

Data Strategies

+ Data Tactics

Thomas Stoeger

thomas.stoeger@northwestern.edu

advertising (postdoc and grad student symposium): August 22



<https://www.nupostdocs.org/symposium.html>

aim: prepared to approach data

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

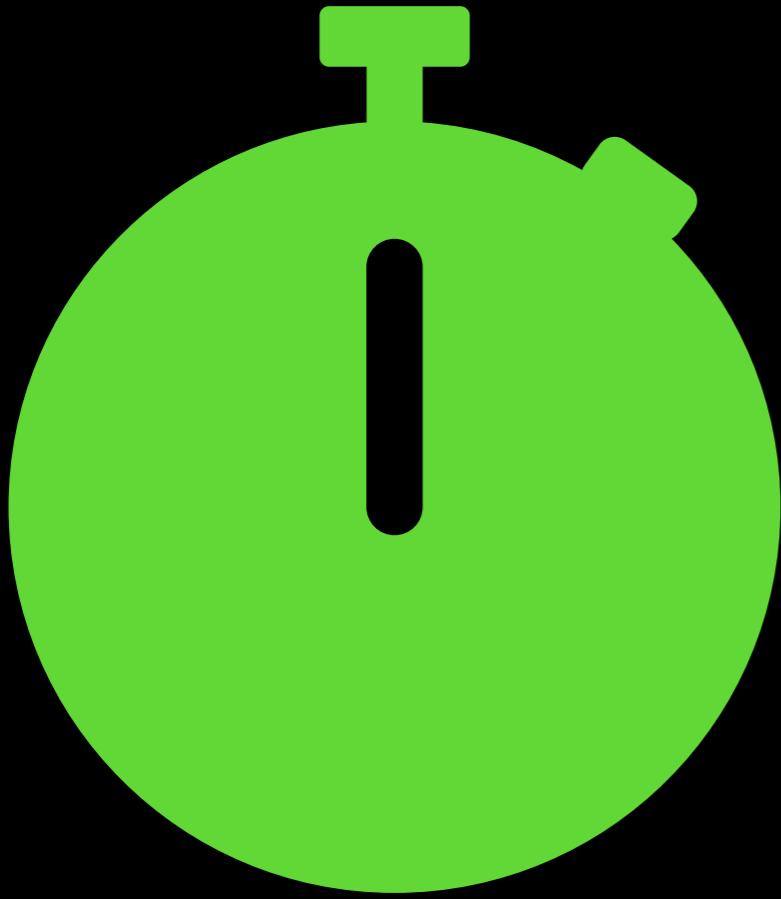
Setup

- Computer
- Organizing data
- Coding

break

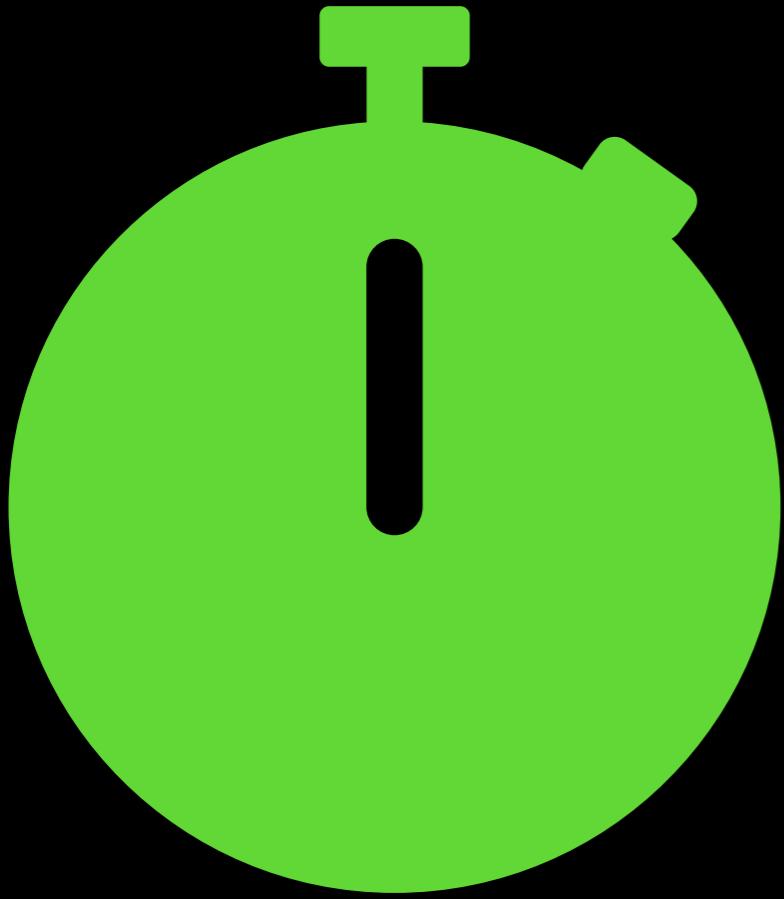


Two minute exercise.



Two minute exercise.

**Find a neighbor, or
neighbors, whom
you do not know.**



Two minute exercise.

Find a neighbor, or
neighbors, whom
you do not know.

**Then talk about the worst
data you have encountered.**

How can we anticipate problems?

How can we anticipate
problems?

“just”, “easy”, “simple”

How can we anticipate problems?

“just”, “easy”, “simple”

Can not share data or code.

Understand the past history of the data.

External:



Read paper
or documentation

Understand the past history of the data.

External:

Read paper
or documentation

**“is generated by
program”**

Read documentation

Internal:

Ask

Custom

Communication



Understand the past history of the data.

External:

↓
Read paper
or documentation

Internal:

↓
Ask
Custom
↓
Communication

“is generated by
program”



↓
Read documentation

**Write consistency checks!
(programmatically)**

Understand the past history of the data.

External:

↓
Read paper
or documentation

Internal:

↓
Ask
Custom
↓
Communication

“is generated by
program”



↓
Read documentation

**Write consistency checks!
(programmatically)**

**Expect to spend
days or weeks**

Understand the past history of the data.

External:

↓
Read paper
or documentation

“is generated by
program”

↓
Read documentation

Internal:

↓
Ask

↓
Custom

↓
Communication

**Write consistency checks!
(programmatically)**

**Expect to spend
days or weeks**

**Tips for getting
all information:**

- be decent person
- remove fear that you will scoop

Workflow

- Anticipate problems
- **Get overview**
- Clean data
- Find connections

Setup

- Environment
- Organizing data
- Coding

break

Introducing Kaggle ML and Data Science Survey, 2017

An example data set, which we will be using

The screenshot shows a web browser displaying the Kaggle website at the URL <https://www.kaggle.com/kaggle/kaggle-survey-2017>. The page title is "Kaggle ML and Data Science Survey, 2017". The page content includes a brief description: "A big picture view of the state of data science and machine learning." Below the description is a "Kaggle" logo and the text "last updated 9 months ago". On the right side of the page, there is a "641 voters" button. At the bottom of the page, there are navigation links: "Data" (which is underlined), "Overview", "Kernels", "Discussion", "Activity", "Download (4 MB)", and a blue "New Kernel" button.

Introducing Kaggle ML and Data Science Survey, 2017

An example data set, which we will be using

The screenshot shows a web browser displaying the Kaggle website at the URL <https://www.kaggle.com/kaggle/kaggle-survey-2017>. The page title is "Kaggle ML and Data Science Survey, 2017". The page content includes a brief description: "A big picture view of the state of data science and machine learning." Below the description is a "Kaggle" logo and the text "last updated 9 months ago". A large red heart-shaped graphic on the right contains the text "super nice example".

super nice example

Reviewed Dataset

Kaggle ML and Data Science Survey, 2017

A big picture view of the state of data science and machine learning.

Kaggle • last updated 9 months ago

Data Overview Kernels Discussion Activity

Download Kernel

The screenshot shows a file browser window with the following interface elements:

- Toolbar:** Includes icons for grid view, list view (selected), card view, search, filter, sort, and upload.
- Name:** A column header for the list of files.
- File List:** A list of five files with their corresponding icons:
 - conversionRates.csv
 - freeformResponses.csv
 - multipleChoiceResponses.csv
 - RespondentTypeREADME.txt
 - schema.csv

Name
conversionRates.csv
freeformResponses.csv
multipleChoiceResponses.csv
RespondentTypeREADME.txt
schema.csv

All: Every respondent was shown this question

Non-worker: Respondents who indicated that they were "Not employed, and not looking for work" or "I prefer not to say"

