	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island	0 

Clicker Lightning Round!

TAS

 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI 	NULL

ID 	Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island 	0 

How will the dirty data affect the results of this query?





(a) Too high








(b) Too low

(c) Unaffected

Clicker Lightning Round!

TAS

 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI 	NULL

ID 	Name	City	State	Hours
 1 	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island 	0 

How will the dirty data affect the results of this query?

(a) Too high

(b) Too low





(c) Unaffected

How many TAs are there?

```
SELECT COUNT ( * )
FROM TAS
```








Clicker Lightning Round!

TAS

 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI 	NULL

How will the dirty data affect the results of this query?

- (a) Too high
- (b) Too low
- (c) Unaffected

ID 	Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island 	0 





How many TAs are there?






```
SELECT COUNT ( * )
FROM TAS
```

Duplicates ->
Double Counting

Clicker Lightning Round!

TAS

 ID	 Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI	 NULL

ID	 Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island	0 

How will the dirty data affect the results of this query?

(a) Too high

(b) Too low





(c) Unaffected

How many TAs have worked zero hours?

```
SELECT COUNT (*)  
FROM TAS  
WHERE Hours = 0
```

Clicker Lightning Round!

TAS






 ID	 Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI	 NULL

How will the dirty data affect the results of this query?

(a) Too high

(b) Too low

(c) Unaffected

ID	 Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island	0 

How many TAs have worked zero hours?





```
SELECT COUNT (*)
FROM TAS
WHERE Hours = 0
```

35

NULLS aren't included
in the where clause

Clicker Lightning Round!

TAS






 ID	 Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI	 NULL

How will the dirty data affect the results of this query?

(a) Too high

(b) Too low

(c) Unaffected





ID	 Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island	0 

How many hours do my commuter TAs work?

```
SELECT SUM(Hours)
FROM TAS
WHERE City != "Providence"
```







Clicker Lightning Round!

TAS

 ID 	Name	City	State	Hours
1	J Meltzer	Providence	RI	42
2	Erin Bugbee	Pawtucket	RI	30
3	David Wang	PVD	Rhode Island	19
4	E Bugbe	Pawtucket	NULL	300
5	Dave Wang	Providence	Rhode Island	19
6	Jacob Meltzer	PVD	Rhode Island	42 
7	Haomo Ni	Warwick	RI 	NULL

How will the dirty data affect the results of this query?

- (a) Too high
- (b) Too low
- (c) Unaffected

ID 	Name	City	State	Hours
 1	J Meltzer	Providence	Rhode Island	42
2 	Erin Bugbee	Pawtucket	Rhode Island	30
3	David Wang	Providence	Rhode Island	38 
7	Haomo Ni	Warwick	Rhode Island 	0 

Inconsistent
names, typos,
and
duplicates...

How many hours do my commuter TAs work?

```
SELECT SUM(Hours)
FROM TAS
WHERE City != "Providence"
```

What's to be done?

What's to be done?

- Look at your data!

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things
- Look at your data

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things
- Look at your data
- When you issue a query, don't take the answer as gospel.

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things
- Look at your data
- When you issue a query, don't take the answer as gospel. Instead...

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things
- Look at your data
- When you issue a query, don't take the answer as gospel. Instead...wait for it...

What's to be done?

- Look at your data!
- Maybe set (sensible) defaults
- Maybe remove outliers
- Look at your data
- Maybe machine learn some of the things
- Look at your data
- When you issue a query, don't take the answer as gospel. Instead...wait for it...look at your data!

Look at your data

Look at your data

```
SELECT City, COUNT(*) as pop  
FROM PEOPLE  
GROUP BY Zip_Code  
ORDER BY pop
```

Look at your data

```
SELECT City, COUNT(*) as pop
FROM PEOPLE
GROUP BY Zip_Code
ORDER BY pop
```

City	Count(*)
Schenectady	2,500
New York City	2,200
Los Angeles	1,900
Dallas	1,400

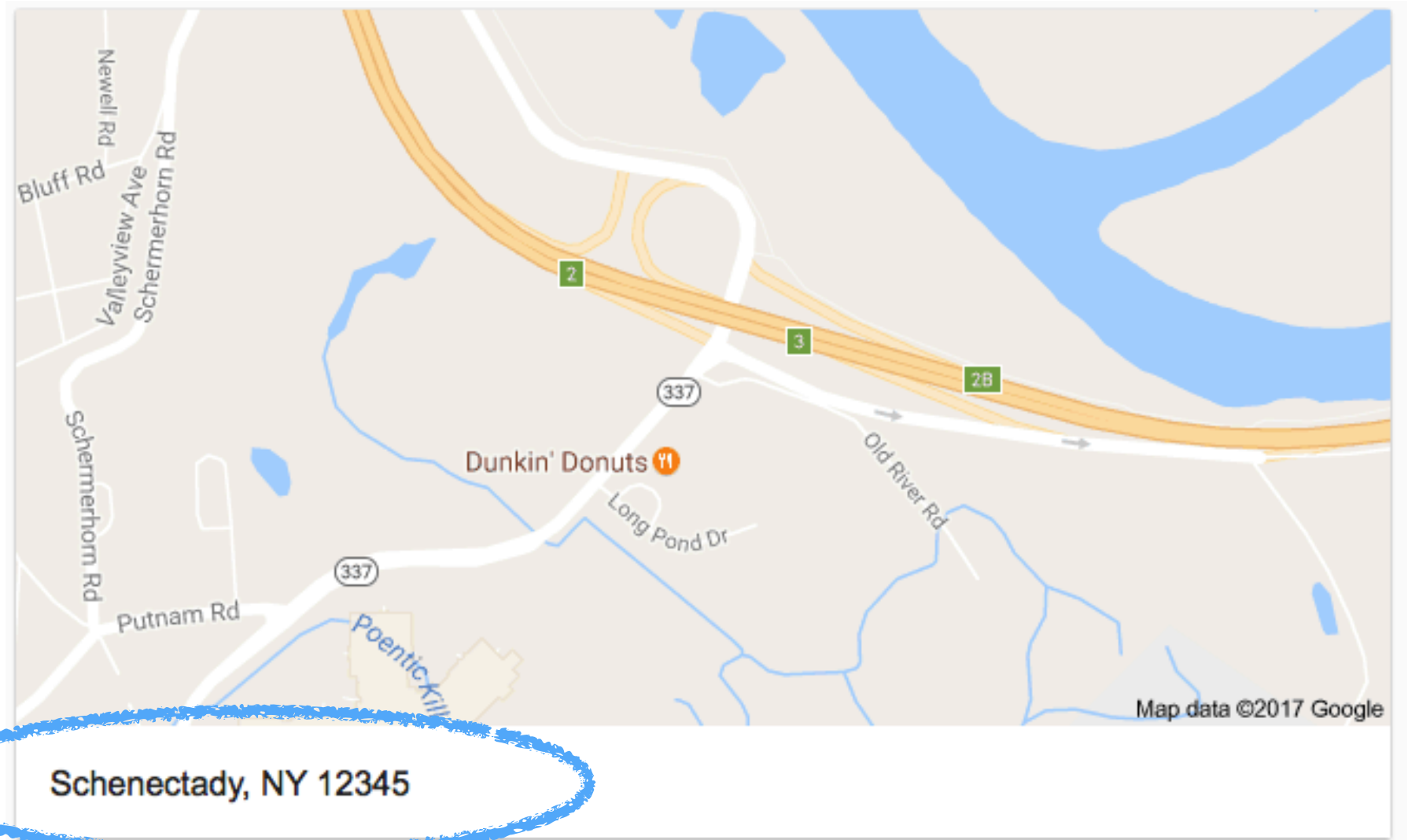
Look at your data

```
SELECT City, COUNT(*) as pop
FROM PEOPLE
GROUP BY Zip_Code
ORDER BY pop
```

City	Count(*)
Schenectady	2,500
New York City	2,200
Los Angeles	1,900
Dallas	1,400

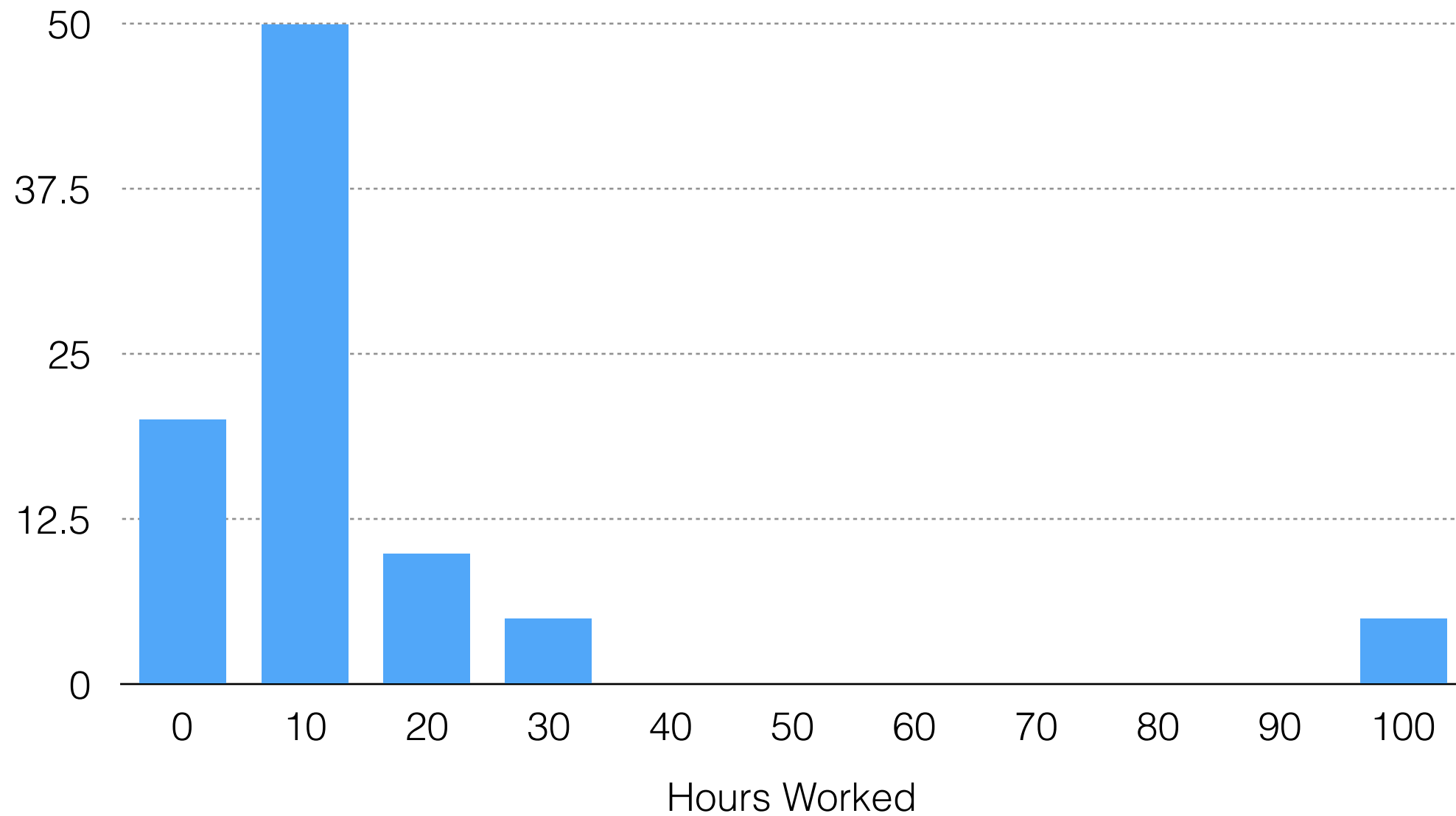
LOC

```
SELECT (
FROM PE
GROUP B
ORDER B
```