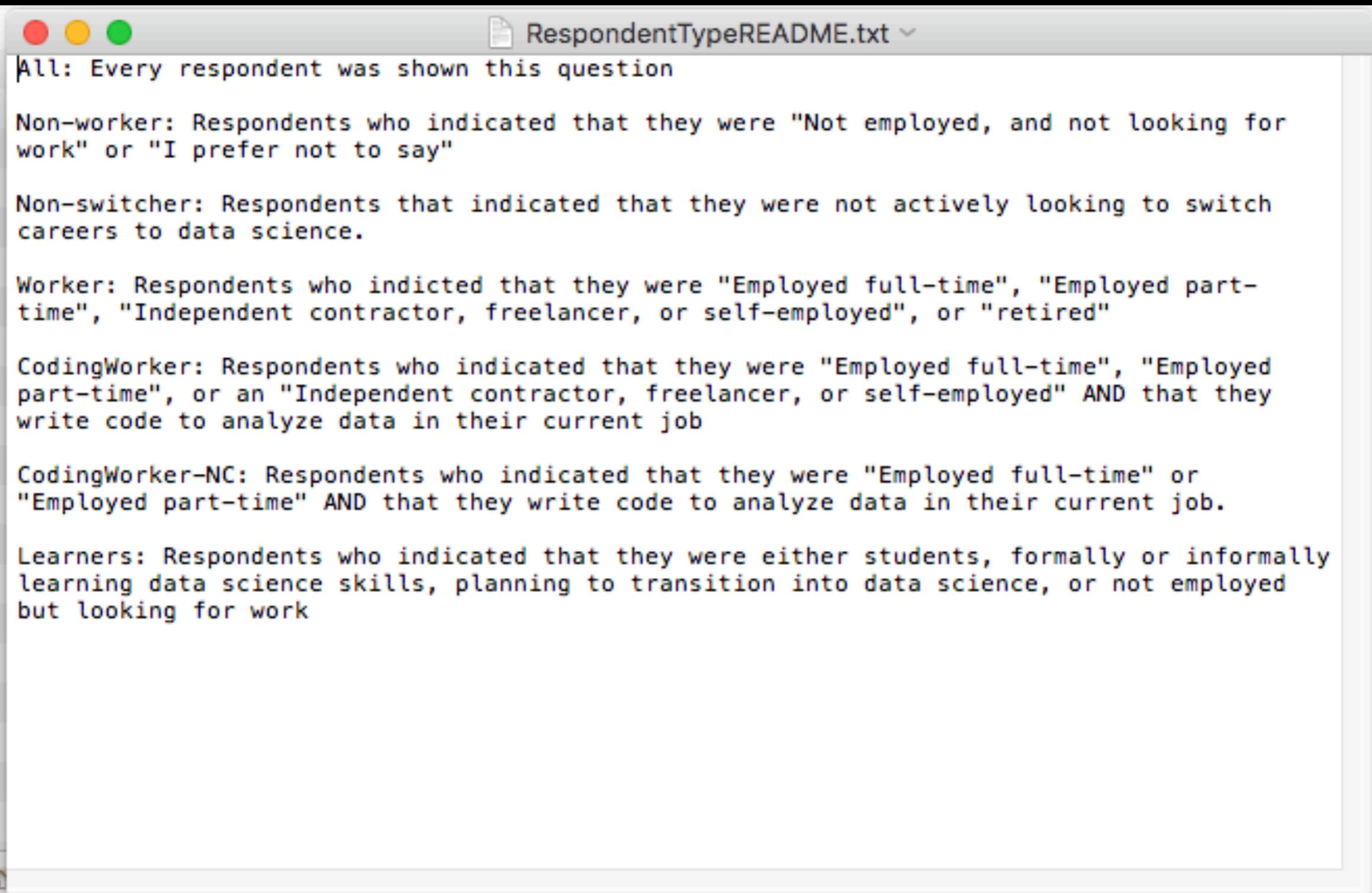
Non-switcher: Respondents that indicated that they were not actively looking to switch careers to data science.

Worker: Respondents who indicated that they were "Employed full-time", "Employed part-time", "Independent contractor, freelancer, or self-employed", or "retired"

CodingWorker: Respondents who indicated that they were "Employed full-time", "Employed part-time", or an "Independent contractor, freelancer, or self-employed" AND that they write code to analyze data in their current job

CodingWorker-NC: Respondents who indicated that they were "Employed full-time" or "Employed part-time" AND that they write code to analyze data in their current job.

Learners: Respondents who indicated that they were either students, formally or informally learning data science skills, planning to transition into data science, or not employed but looking for work



All: Every respondent was shown this question

Non-worker: Respondents who indicated that they were "Not employed, and not looking for work" or "I prefer not to say"

Non-switcher: Respondents that indicated that they were not actively looking to switch careers to data science.

Worker: Respondents who indicated that they were "Employed full-time", "Employed part-time", "Independent contractor, freelancer, or self-employed", or "retired"

CodingWorker: Respondents who indicated that they were "Employed full-time", "Employed part-time", or an "Independent contractor, freelancer, or self-employed" AND that they write code to analyze data in their current job

CodingWorker-NC: Respondents who indicated that they were "Employed full-time" or "Employed part-time" AND that they write code to analyze data in their current job.

Learners: Respondents who indicated that they were either students, formally or informally learning data science skills, planning to transition into data science, or not employed but looking for work

Understand terminology!

Understand terminology!

schema

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format.

Save As...

	A	B	C	D
1	Column	Question		Asked
2	GenderSelect	Select your gender identity. - Selected Choice		All
3	GenderFreeForm	Select your gender identity. - A different identity - Text		All
4	Country	Select the country you currently live in.		All
5	Age	What's your age?		All
6	EmploymentStatus	What's your current employment status?		All
7	StudentStatus	Are you currently enrolled as a student at a degree granting school?		Non-worker
8	LearningDataScience	Are you currently focused on learning data science skills either formally or informally?		Non-worker
9	KaggleMotivationFreeForm	What's your motivation for being a Kaggle user?		Non-switcher
10	CodeWriter	Do you write code to analyze data in your current job, freelance contracts, or most recent job if retired?		Worker1
11	CareerSwitcher	Are you actively looking to switch careers to data science?		Worker1
12	CurrentJobTitleSelect	Select the option that's most similar to your current job/professional title (or most recent title if retired). - Selected Choice		Worker1
13	CurrentJobTitleFreeForm	Select the option that's most similar to your current job/professional title (or most recent title if retired). - Other - Text		Worker1
14	TitleFit	How adequately do you feel your title describes what you do (or what you did if retired)?		Worker1
15	CurrentEmployerType	Which of the following describe your current employer (or most recent employer if retired)? (Select all that apply)		Worker1
16	MLToolNextYearSelect	Which tool or technology are you most excited about learning in the next year? (Select one option) - Selected Choice		All
17	MLToolNextYearFreeForm	Which tool or technology are you most excited about learning in the next year? (Select one option) - Other - Text		All
18	MLMethodNextYearSelect	Which ML/DS method are you most excited about learning in the next year? (Select one option) - Selected Choice		All
19	MLMethodNextYearFreeForm	Which ML/DS method are you most excited about learning in the next year? (Select one option) - Other - Text		All
20	LanguageRecommendationSelect	What programming language would you recommend a new data scientist learn first? (Select one option) - Selected Choice		All
21	LanguageRecommendationFreeForm	What programming language would you recommend a new data scientist learn first? (Select one option) - Other - Text		All
22	PublicDatasetsSelect	Where do you find public datasets to practice data science skills? (Select all that apply) - Selected Choice		All
23	PublicDatasetsFreeForm	Where do you find public datasets to practice data science skills? (Select all that apply) - Other - Text		All
24	PersonalProjectsChallengeFreeForm	What is your biggest challenge with the public datasets you find for personal projects?		All
25	LearningPlatformSelect	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Selected Choice		All
26	LearningPlatformCommunityFreeForm	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Non-Kaggle online communities - Text		All
27	LearningPlatformFreeForm1	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text1		All
28	LearningPlatformFreeForm2	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text2		All
29	LearningPlatformFreeForm3	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text3		All
30	LearningPlatformUsefulnessArxiv	How useful did you find these platforms & resources for learning data science skills? - Arxiv		All
31	LearningPlatformUsefulnessBlogs	How useful did you find these platforms & resources for learning data science skills? - Blogs		All
32	LearningPlatformUsefulnessCollege	How useful did you find these platforms & resources for learning data science skills? - College/University		All
33	LearningPlatformUsefulnessCompany	How useful did you find these platforms & resources for learning data science skills? - Company internal community		All
34	LearningPlatformUsefulnessConferences	How useful did you find these platforms & resources for learning data science skills? - Conferences		All
35	LearningPlatformUsefulnessFriends	How useful did you find these platforms & resources for learning data science skills? - Friends network		All
36	LearningPlatformUsefulnessKaggle	How useful did you find these platforms & resources for learning data science skills? - Kaggle		All
37	LearningPlatformUsefulnessNewsletters	How useful did you find these platforms & resources for learning data science skills? - Newsletters		All
38	LearningPlatformUsefulnessCommunities	How useful did you find these platforms & resources for learning data science skills? - Non-Kaggle online communities		All

schema

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format.

Save As...