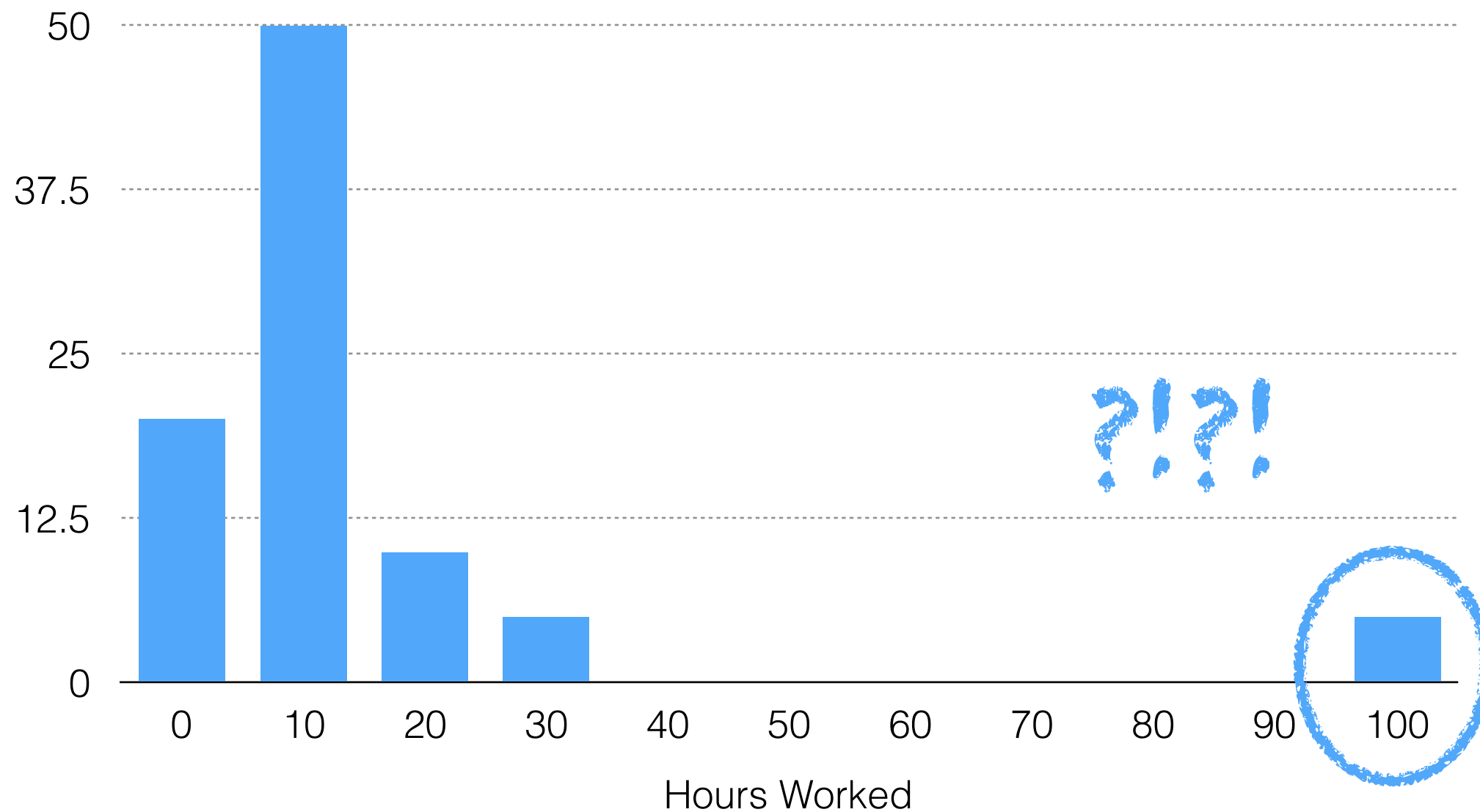


City	Count(*)
12345	2,500
10001	2,2000
90001	1,900
75001	1,400

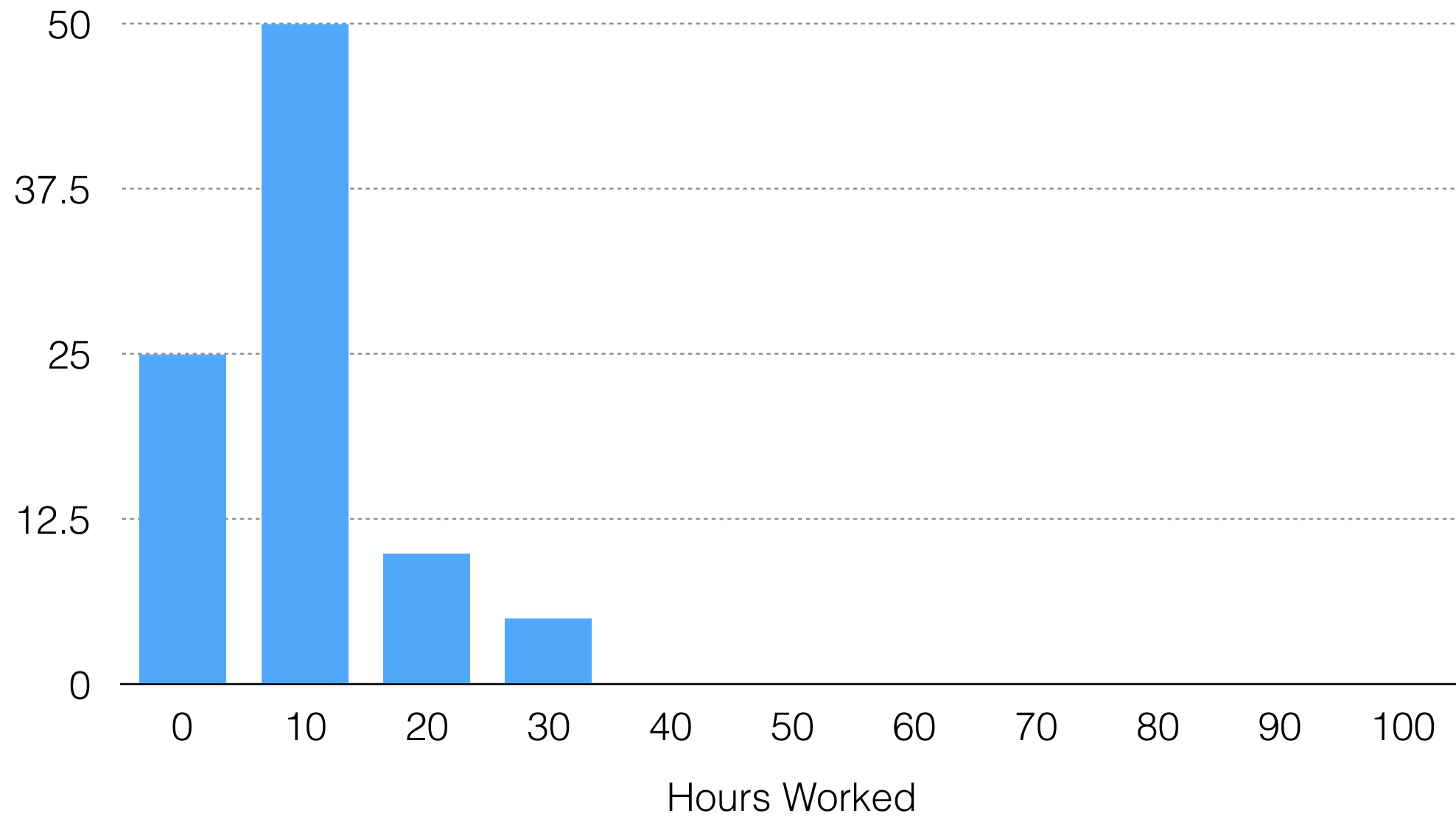
Set Defaults/Remove Outliers



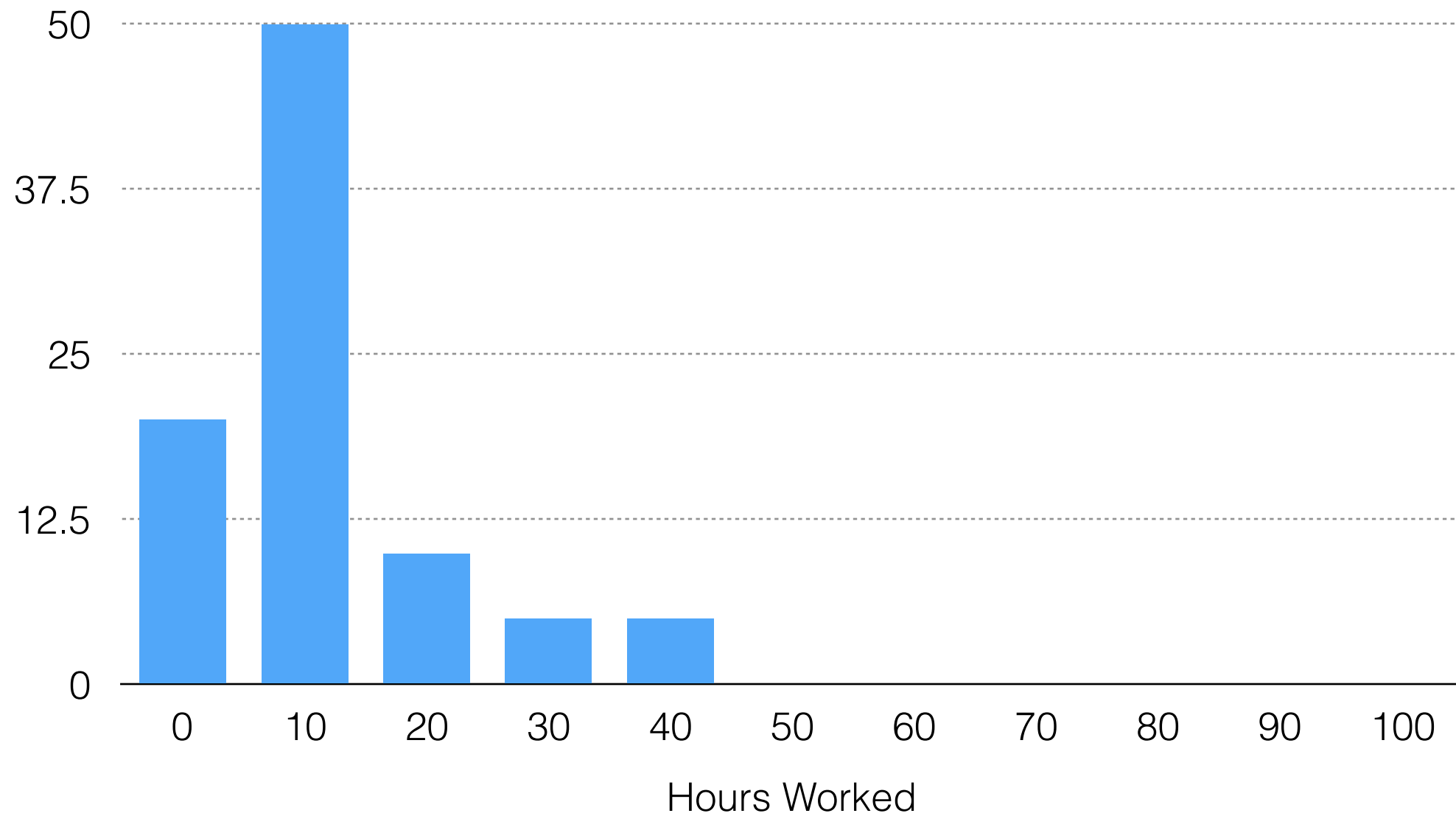
Set Defaults/Remove Outliers



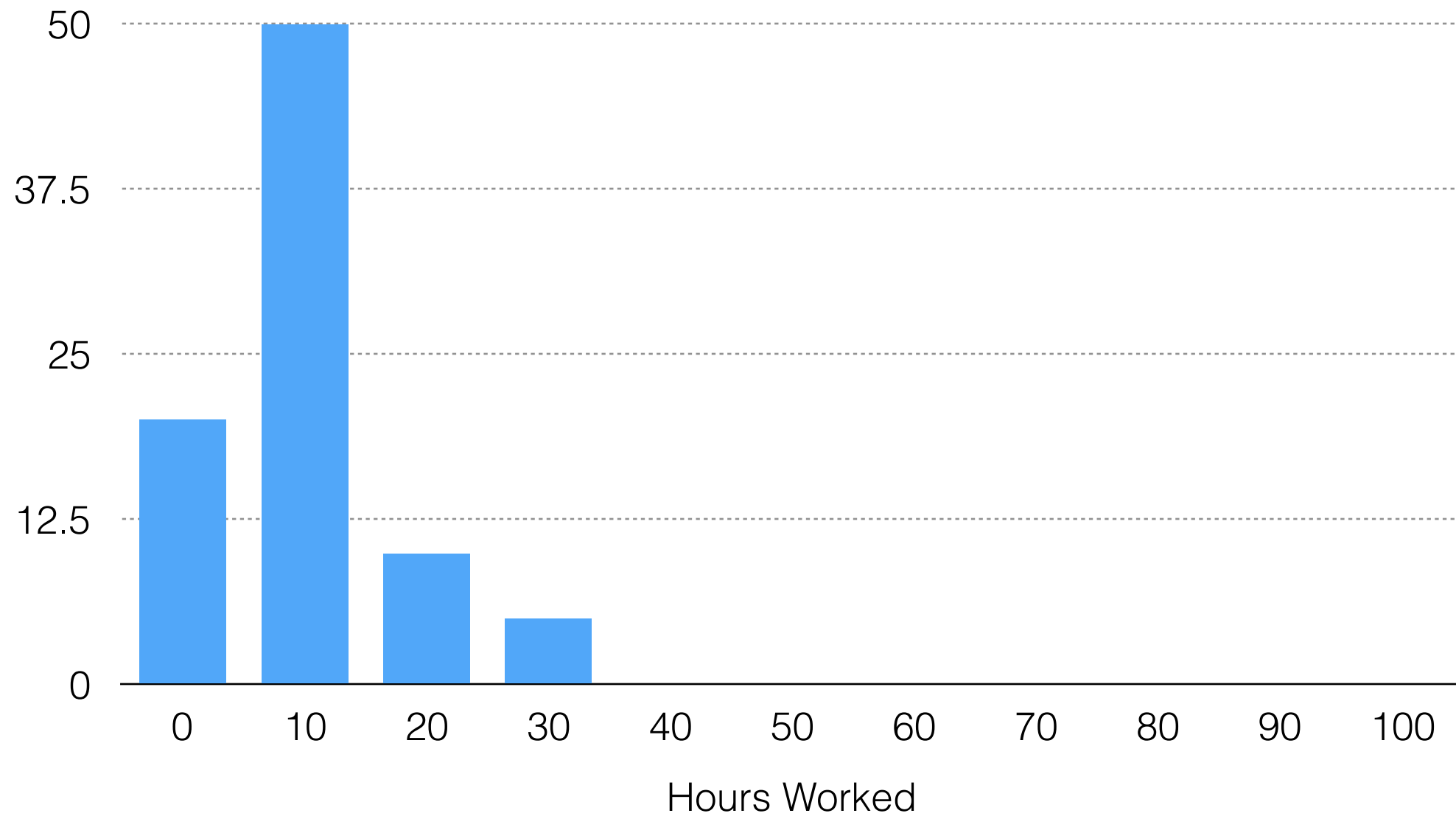
Set Defaults/Remove Outliers



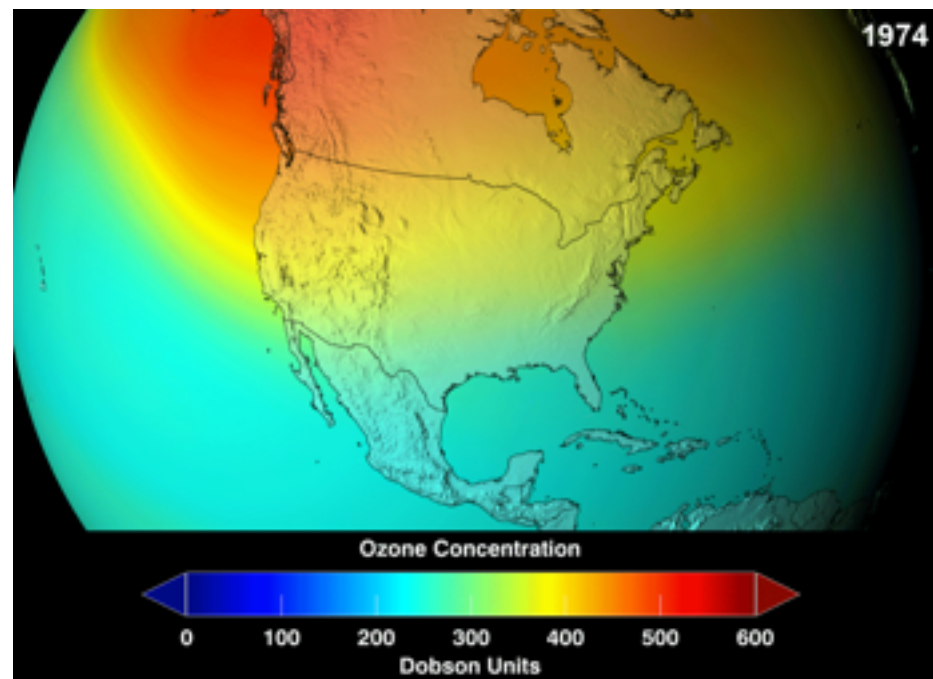
Set Defaults/Remove Outliers



Set Defaults/Remove Outliers



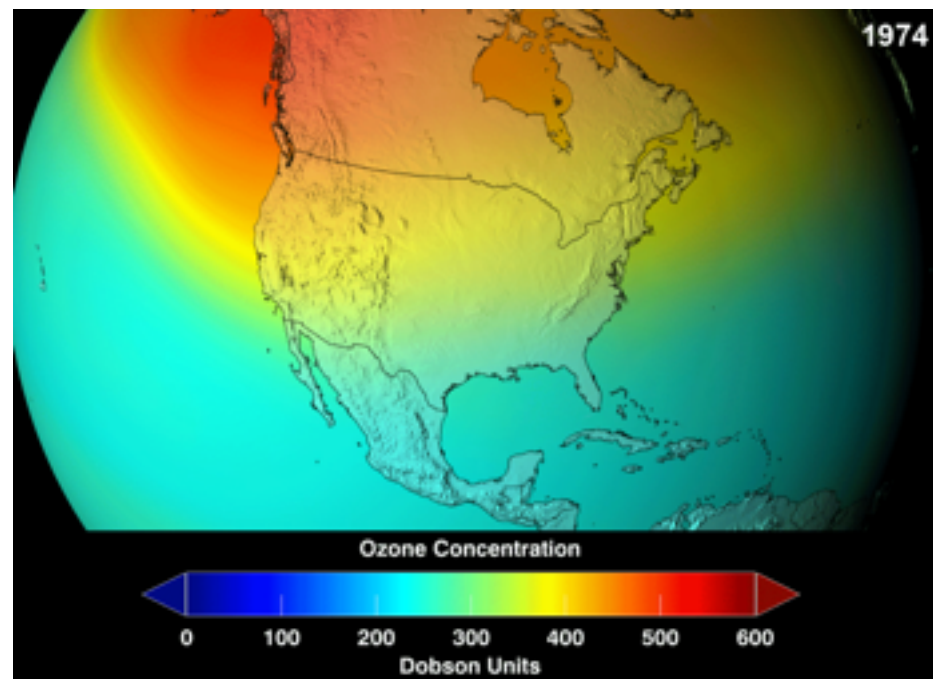
Set Defaults/Remove Outliers



The discovery of the Antarctic "ozone hole" by British Antarctic Survey scientists Farman, Gardiner and Shanklin...came as a shock to the scientific community...[The data] were initially rejected as unreasonable by data quality control algorithms (they were filtered out as errors since the values were unexpectedly low); the ozone hole was detected only in satellite data when the raw data was reprocessed following evidence of ozone depletion in in situ observations. When the software was rerun without the flags, the ozone hole was seen as far back as 1976.

https://en.wikipedia.org/wiki/Ozone_depletion#Antarctic_ozone_hole

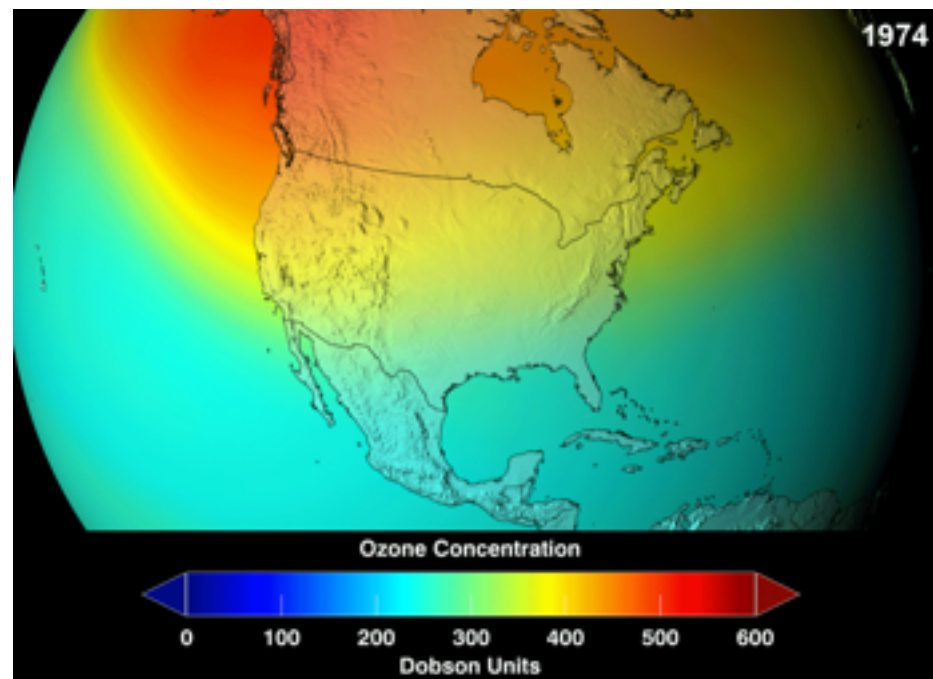
Set Defaults/Remove Outliers



The discovery of the Antarctic "ozone hole" by British Antarctic Survey scientists Farman, Gardiner and Shanklin...came as a shock to the scientific community...[The data] were initially rejected as unreasonable by data quality control algorithms (they were filtered out as errors since the values were unexpectedly low); the ozone hole was detected only in satellite data when the raw data was reprocessed following evidence of ozone depletion in in situ observations. When the software was rerun without the flags, the ozone hole was seen as far back as 1976.

https://en.wikipedia.org/wiki/Ozone_depletion#Antarctic_ozone_hole

Set Defaults/Remove Outliers



The discovery of the Antarctic "ozone hole" by British Antarctic Survey scientists Farman, Gardiner and Shanklin...came as a shock to the scientific community...[The data] were initially rejected as unreasonable by data quality control algorithms (they were filtered out as errors since the values were unexpectedly low); the ozone hole was detected only in satellite data when the raw data was reprocessed following evidence of ozone depletion in in situ observations. When the software was rerun without the flags, the ozone hole was seen as far back as 1976.

Always always
always! Look at
the data!

https://en.wikipedia.org/wiki/Ozone_depletion#Antarctic_ozone_hole

String Similarity

String Similarity: Edit Distance

Minimal number of edits (inserts, deletes, substitutions) needed to transform A into B.

https://en.wikipedia.org/wiki/Levenshtein_distance

String Similarity: Edit Distance

$$d_{i0} = \sum_{k=1}^i w_{\text{del}}(b_k), \quad \text{for } 1 \leq i \leq m$$

$$d_{0j} = \sum_{k=1}^j w_{\text{ins}}(a_k), \quad \text{for } 1 \leq j \leq n$$

$$d_{ij} = \begin{cases} d_{i-1,j-1} & \text{for } a_j = b_i \\ \min \begin{cases} d_{i-1,j} + w_{\text{del}}(b_i) \\ d_{i,j-1} + w_{\text{ins}}(a_j) \\ d_{i-1,j-1} + w_{\text{sub}}(a_j, b_i) \end{cases} & \text{for } a_j \neq b_i \end{cases} \quad \text{for } 1 \leq i \leq m, 1 \leq j \leq n.$$

https://en.wikipedia.org/wiki/Levenshtein_distance

String Similarity: Edit Distance

11⁵th Waterman St., Providence, RI

11⁰th Waterman St., Providence, RI

EditDistance = 1

String Similarity: Edit Distance

Waterman St~~reet~~, Providence, RI

Waterman St, Providence, RI

EditDistance = 4

String Similarity: Edit Distance

Problems?

String Similarity: Edit Distance

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

String Similarity: Edit Distance

Edit Distance = 0 

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

String Similarity: Edit Distance

Edit Distance = 0 

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

NE 148th Ave, Redmond, WA

String Similarity: Edit Distance

Edit Distance = 0 

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

148th Ave NE, Redmond, WA

NE 148th Ave, Redmond, WA

Edit Distance = 4

String Similarity: Jaccard Similarity

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|}$$

String Similarity: Jaccard Similarity

148th Ave NE, Redmond, WA

140th Ave NE, Redmond, WA

String Similarity: Jaccard Similarity

148th Ave NE, Redmond, WA

140th Ave NE, Redmond, WA

String Similarity: Jaccard Similarity

148th Ave NE, Redmond, WA

140th Ave NE, Redmond, WA

$$\text{Jaccard} = 4 / 6 = .67$$

String Similarity: Jaccard Similarity

148th Ave NE, Redmond, WA
NE 148th Ave, Redmond, WA

Jaccard = ???

https://en.wikipedia.org/wiki/Jaccard_index

String Similarity: Jaccard Similarity

148th Ave NE, Redmond, WA
NE 148th Ave, Redmond, WA

Jaccard = 1

https://en.wikipedia.org/wiki/Jaccard_index

Clicker Question!

iPad Two 16GB WiFi White

iPad 2nd generation 16GB WiFi White

What's the Jaccard Similarity?

(a) $3/8$

(b) $4/11$

(c) $4/7$

Clicker Question!

iPad Two 16GB WiFi White

iPad 2nd generation 16GB WiFi White

What's the Jaccard Similarity?

(a) $3/8$

(b) $4/11$

(c) $4/7$

$\#(\text{iPad, 16GB, Wifi, White})$

$\#(\text{iPad, Two, 2nd, generation, 16GB, Wifi, White})$

String Similarity: Jaccard Similarity

Michigan State University
Michigan State Univ.

Michigan State University
Ohio State University

https://en.wikipedia.org/wiki/Jaccard_index

String Similarity: Jaccard Similarity

Jaccard = 0.5

Michigan State University

Michigan State Univ.



Jaccard = 0.5

Michigan State University

Ohio State University



String Similarity: (Weighted) Jaccard Similarity

3

Jaccard = 0.5

Michigan¹ State¹ University

Michigan State Univ.

Jaccard = 0.25

Michigan State University

Ohio State University

String Similarity: (Weighted) Jaccard Similarity

3

Jaccard = 0.5

Michigan¹ State¹ University

Michigan State Univ.