	A	B	C	D
1	Column	Question		Asked
2	GenderSelect	Select your gender identity. - Selected Choice		All
3	GenderFreeForm	Select your gender identity. - A different identity - Text		All
4	Country	Select the country you currently live in.		All
5	Age	What's your age?		All
6	EmploymentStatus	What's your current employment status?		All
7	StudentStatus	Are you currently enrolled as a student at a degree granting school?		Non-worker
8	LearningDataScience	Are you currently focused on learning data science skills either formally or informally?		Non-worker
9	KaggleMotivationFreeForm	What's your motivation for being a Kaggle user?		Non-switcher
10	CodeWriter	Do you write code to analyze data in your current job, freelance contracts, or most recent job if retired?		Worker1
11	CareerSwitcher	Are you actively looking to switch careers to data science?		Worker1
12	CurrentJobTitleSelect	Select the option that's most similar to your current job/professional title (or most recent title if retired). - Selected Choice		Worker1
13	CurrentJobTitleFreeForm	Select the option that's most similar to your current job/professional title (or most recent title if retired). - Other - Text		Worker1
14	TitleFit	How adequately do you feel your title describes what you do (or what you did if retired)?		Worker1
15	CurrentEmployerType	Which of the following describe your current employer (or most recent employer if retired)? (Select all that apply)		Worker1
16	MLToolNextYearSelect	Which tool or technology are you most excited about learning in the next year? (Select one option) - Selected Choice		All
17	MLToolNextYearFreeForm	Which tool or technology are you most excited about learning in the next year? (Select one option) - Other - Text		All
18	MLMethodNextYearSelect	Which ML/DS method are you most excited about learning in the next year? (Select one option) - Selected Choice		All
19	MLMethodNextYearFreeForm	Which ML/DS method are you most excited about learning in the next year? (Select one option) - Other - Text		All
20	LanguageRecommendationSelect	What programming language would you recommend a new data scientist learn first? (Select one option) - Selected Choice		All
21	LanguageRecommendationFreeForm	What programming language would you recommend a new data scientist learn first? (Select one option) - Other - Text		All
22	PublicDatasetsSelect	Where do you find public datasets to practice data science skills? (Select all that apply) - Selected Choice		All
23	PublicDatasetsFreeForm	Where do you find public datasets to practice data science skills? (Select all that apply) - Other - Text		All
24	PersonalProjectsChallengeFreeForm	What is your biggest challenge with the public datasets you find for personal projects?		All
25	LearningPlatformSelect	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Selected Choice		All
26	LearningPlatformCommunityFreeForm	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Non-Kaggle online communities - Text		All
27	LearningPlatformFreeForm1	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text1		All
28	LearningPlatformFreeForm2	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text2		All
29	LearningPlatformFreeForm3	What platforms & resources have you used to continue learning data science skills? (Select all that apply) - Other - Text3		All
30	LearningPlatformUsefulnessArxiv	How useful did you find these platforms & resources for learning data science skills? - Arxiv		All
31	LearningPlatformUsefulnessBlogs	How useful did you find these platforms & resources for learning data science skills? - Blogs		All
32	LearningPlatformUsefulnessCollege	How useful did you find these platforms & resources for learning data science skills? - College/University		All
33	LearningPlatformUsefulnessCompany	How useful did you find these platforms & resources for learning data science skills? - Company internal community		All
34	LearningPlatformUsefulnessConferences	How useful did you find these platforms & resources for learning data science skills? - Conferences		All
35	LearningPlatformUsefulnessFriends	How useful did you find these platforms & resources for learning data science skills? - Friends network		All
36	LearningPlatformUsefulnessKaggle	How useful did you find these platforms & resources for learning data science skills? - Kaggle		All
37	LearningPlatformUsefulnessNewsletters	How useful did you find these platforms & resources for learning data science skills? - Newsletters		All
38	LearningPlatformUsefulnessCommunities	How useful did you find these platforms & resources for learning data science skills? - Non-Kaggle online communities		All

Understand possibilities!

×

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve th

A1	Column	A
1	Column	Question
2	GenderSelect	Select your gender identity. - Selected Choice
3	GenderFreeForm	Select your gender identity. - A different identity - Text
4	Country	Select the country you currently live in.
5	Age	What's your age?
6	EmploymentStatus	What's your current employment status?
7	StudentStatus	Are you currently enrolled as a student at a degree granting school?
8	LearningDataScience	Are you currently focused on learning data science skills either formally or informally?
9	KaggleMotivationFreeForm	What's your motivation for being a Kaggle user?
10	CodeWriter	Do you write code to analyze data in your current job, freelance contracts, or most recently?
11	CareerSwitcher	Are you actively looking to switch careers to data science?
12	CurrentJobTitleSelect	Select the option that's most similar to your current job/professional title (or most recent employer)
13	CurrentJobTitleFreeForm	Select the option that's most similar to your current job/professional title (or most recent employer)
14	TitleFit	How adequately do you feel your title describes what you do (or what you did if retired)?
15	CurrentEmployerType	Which of the following describe your current employer (or most recent employer if retired)?
16	MLToolNextYearSelect	Which tool or technology are you most excited about learning in the next year? (Select all that apply)
17	MLToolNextYearFreeForm	Which tool or technology are you most excited about learning in the next year? (Select all that apply)
18	MLMethodNextYearSelect	Which ML/DS method are you most excited about learning in the next year? (Select all that apply)
19	MLMethodNextYearFreeForm	Which ML/DS method are you most excited about learning in the next year? (Select all that apply)
20	LanguageRecommendationSelect	What programming language would you recommend a new data scientist learn first?
21	LanguageRecommendationFreeForm	What programming language would you recommend a new data scientist learn first?
22	PublicDatasetsSelect	Where do you find public datasets to practice data science skills? (Select all that apply)
23	PublicDatasetsFreeForm	Where do you find public datasets to practice data science skills? (Select all that apply)
24	PersonalProjectsChallengeFreeForm	What is your biggest challenge with the public datasets you find for personal projects?
25	LearningPlatformSelect	What platforms & resources have you used to continue learning data science skills?
26	LearningPlatformCommunityFreeForm	What platforms & resources have you used to continue learning data science skills?
27	LearningPlatformFreeForm1	What platforms & resources have you used to continue learning data science skills?

multipleChoiceResponses

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format.