Jaccard = 0.5

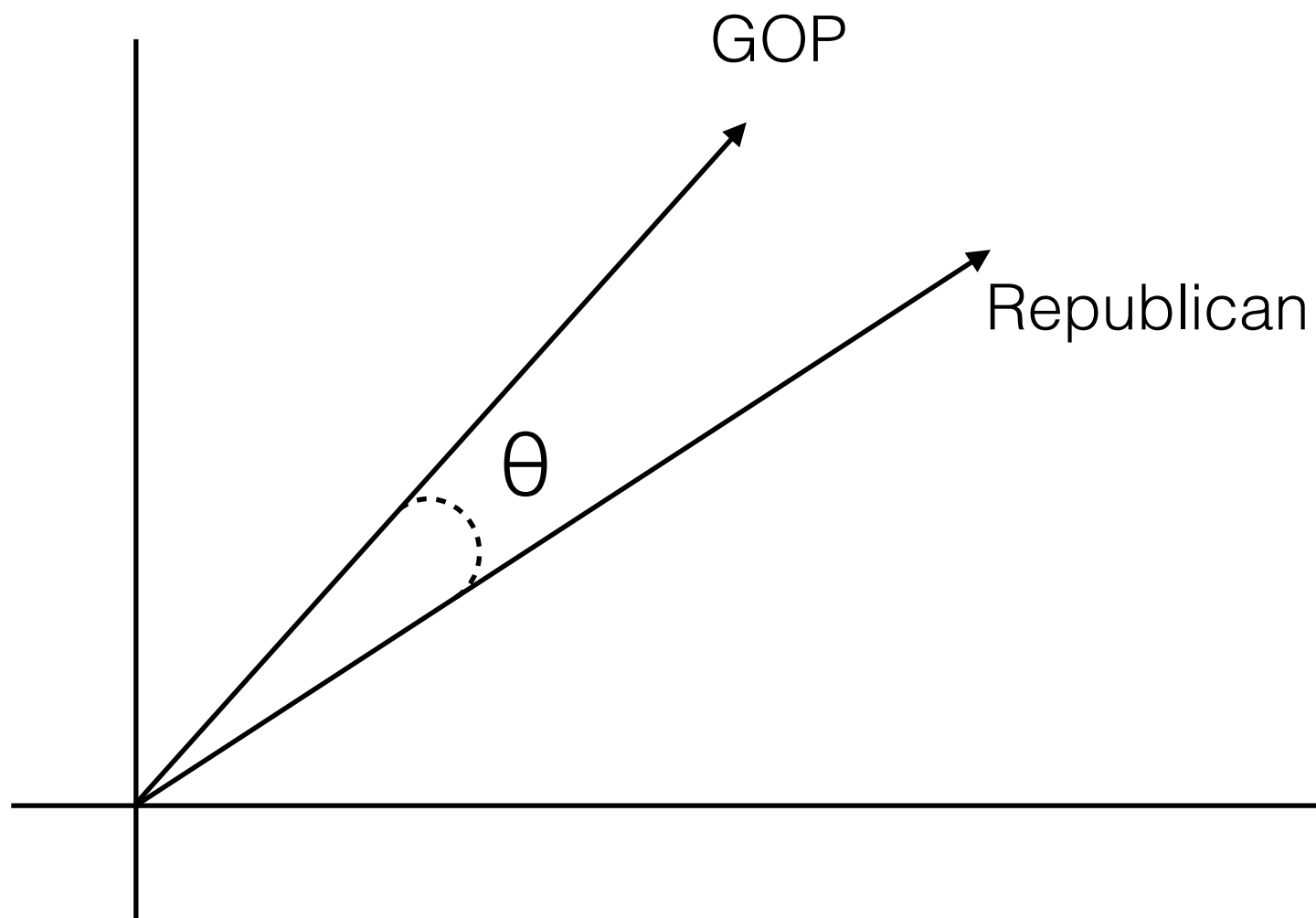
Michigan State University

University of Michigan

String Similarity: Cosine Similarity

	Senator	Washington	announced	party	primary	chairman
GOP	1002	41	502	700	400	3
Republican	800	35	521	698	423	10

String Similarity: Cosine Similarity



Clicker Question!

Brown
Brown Uni.

**Which metric would (likely)
consider the above words more
similar?**

- (a) Jaccard**
- (b) Cosine**

Clicker Question!

Brown
Brown Uni.

**Which metric would (likely)
consider the above words more
similar?**

- (a) Jaccard**
- (b) Cosine**

Clicker Question!

Motown
Detroit

**Which metric would (likely)
consider the above words more
similar?**

- (a) Jaccard**
- (b) Cosine**

Clicker Question!

Motown
Detroit

**Which metric would (likely)
consider the above words more
similar?**

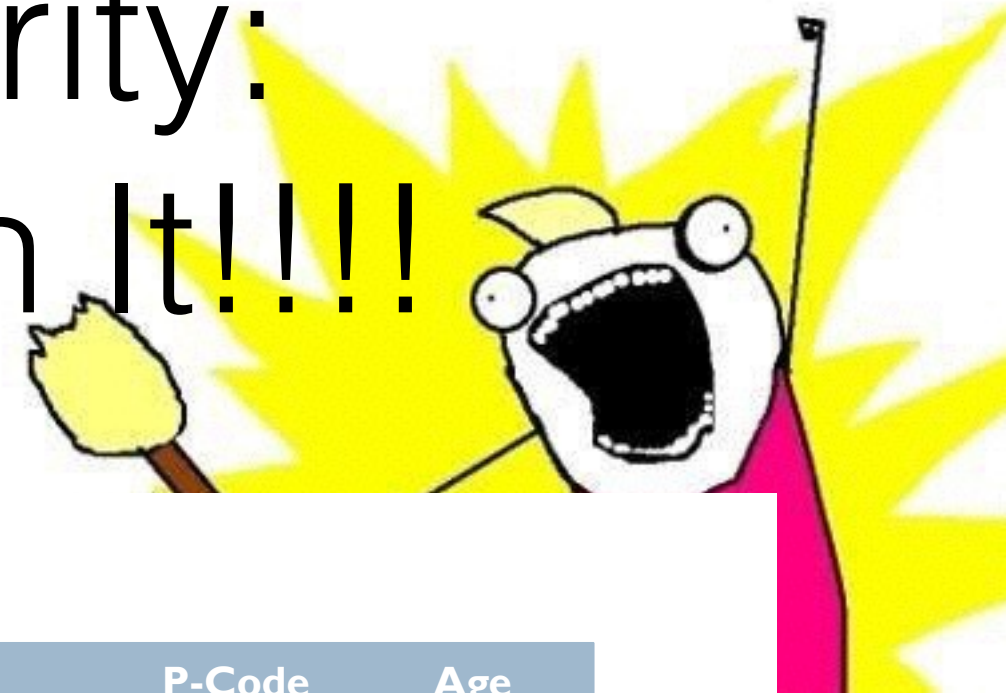
(a) Jaccard

 (b) Cosine

String Similarity: Machine Learn It!!!!



String Similarity: Machine Learn It!!!!



Customer

Id	Name	Street	City	State	P-Code	Age
1	J Smith	123 University Ave	Seattle	Washington	98106	42
2	Mary Jones	245 3rd St	Redmond	WA	98052-1234	30
3	Bob Wilson	345 Broadway	Seattle	Washington	98101	19
4	M Jones	245 Third Street	Redmond	NULL	98052	299
5	Robert Wilson	345 Broadway St	Seattle	WA	98101	19
6	James Smith	123 Univ Ave	Seatl	WA	NULL	41
7	J Widom	123 University Ave	Palo Alto	CA	94305	NULL
...