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T
1	GenderSelected	Country	Age	Employment	StudentStatus	LearningData	CodeWriter	CareerSwitch	CurrentJobTitle	TitleFit	CurrentEmpl	MLToolNext	MLMethod	LanguageReq	PublicDataset	LearningPlatform	LearningPlatform	LearningPlatform	LearningPlatform	LearningPlatform
2	Non-binary, genderqueer, NA			Employed full-time			Yes		DBA/Database	Fine	Employed by	SAS Base	Random Forest	F#	Dataset	aggr	College/University,Conferences,Podcasts,Trade book			Very useful
3	Female	United States	30	Not employed, but looking for work							Python	Random Forest	Python		Dataset	aggr	Kaggle			
4	Male	Canada	28	Not employed, but looking for work							Amazon Web	Deep learning	R		Dataset	aggr	Arxiv,College	Very useful		Somewhat useful
5	Male	United States	56	Independent contractor, freelancer, or self-employed	Yes				Operations	Poorly	Self-employed	TensorFlow	Neural Nets	Python	I collect my data	on Blogs,College/University,C	Very useful	Very useful		Very useful
6	Male	Taiwan	38	Employed full-time			Yes		Computer Science	Fine	Employed by	TensorFlow	Text Mining	Python	GitHub	Arxiv,Conferen	Very useful			Somewhat useful
7	Male	Brazil	46	Employed full-time			Yes		Data Scientist	Fine	Employed by	TensorFlow	Genetic & Evolutionary	Python	Dataset	aggr	Kaggle,Online courses,Stack Overflow	Q&A,Textbook		
8	Male	United States	35	Employed full-time			Yes		Computer Science	Fine	Employed by	TensorFlow	Text Mining	R	Dataset	aggr	Arxiv,Blogs,K	Somewhat useful		
9	Female	India	22	Employed full-time			No	Yes	Software Dev	Fine	Employed by	Google Cloud	Deep learning	SQL	Dataset	aggr	College/University,Kaggle,Online cours	Very useful		
10	Female	Australia	43	Employed full-time			Yes		Business Analytics	Fine	Employed by	Microsoft Excel	Link Analysis	Python	University/N	Blogs,Company internal co	Very useful			Very useful
11	Male	Russia	33	Employed full-time			Yes		Software Development	Fine	Employed by	C/C++	Deep learning	Python	Dataset	aggr	Arxiv,Blogs,C	Somewhat useful		Somewhat useful
12	Female	Russia	20	Not employed	Yes	Yes, I'm focused on learning mostly data science skills					Python	Neural Nets	Python	Dataset	aggr	Kaggle,Online courses				
13	Male	India	27	Employed full-time			Yes		Data Scientist	Fine	Employed by	Other	Deep learning	Python	Dataset	aggr	Kaggle,Non-Kaggle online communities,Personal Projects,YouTube Vid			
14	Male	Brazil	26	Employed full-time			No	Yes	Engineer	Fine	Employed by	DataRobot	Deep learning	R	Dataset	aggr	College/University,Conferences,Kaggle,	Somewhat useful		Somewhat useful
15	Male	Netherlands	54	Employed full-time			No	No												
16	Male	Taiwan	26	Employed full-time			Yes		Software Development	Fine	Employed by	TensorFlow	Deep learning	Python	Dataset	aggr	Blogs,Conferences,Kaggle,	Very useful		Very useful
17	Male	United States	58	Independent contractor, freelancer, or self-employed	Yes				DBA/Database	Poorly	Employed by	Python	Rule Induction	R	Dataset	aggr	Kaggle,Personal Projects,Podcasts,Stack Overflow	Q&A,Trade book		
18	Male	Italy	58	Employed full-time			No	No												
19	Male	United Kingdom	24	Employed full-time			No	No												
20	Male	United States	26	Not employed, but looking for work							TensorFlow	Regression	Python		GitHub	Textbook				
21	Male	Brazil	39	Not employed, but looking for work							Python		Python		University/N	College/University,Textbook,Tutoring/n	Very useful			
22	Male	United States	49	Independent contractor, freelancer, or self-employed	No		Yes		Scientist/Researcher	Fine	Self-employed	Amazon Mac	Proprietary API	Java	Google Searc	Online courses,Podcasts				
23	Male	United States	25	Employed part-time			Yes		Researcher	Fine	Employed by	Amazon Mac	Deep learning	Python	Dataset	aggr	Arxiv,College	Somewhat useful		Very useful
24	Male	United States	33	Employed full-time			Yes		Scientist/Researcher	Perfectly	Employed by	R	Deep learning	Matlab	I collect my da	Arxiv,Blogs,C	Somewhat useful	Not Useful		Not Useful
25	Male	Czech Republic	21	Employed part-time			Yes		Other	Fine	Employed by	R	Deep learning	Python	I collect my da	Blogs,College/University,k	Somewhat useful	Not Useful		Not Useful
26	Male	United States	NA	Employed full-time			Yes		Software Development	Fine	Employed by	Spark / MLLib	Deep learning	Matlab	Dataset	aggr	Online courses,Personal Projects,Textbook			
27	Male	Russia	22	Employed full-time			Yes		Data Analyst	Fine	Employed by	TensorFlow	Genetic & Evolutionary	Python	Dataset	aggr	Arxiv,College	Very useful		Somewhat useful
28	Male	Netherlands	51	Employed full-time			Yes		Engineer	Poorly	Employed by	I don't plan to work	I don't plan to work	R	I collect my da	Blogs,Tutoring/mentoring	Very useful			
29	Male	Colombia	34	Employed full-time			Yes		Data Scientist	Fine	Employed by	Spark / MLLib	Ensemble Methods	Python	Google Searc	Online courses,Personal Projects,Stack Overflow	Q&A			
30	Male	Germany	41	Independent contractor, freelancer, or self-employed	Yes				Data Scientist	Fine	Self-employed	I don't plan to work	Factor Analysis	Python	Dataset	aggr	Arxiv,Blogs,C	Very useful	Very useful	Somewhat useful
31	Female	Canada	32	Not employed	Yes	Yes, I'm focused on learning mostly data science skills					Amazon Web	Genetic & Evolutionary	C/C++/C#		Dataset	aggr	Arxiv,College	Somewhat useful		Very useful
32	Male	Denmark	53	Employed full-time			Yes		Business Analyst	Fine	Employed by	professional	Proprietary API	Python	Dataset	aggr	Blogs,Friends network,Personal	Very useful		
33	Male	Poland	29	Employed full-time			Yes		Software Development	Fine	Employed by	TensorFlow	Deep learning	Python	Dataset	aggr	Kaggle,Online courses,Personal Projects,Stack Overflow	Q&A,Textbook		
34	Male	United Kingdom	36	Employed full-time			Yes		Data Scientist	Poorly	Employed by	Microsoft Azure	Proprietary API	Python	University/N	Arxiv,Blogs,S	Very useful	Somewhat useful		Very useful
35	Male	Russia	34	Employed full-time			Yes		Machine Learning	Perfectly	Employed by	Python	Deep learning	Python	Google Searc	Arxiv,Conferen	Somewhat useful			
36	Male	United States	35	Employed full-time			Yes		Engineer	Fine	Employed by	Spark / MLLib	Deep learning	Python	Dataset	aggr	Kaggle,Online courses,Personal Projects,Stack Overflow	Q&A,Textbook		

Some answers are well defined (multiple choice).

AutoSave OFF

multipleC

Home Insert Page Layout Formulas Data Review View

Cut Copy Paste Format

Calibri (Body) 12 A A = = = ab Wrap Text

B I U Gridlines A Merge & Center

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To pres

A1 Formula Bar GenderSelect

	A	B	C	D	E	F	G	H	I	J	K
1	GenderSelected	Country	Age	Employment	StudentStatus	LearningData	CodeWriter	CareerSwitch	CurrentJobTitle	TitleFit	Comments
2	Non-binary, genderqueer, or agender	NA		Employed full-time			Yes		DBA/Database Admin	Fine	
3	Female	United States		30	Not employed, but looking for work						
4	Male	Canada		28	Not employed, but looking for work						
5	Male	United States		56	Independent contractor, freelancer, or self-employed	Yes			Operations Manager	Poorly	
6	Male	Taiwan		38	Employed full-time		Yes		Computer Science	Fine	
7	Male	Brazil		46	Employed full-time		Yes		Data Scientist	Fine	
8	Male	United States		35	Employed full-time		Yes		Computer Science	Fine	
9	Female	India		22	Employed full-time		No	Yes	Software Developer	Fine	
10	Female	Australia		43	Employed full-time		Yes		Business Analyst	Fine	
11	Male	Russia		33	Employed full-time		Yes		Software Developer	Fine	
12	Female	Russia		20	Not employed	Yes	Yes, I'm focused on learning mostly data science skills				
13	Male	India		27	Employed full-time		Yes		Data Scientist	Fine	
14	Male	Brazil		26	Employed full-time		No	Yes	Engineer	Fine	
15	Male	Netherlands		54	Employed full-time		No	No			
16	Male	Taiwan		26	Employed full-time		Yes		Software Developer	Fine	
17	Male	United States		58	Independent contractor, freelancer, or self-employed	Yes			DBA/Database Admin	Poorly	
18	Male	Italy		58	Employed full-time		No	No			
19	Male	United Kingdom		24	Employed full-time		No	No			

multipleChoiceResponses

 Wrap Text



Merge & Center ▾

General ▾

\$

▼

%

,

← .0
.00

.00
→ .0

-delimited (.csv) format. To preserve these features, save it in an E

I

J

K

L

M

\switch CurrentJobTitleFit

CurrentEmployeeMLToolNext\MLMethodN

AutoSave OFF

freeformResponses

Possible Data Loss Some features might be lost if you save this workbook in the comma-delimited (.csv) format. To preserve these features, save it in an Excel file format.

A1 GenderFreeForm

	A	B	C	D	E	F	G	H	I	J	K	L	M	N	O	P	Q	R	S	T	JobS
1	GenderFree	KaggleMotiv	CurrentJobTi	MLToolNext\	MLMethodN	LanguageRe	PublicDatabase	PersonalProj	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	LearningPlat	BlogsPodcas	Jobs
2									Data manipulation			NA		NA		NA		NA		NA	None
3									I can't find time to practice consistently			NA		NA		NA		NA		NA	
4										Meetups			NA		NA		NA		NA		
5										Connectivity/data fusion			NA		NA		NA		NA		
6													NA		NA		NA		NA		
7										kdnuggets	Prepping data		NA		NA		NA		NA		
8													NA		NA		NA		NA		
9											Stanford SNAP		NA		NA		NA		NA		
10													NA		NA		NA		NA		
11											PyTorch		NA		NA		NA		NA		
12													NA		NA		NA		NA		
13													NA		NA		NA		NA		
14													NA		NA		NA		NA		
15													NA		NA		NA		NA		
16													NA		NA		NA		NA		
17													NA		NA		NA		NA		
18													NA		NA		NA		NA		
19													NA		NA		NA		NA		
20											Curious		NA		NA		NA		NA		
21											Hydrographic Surveyor		NA		NA		NA		NA		
22													NA		NA		NA		NA		
23													Poor data quality / lack of documentation		NA		NA		NA		
24														NA		NA		NA			
25											mechanical engineer	don't know		NA		NA		NA		NA	
26											Promote our	Technical support engineer		NA		NA		NA		NA	
27														NA		NA		NA		NA	
28														NA		NA		NA		NA	
29													Crawling Airbnb		NA		NA		NA		
30														NA		NA		NA			
31													Amount and quality of data		NA		NA		NA		
32														NA		NA		NA			
33														NA		NA		NA			
34											half man - half dog		None in specific.		NA		NA		NA		
35														NA		NA		NA			

Some answers are not well defined (free form).

24

25 mechanical engineer

26 Promote our Technical support engineer

27

28

29

30

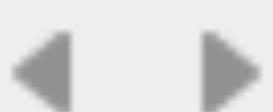
31

32 Gender (free form).

33

34 half man - half dog

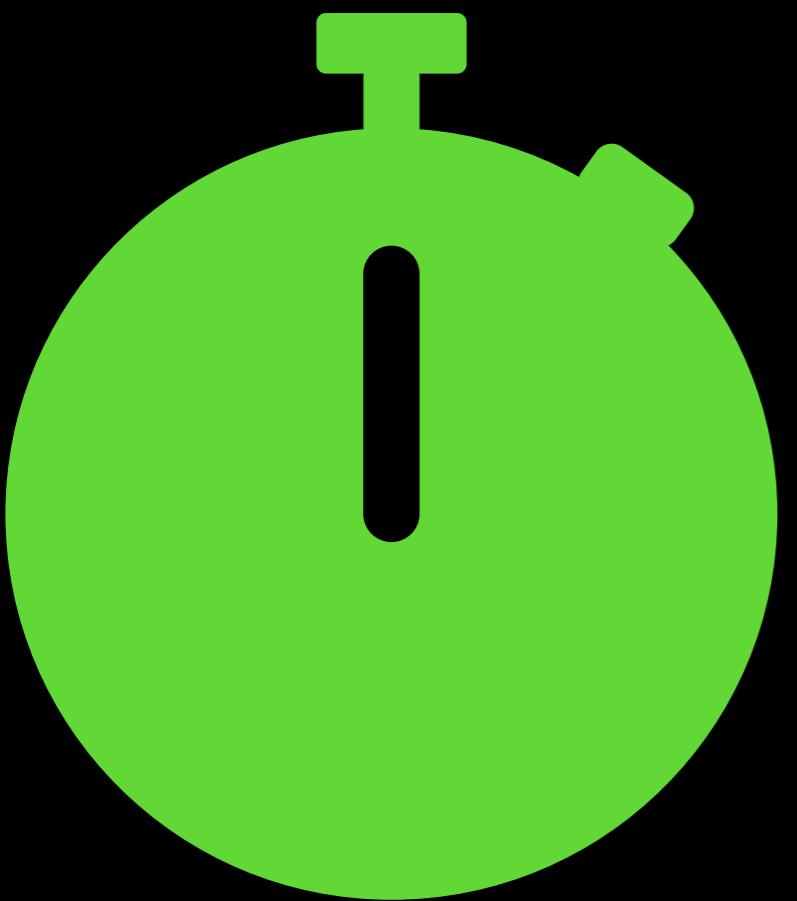
35

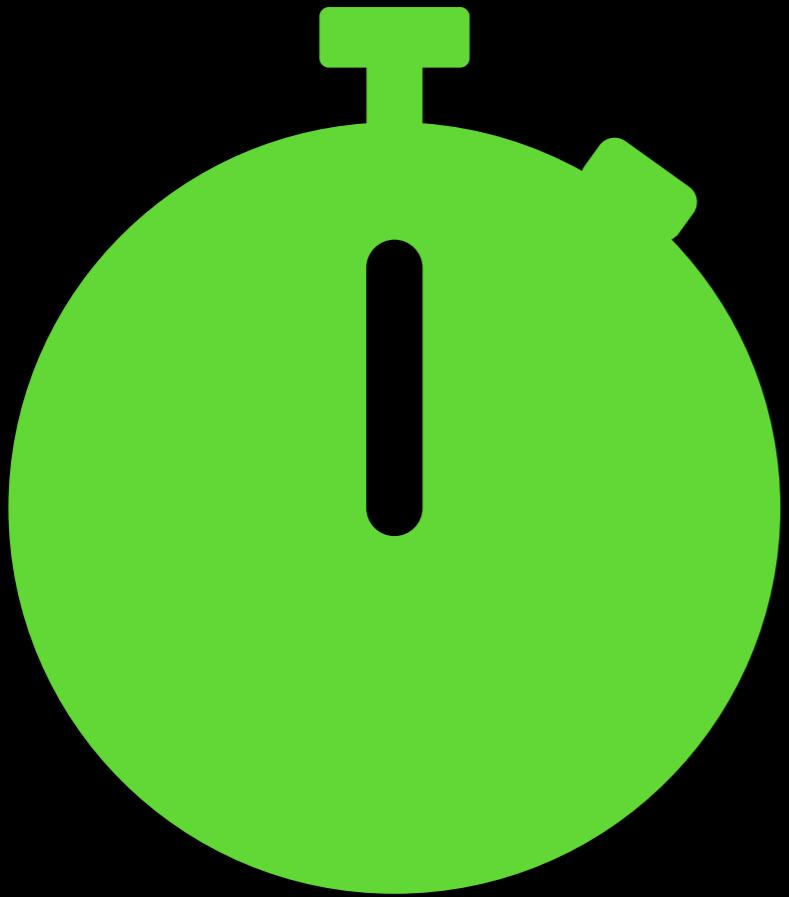


freeformResponses



Ready





Find another neighbor.

How does the Kaggle dataset help users to understand the meaning of the data?



Find another neighbor.

How does the Kaggle dataset help users to understand the meaning of the data?

Do academic datasets usually come with good descriptions?



Let us pretend we don't know much about the data.

```
In [1]: %matplotlib inline
```

```
In [2]: cd '/Users/tstoeger/Dropbox/Work/kaggle_survey'  
/Users/tstoeger/Dropbox/Work/kaggle_survey
```

```
In [3]: import pandas as pd  
import seaborn as sns
```

```
In [1]: %matplotlib inline
```

```
In [2]: cd '/Users/tstoeger/Dropbox/Work/kaggle_survey'
```

/Users/tstoeger/Dropbox/Work/kaggle_survey

```
In [3]: import pandas as pd  
import seaborn as sns
```

```
In [4]: df_multiple_choice = pd.read_csv(  
    './multipleChoiceResponses.csv',  
    encoding='latin-1', # files have some special characters  
    low_memory=False) # some columns have mixed type
```

```
In [1]: %matplotlib inline
```

```
In [2]: cd '/Users/tstoeger/Dropbox/Work/kaggle_survey'
```

/Users/tstoeger/Dropbox/Work/kaggle_survey

```
In [3]: import pandas as pd  
import seaborn as sns
```

```
In [4]: df_multiple_choice = pd.read_csv(  
    './multipleChoiceResponses.csv',  
    encoding='latin-1', # files have some special characters  
    low_memory=False) # some columns have mixed type
```

```
In [5]: df_multiple_choice.shape # get idea of number of datasets
```

```
Out[5]: (16716, 228)
```

```
In [1]: %matplotlib inline
```

```
In [2]: cd '/Users/tstoeger/Dropbox/Work/kaggle_survey'
```

/Users/tstoeger/Dropbox/Work/kaggle_survey

```
In [3]: import pandas as pd  
import seaborn as sns
```

```
In [4]: df_multiple_choice = pd.read_csv(  
    './multipleChoiceResponses.csv',  
    encoding='latin-1', # files have some special characters  
    low_memory=False) # some columns have mixed type
```

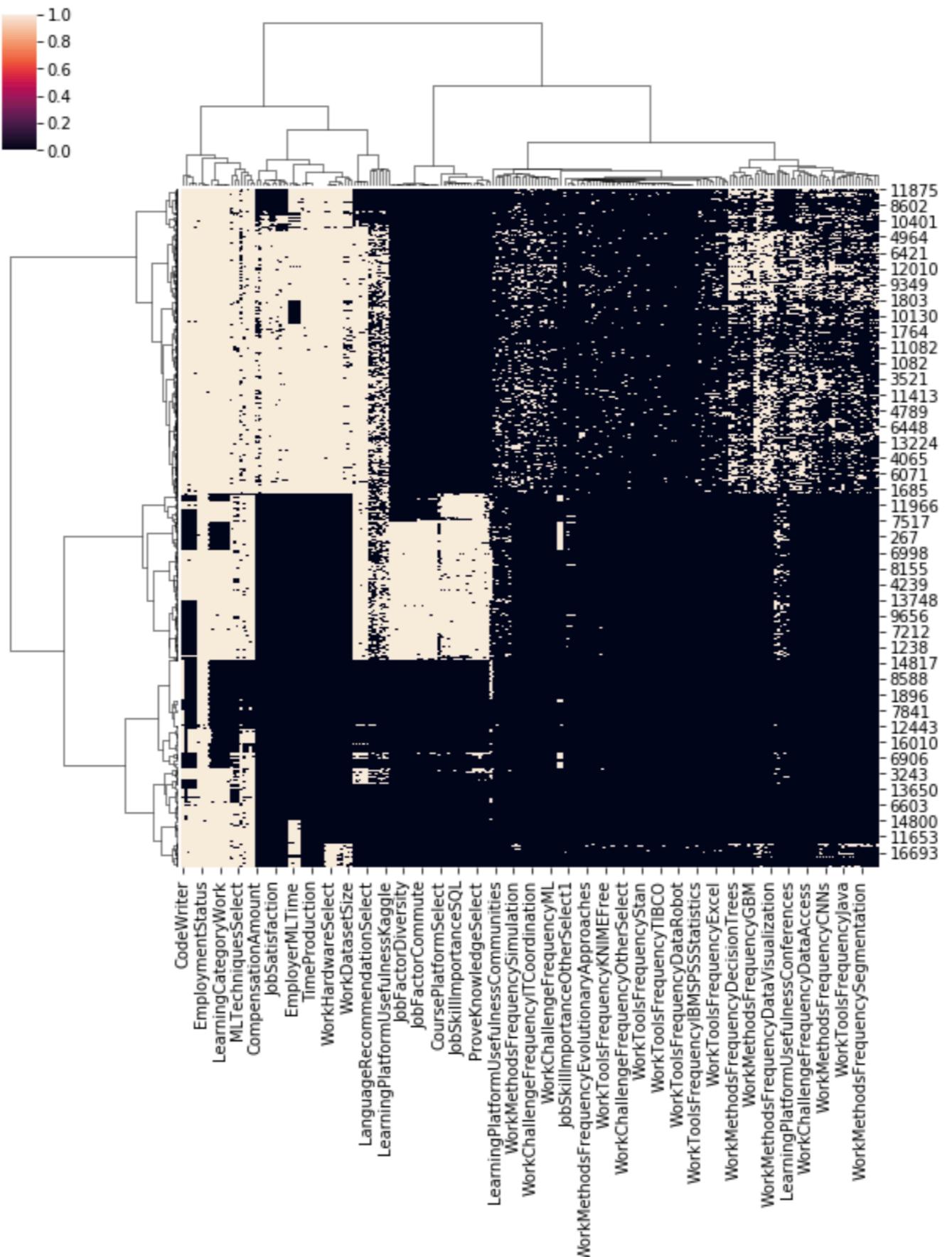
```
In [5]: df_multiple_choice.shape # get idea of number of datasets
```