$WtJaccard = 0.57$

0.91

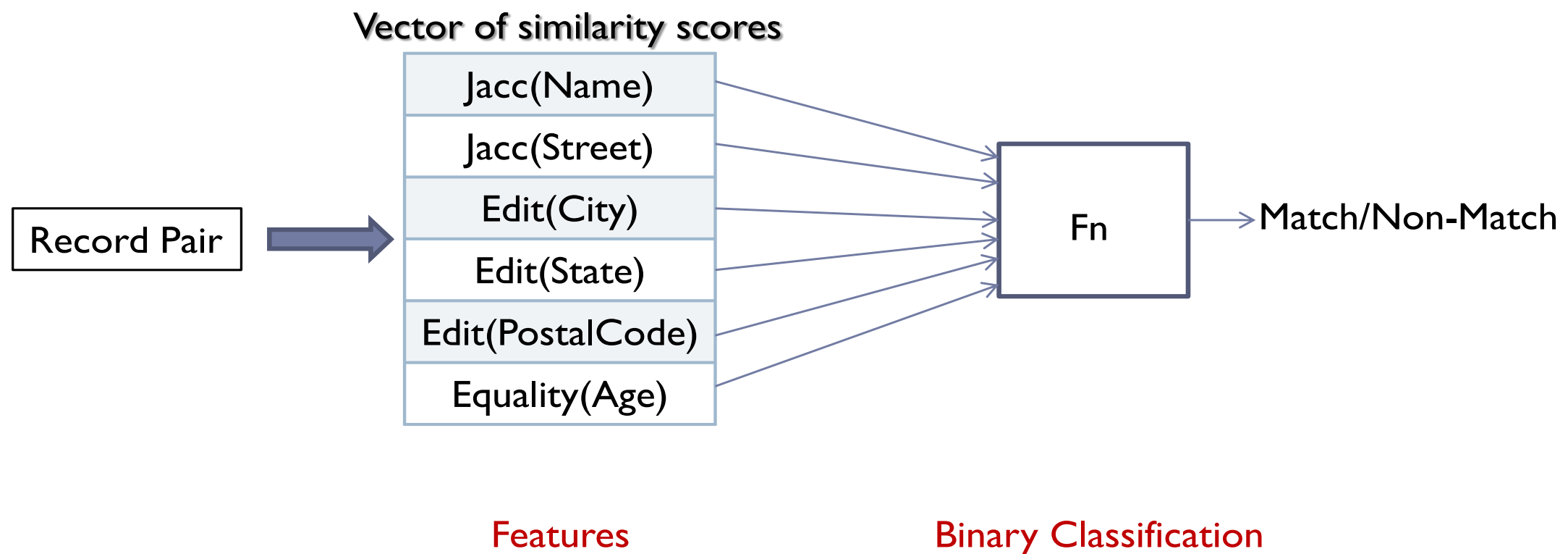
1.0

0.0

1.0

1.0

String Similarity: Machine Learn It!!!!



String Similarity: Machine Learn It!!!!



Bob Wilson	345 Broadway	Seattle	Washington	98101	19
Robert Wilson	345 Broadway St	Seattle	WA	98101	19

Match

B Wilson	123 Broadway	Boise	Idaho	83712	19
Robert Wilson	345 Broadway St	Seattle	WA	98101	19

Non-Match

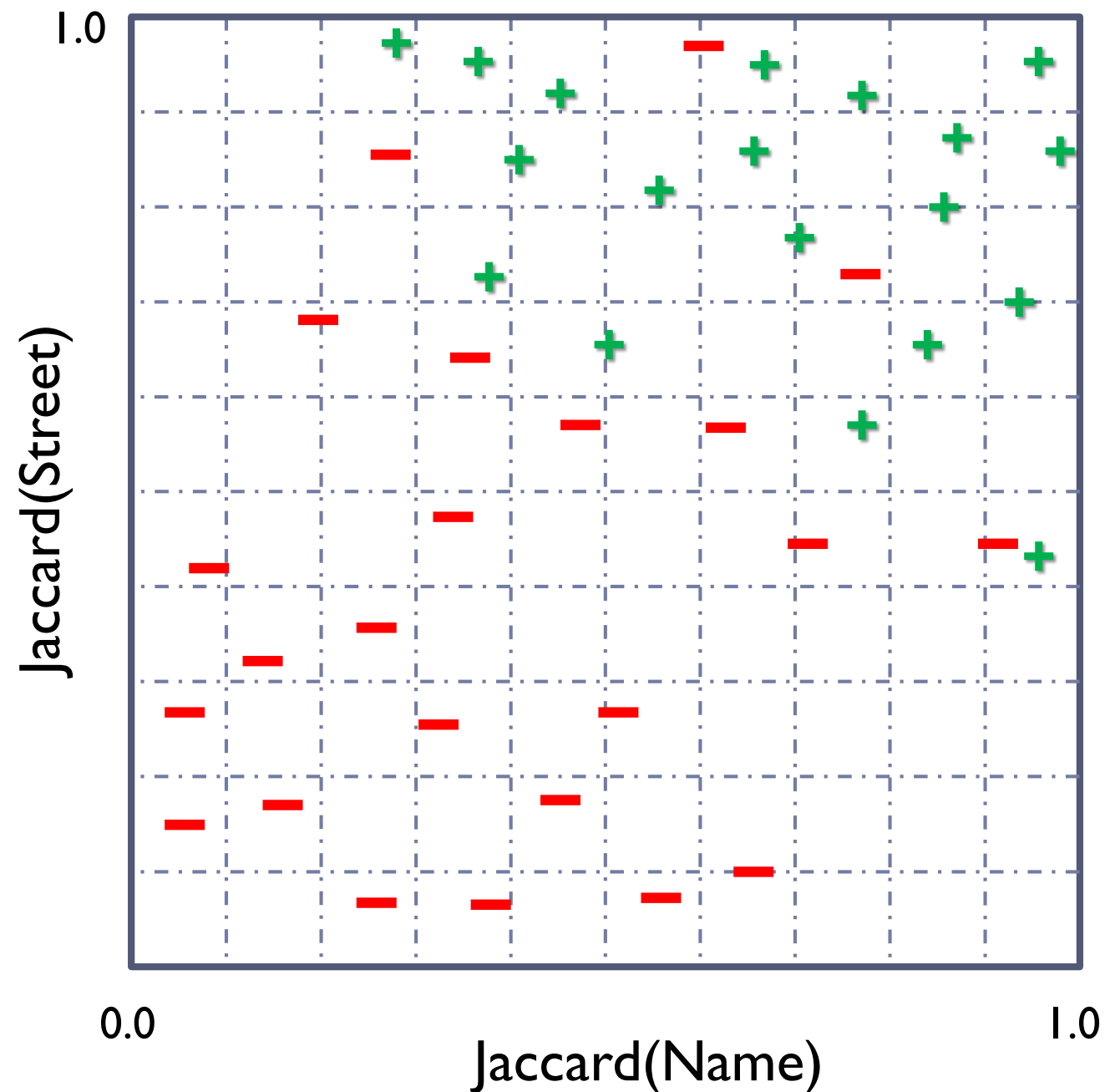
Mary Jones	245 3rd St	Redmond	WA	98052-1234	30
M Jones	245 Third Street	Redmond	NULL	98052	299