Out[5]: (16716, 228)

```
In [6]: sns.clustermap( # visualize presence of data  
    df_multiple_choice.notnull(),  
    method='ward')
```

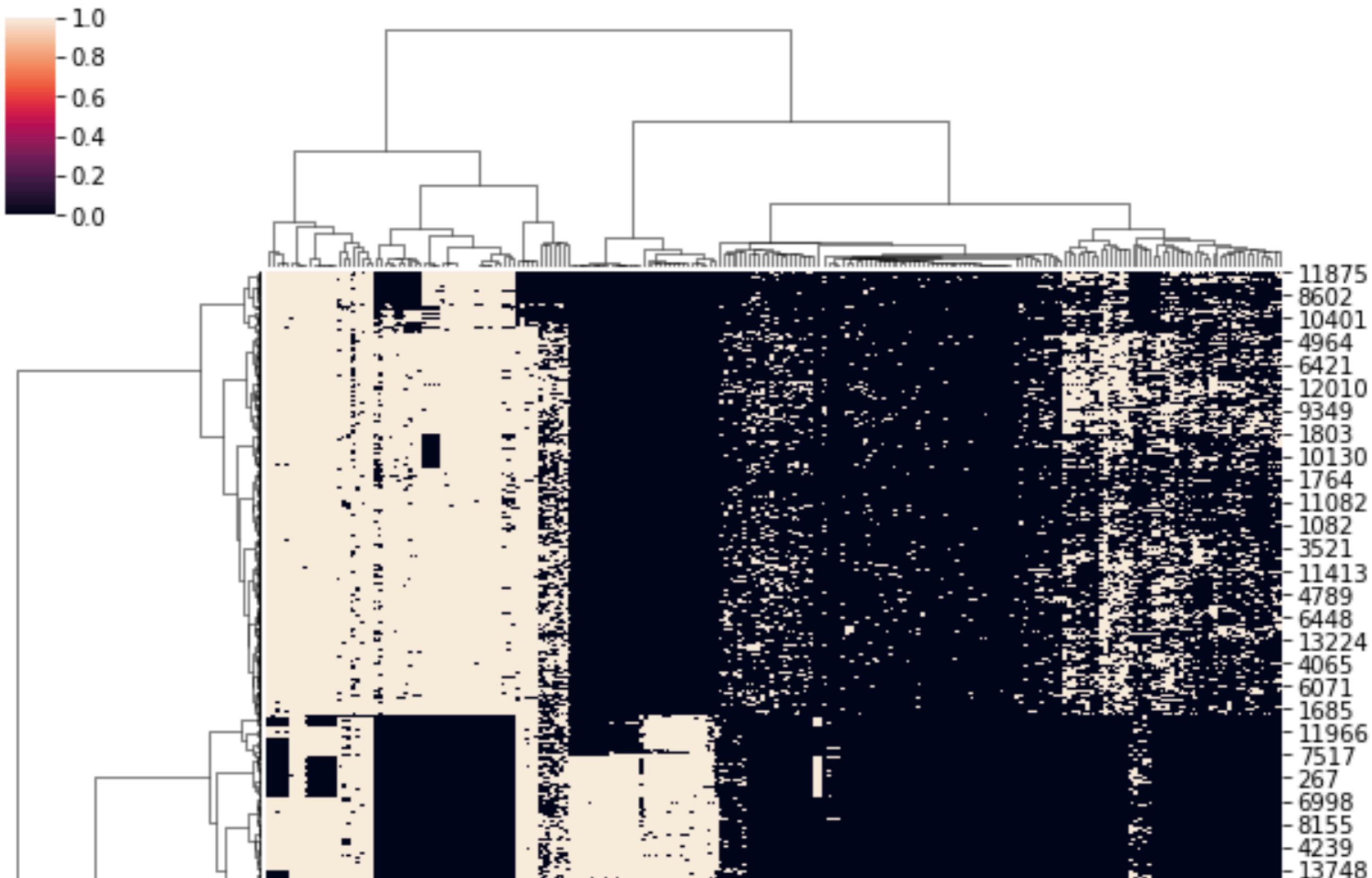
```
In [6]: sns.clustermap(          # visualize presence of data
    df_multiple_choice.notnull(),
    method='ward')
```

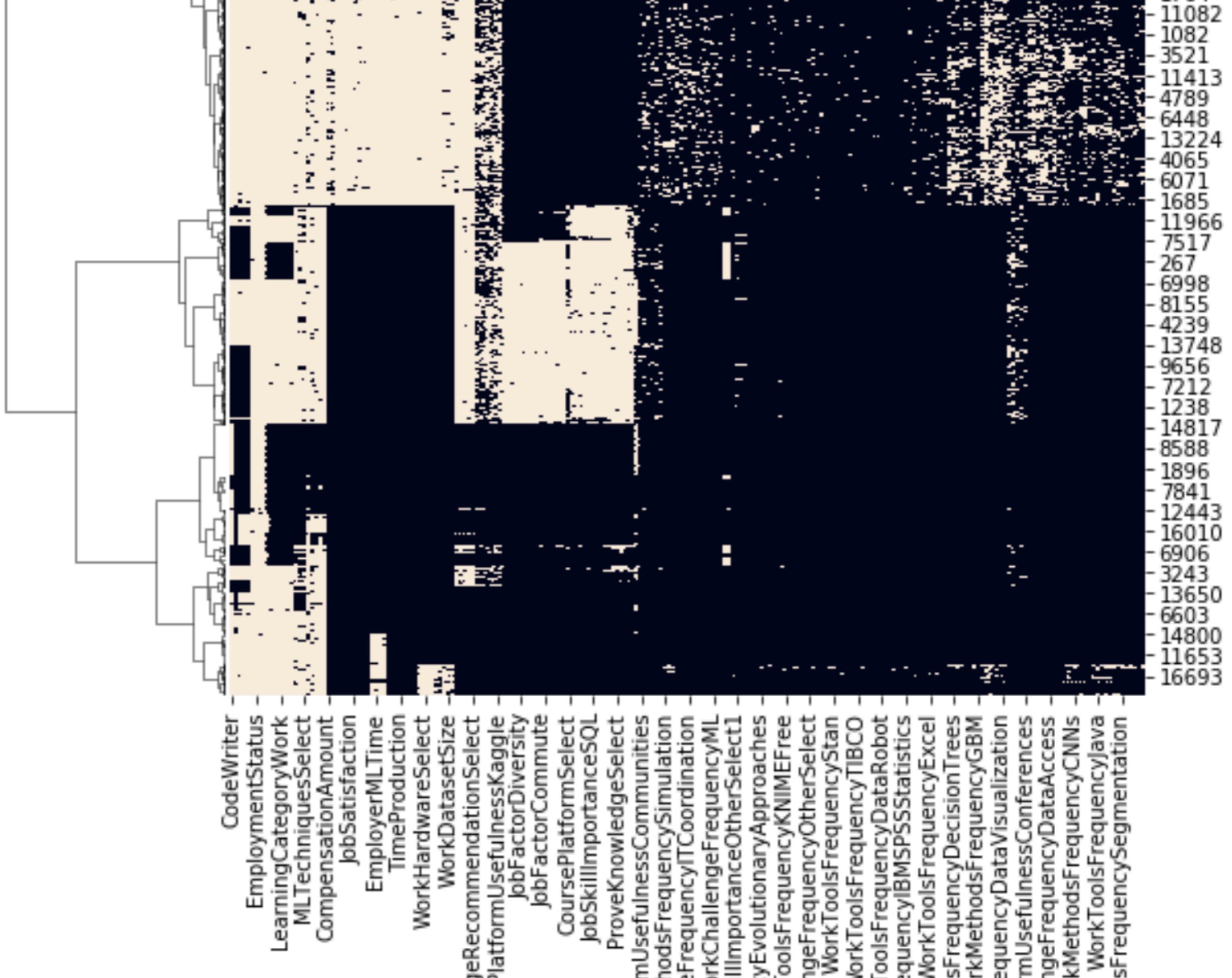
```
Out[6]: <seaborn.matrix.ClusterGrid at 0x10cd02a58>
```



```
sns.clustermap(          # visualize presence of data
    df_multiple_choice.notnull(),
    method='ward')
```

```
<seaborn.matrix.ClusterGrid at 0x10cd02a58>
```





Always look at data!

In [7]: df_multiple_choice

Out[7]:

	GenderSelect	Country	Age	EmploymentStatus	StudentStatus	LearningDataScience
0	Non-binary, genderqueer, or gender non- conforming		NaN	NaN	Employed full-time	NaN
1	Female	United States	30.0	Not employed, but looking for work	NaN	NaN
2	Male	Canada	28.0	Not employed, but looking for work	NaN	NaN
3	Male	United States	56.0	Independent contractor, freelancer, or self- em...	NaN	NaN
4	Male	Taiwan	38.0	Employed full-time	NaN	NaN
5	Male	Brazil	46.0	Employed full-time	NaN	NaN
6	Male	United States	35.0	Employed full-time	NaN	NaN
7	Female	India	22.0	Employed full-time	NaN	NaN
8	Female	Australia	43.0	Employed full-time	NaN	NaN
9	Male	Russia	33.0	Employed full-time	NaN	NaN

Always look at data!

```
In [7]: df_multiple_choice
```

```
Out[7]:
```

	GenderSelect	Country	Age	EmploymentStatus	StudentStatus	LearningDataScience
0	Non-binary, genderqueer, or gender non- conforming		NaN	NaN	Employed full-time	NaN
1	Female	United States	30.0			NaN
2	Male	Canada	29.0			NaN
3	Male	United States	56.0	free-lance		NaN
4	Male	Taiwan	38.0	Employed full-time		NaN

In a nicely organized data set, columns are organized in a meaningful way

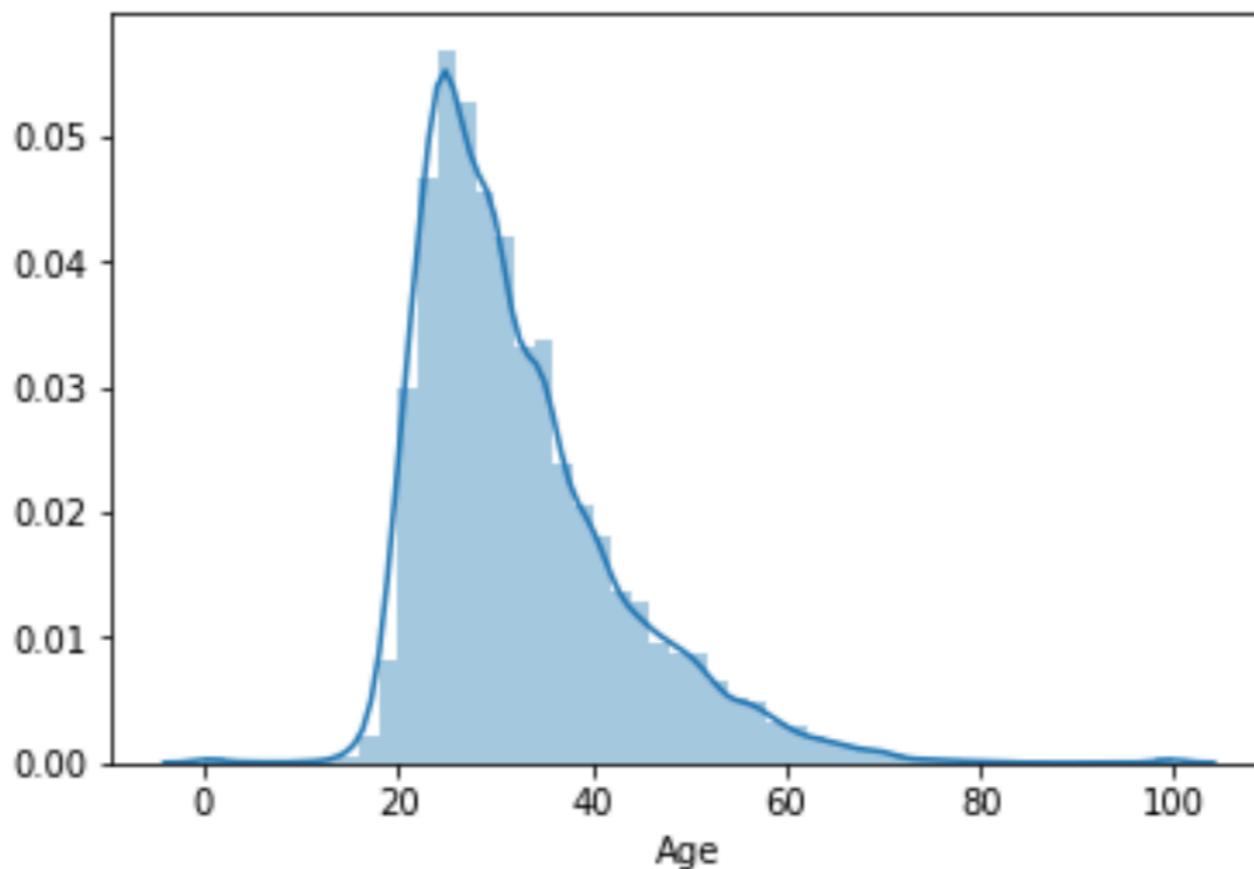
7	Female	India	22.0	Employed full-time	NaN	NaN
8	Female	Australia	43.0	Employed full-time	NaN	NaN
9	Male	Russia	33.0	Employed full-time	NaN	NaN

```
In [8]: sns.distplot(  
    df_multiple_choice['Age'].dropna(),  
)
```

Sanity-check numbers, and
test for unusual values.

```
In [8]: sns.distplot(  
    df_multiple_choice[ 'Age' ].dropna() ,  
)
```

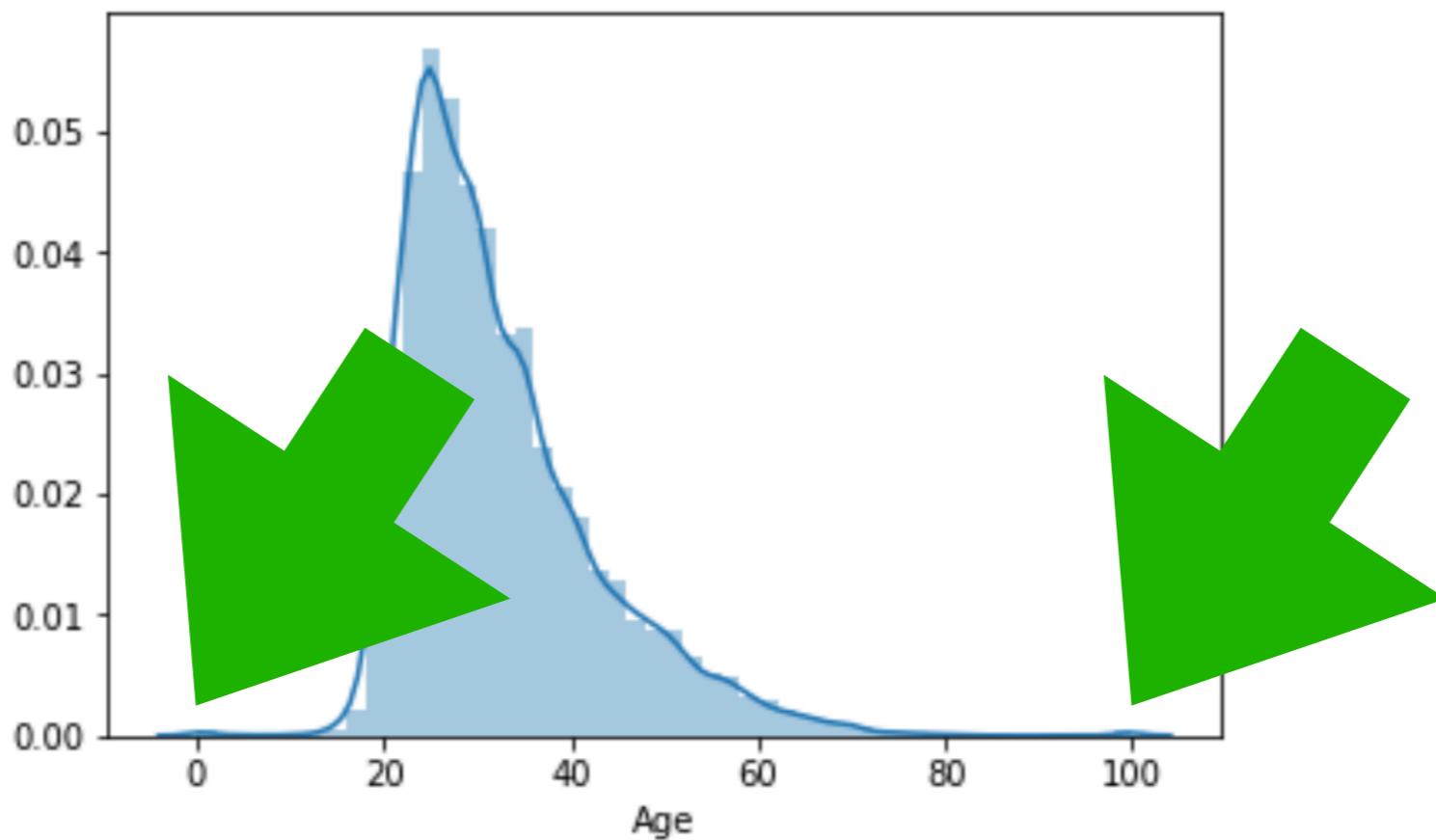
```
Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x1a124307f0>
```



Sanity-check numbers.

```
In [8]: sns.distplot(  
    df_multiple_choice[ 'Age' ].dropna() ,  
)
```

```
Out[8]: <matplotlib.axes._subplots.AxesSubplot at 0x1a124307f0>
```



Sanity-check numbers.