Match

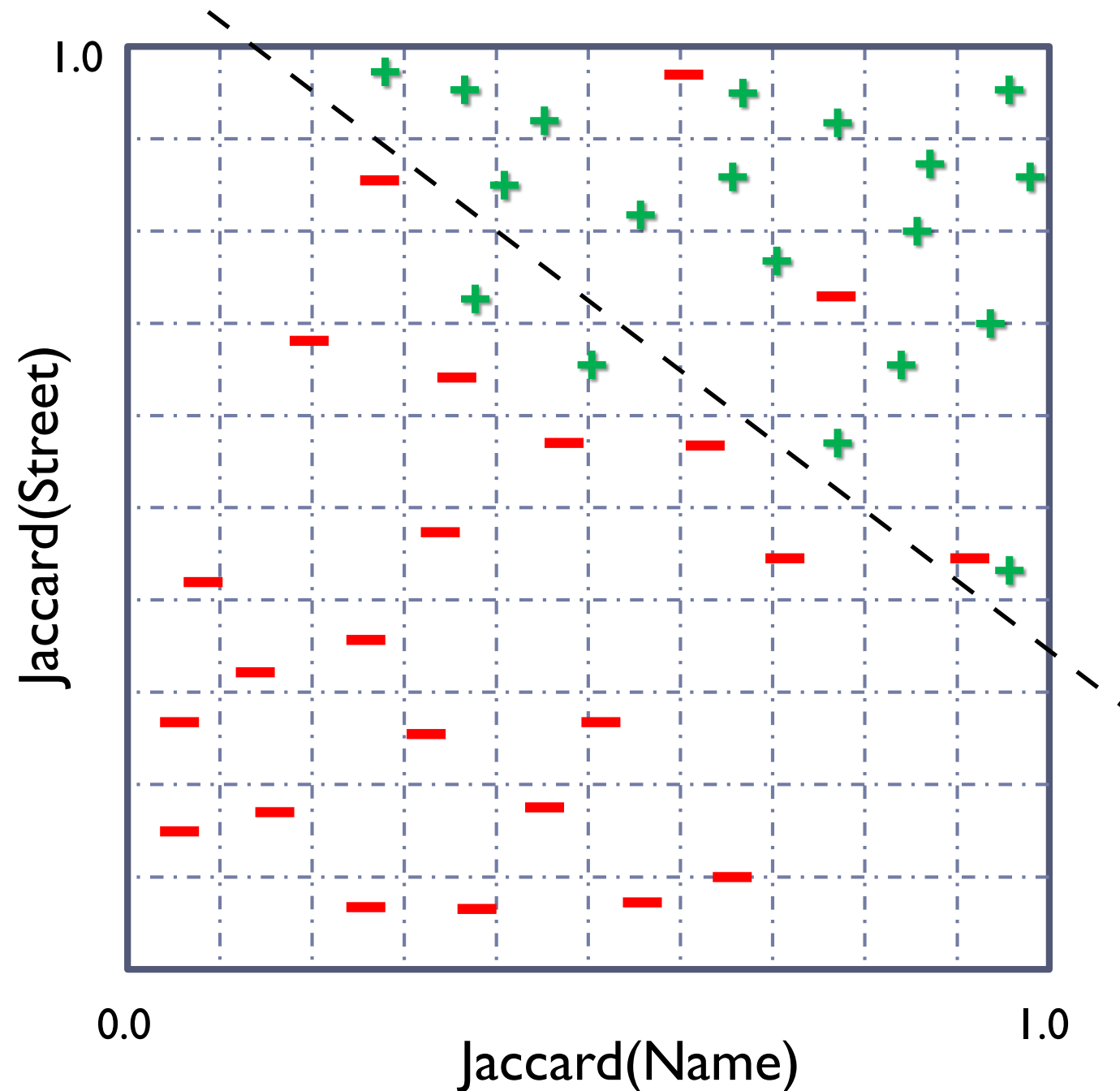
Mary Jones	245 3rd St	Redmond	WA	98052-1234	30
Robert Wilson	345 Broadway St	Seattle	WA	98101	19

Non-Match

String Similarity: Machine Learn It!!!!



String Similarity: Machine Learn It!!!!



Code-along!



```
cat data.txt | cut -f 2,4 | sort | uniq -c | sort -nr | head
```


Check in

How comfortable would you say you are using the command line?

- (a) Not at all...I am actually not even sure what the question means.
- (b) I've used it for turnin and a few other things (e.g. copy-pasting to install things), but thats it
- (c) I can get around, but its not my "home base"
- (d) Command line is my default for developing

Check in



How comfortable would you say you are using the command line?

- (a) Not at all...I am actually not even sure what the question means.
- (b) I've used it for turnin and a few other things (e.g. copy-pasting to install things), but thats it
- (c) I can get around, but its not my "home base"
- (d) Command line is my default for developing

Bash Scripting

<https://cs.brown.edu/people/epavlick/articles.txt>

- | | |
|----------------------|------------------|
| 1. ID | 6. Victim Age |
| 2. City | 7. Shooter Age |
| 3. State | 8. Url |
| 4. Date (YYYY-MM-DD) | 9. Title |
| 5. Time | 10. Article Text |

cat, less, head, tail

- what does this data even look like?

```
# first 10 lines of file  
$ head articles.txt
```

```
# first line of file  
$ head -n 1 articles.txt
```

```
# random 10 lines from file  
$ cat articles.txt | shuf | head
```

WC

- how many articles are there

how many bytes, words, and lines are there?

```
$ wc articles.txt
```

how many lines are there?

```
$ wc -l articles.txt
```

pipe (|), redirect (>)

```
$ head articles.txt | wc -l  
10
```

write output to file called "tmp"

```
$ head articles.txt > tmp
```

```
$ wc -l tmp  
10 tmp
```

```
$ head articles.txt | wc -l > tmp
```

```
$ cat tmp  
10
```

cut

```
$ cat articles.txt | cut -f 1 | head -n  
3
```

```
Antioch  
Greeley  
Bridgeport
```

```
$ cat articles.txt | cut -f 3 | cut -f 1  
-d '-' | head -n 3
```

```
2016  
2015  
2014
```

sort, uniq

```
# print the lowest 3 values (includes duplicates)
```

```
$ cat articles.txt | cut -f 4 | cut -f 1 -d '-' | sort | head -n 3  
1929  
1932  
1932
```

```
# print lowest three values (remove duplicates but count how many  
occurrences of each
```

```
$ cat articles.txt | cut -f 4 | cut -f 1 -d '-' | sort | uniq -c |  
head -n 3  
    1 1929  
    2 1932  
    3 1942
```

```
# find the most frequent years
```

```
$ cat articles.txt | cut -f 4 | cut -f 1 -d '-' | sort | uniq -c |  
sort -r | head  
5091 2015  
1821 2016  
1784 NA
```


sort, uniq

How many duplicated entries are there (using url as the uniq id)?

```
# total number of urls (lines)
```

```
$ cat articles.txt | cut -f 8 | wc -l  
9584
```

```
# number of unique urls
```

```
$ cat articles.txt | cut -f 8 | sort | uniq | wc -l  
7990
```

```
# number of duplicated urls
```

```
$ cat articles.txt | cut -f 8 | sort | uniq -d | wc -l  
981
```

regex (grep, sed, awk)

```
$ cat articles.txt | cut -f 2 | grep "NY" | head -n 5  
NY  
HOMINY  
NYC  
NY  
NY
```

```
$ cat articles.txt | cut -f 2 | grep "^NY$" | head  
NY  
NY  
NY  
NY
```

```
$ cat articles.txt | cut -f 2 | grep "^NY[.]*" | head  
NY  
NYC  
NY  
NY  
NY
```

regex (grep, sed, awk)

```
$ cat articles.txt | cut -f 4 | sed "s/[0-9]/#/g" | head -n 3
```

```
####-##-##
```

```
####-##-##
```

```
####-##-##
```

```
$ cat articles.txt | cut -f 3 | sed "s/[A-Z][A-Z] - //g" | grep -v Unclear |
```

```
head -n 3
```

```
Minnesota
```

```
North Carolina
```

```
Michigan
```

Being all fancy...

```
# replace all non-numeric characters with blanks
$ cat articles.txt | cut -f 6 | sed "s/[^0-9]//g" |
head
```

```
# plot a histogram of all ages
cat articles.txt | cut -f 6 | sed "s/[^0-9]//g" |
grep -v "^$" | python2 -c "import sys,
matplotlib.pyplot as plt; plt.hist([int(i) for i in
sys.stdin]); plt.show()"
```

```
# plot a histogram of all ages, removing outliers
cat articles.txt | cut -f 6 | sed "s/[^0-9]//g" |
grep -v "^$" | python2 -c "import sys,
matplotlib.pyplot as plt; plt.hist([min(int(i), 100)
for i in sys.stdin]); plt.show()"
```

\$

This is funny because it is a
regex joke. Please laugh and
validate me. I will wait.