Workflow

- Anticipate problems
- **Get overview**
- Clean data
- Find connections

Setup

- Environment
- Organizing data
- Coding

break

Work

What are two strategies?

- Anticipate problems
- **Get overview**
- Clean data
- Find connections
- Environment
- Organizing data
- Coding

break

*Data science is: 80% data cleaning and
20% complaining about data cleaning.*



PersonalProjectsChallenge

Download this workbook in the comma-delimited (.csv) format. To preserve these features:

You are not alone.

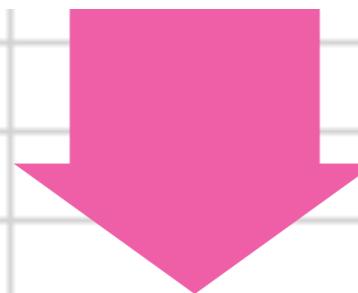


F	G	I	J	K	L
languageRe	PublicDatabase	PersonalProj	LearningPlat	LearningPlat	LearningPlat
		Data manipulation			
		I can't find time to practice consistently			
			Meetups		
		Connectivity/data fusion			
kdnuggets	Prepping data				
	Stanford SNAP				

Personal Projects Challenge

You are not alone.

Don't know



Poor data quality / lack of documentation

Crawling Airbnb

Amount and quality of data

None in specific.

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
- Regular expressions.
- Entity matching.

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
- Regular expressions.
- Entity matching.

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
- Regular expressions.
- Entity matching.

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
- Regular expressions.
- Entity matching.



Salaries of data scientists
(part of Kaggle Survey)

```
In [13]: # For the sake of this example, let us only consider  
# salaries that have been converted to USD  
f = df_multiple_choice['CompensationCurrency'] == 'USD'  
df_multiple_choice = df_multiple_choice[f]
```

```
In [13]: # For the sake of this example, let us only consider  
# salaries that have been converted to USD  
f = df_multiple_choice['CompensationCurrency'] == 'USD'  
df_multiple_choice = df_multiple_choice[f]
```

```
In [14]: # For the sake of this example, let us only consider  
# records, where CompensationAmount is defined  
# (we still keep records, if something else, e.g.  
# StudentStatus, would be not defined)  
df_multiple_choice = df_multiple_choice.dropna(  
    subset=['CompensationAmount'])  
)
```

```
In [13]: # For the sake of this example, let us only consider  
# salaries that have been converted to USD  
f = df_multiple_choice['CompensationCurrency'] == 'USD'  
df_multiple_choice = df_multiple_choice[f]
```

```
In [14]: # For the sake of this example, let us only consider  
# records, where CompensationAmount is defined  
# (we still keep records, if something else, e.g.  
# StudentStatus, would be not defined)  
df_multiple_choice = df_multiple_choice.dropna(  
    subset=['CompensationAmount'])  
)
```

```
In [15]: # Manually inspect the data in CompensationAmount  
df_multiple_choice['CompensationAmount']
```

```
Out[15]: 3      250,000  
21      20000  
22      100000  
34      133000  
37      80000  
61      15000  
75      215000  
86      83500
```

In [16]: # Let us convert these values to a number
df_multiple_choice['CompensationAmount'] = df_multiple_choice['CompensationAmount'].astype(float)

Error

```
# Let us convert these values to a number
df_multiple_choice['CompensationAmount'] = df_multiple_choice[
    'CompensationAmount'].astype(float)

-----
ValueError                                Traceback (most recent call last)
<ipython-input-16-f163d2412d49> in <module>()
      1 df_multiple_choice['CompensationAmount'] = df_multiple_choice[
----> 2     'CompensationAmount'].astype(float)

~/anaconda3/lib/python3.6/site-packages/pandas/util/_decorators.py in wrapper(*args, **kwargs)
   176         else:
   177             kwargs[new_arg_name] = new_arg_value
--> 178         return func(*args, **kwargs)
   179     return wrapper
   180 return _deprecate_kwarg

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals/dtypes.py in astype(self, dtype, copy, errors, **kwargs)
  4995         # else, only a single dtype is given
  4996         new_data = self._data.astype(dtype=dtype, copy=copy, errors=errors,
--> 4997                               **kwargs)
  4998         return self._constructor(new_data).__finalize__(self)
  4999

~/anaconda3/lib/python3.6/site-packages/pandas/core/frame.py in astype(self, dtype, copy, errors, **kwargs)
  3712     def astype(self, dtype, **kwargs):
--> 3714         return self.apply('astype', dtype=dtype, **kwargs)
  3715
  3716     def convert(self, **kwargs):

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in apply(self, f, axis, skipna, convert_dtype, convert_index, copy, do_integrity_check, consolidate, **kwargs)
  3579
  3580         kwargs['mgr'] = self
--> 3581         applied = getattr(b, f)(**kwargs)
  3582         result_blocks = _extend_blocks(applied, result_blocks)
  3583

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in astype(self, dtype, copy, errors, values, **kwargs)
  573     def astype(self, dtype, copy=False, errors='raise', values=None):
  574         return self._astype(dtype, copy=copy, errors=errors,
--> 575                           **kwargs)
  576
  577     def _astype(self, dtype, copy=False, errors='raise', values=None):

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in _astype(self, dtype, copy, errors, values, klass, mgr, **kwargs)
  662
  663         # _astype_nansafe works fine here
--> 664         values = astype_nansafe(values, dtype=dtype)
  665         values = values.reshape(self.shape)
  666

~/anaconda3/lib/python3.6/site-packages/pandas/core/dtypes/creation.py in astype(arr, dtype, copy)
  728
  729     if copy:
--> 730         return arr.astype(dtype, copy=True)
  731     return arr.view(dtype)
  732

ValueError: could not convert string to float: '85,000'
```



Only fix what you want to fix. (85,000 and similar)

```
ValueError: could not convert string to float: '85,000'
```

```
In [17]: f = df_multiple_choice['CompensationAmount'].str.contains(  
    '[0-9]*,[0-9]{3}$$') # create a highly specific regular expression
```

[0-9]*,[0-9]{3}\$\$

Only fix what you want to fix. (85,000 and similar)

```
ValueError: could not convert string to float: '85,000'
```

```
In [17]: f = df_multiple_choice['CompensationAmount'].str.contains(  
    '[0-9]*,[0-9]{3}$$') # create a highly specific regular expression
```

Can someone translate this?

[0-9]*,[0-9]{3}\$\$

Only fix what you want to fix. (85,000 and similar)

```
ValueError: could not convert string to float: '85,000'
```

```
In [17]: f = df_multiple_choice['CompensationAmount'].str.contains(  
    '[0-9]*,[0-9]{3}$') # create a highly specific regular expression
```

```
In [18]: df_multiple_choice.loc[f, 'CompensationAmount'] = df_multiple_choice.loc[  
    f, 'CompensationAmount'].str.replace(',', '')
```

Can someone translate this?

[0-9]*,[0-9]{3}\$

Only fix what you want to fix. (85,000 and similar)

```
ValueError: could not convert string to float: '85,000'
```

```
In [17]: f = df_multiple_choice['CompensationAmount'].str.contains(  
    '[0-9]*,[0-9]{3}$') # create a highly specific regular expression
```

```
In [18]: df_multiple_choice.loc[f, 'CompensationAmount'] = df_multiple_choice.loc[  
    f, 'CompensationAmount'].str.replace(',', '')
```

Can someone translate this?

[0-9]*,[0-9]{3}\$

Trick question: Could this be shorter to
solve our problem?

Cleaning data is an interactive process

```
In [19]: df_multiple_choice['CompensationAmount'] = df_multiple_choice['CompensationAmount'].astype(float)

-----
ValueError                                                 Traceback (most recent call last)
<ipython-input-19-f163d2412d49> in <module>()
      1 df_multiple_choice['CompensationAmount'] = df_multiple_choice[
----> 2     'CompensationAmount'].astype(float)

~/anaconda3/lib/python3.6/site-packages/pandas/util/_decorators.py in wrapper(*args, **kwargs)
    176         else:
    177             kwargs[new_arg_name] = new_arg_value
--> 178     return func(*args, **kwargs)
    179     return wrapper
    180 return _deprecate_kwarg

~/anaconda3/lib/python3.6/site-packages/pandas/core/generic.py in astype(self, dtype, copy, errors, **kwargs)
    4995         # else, only a single dtype is given
    4996         new_data = self._data.astype(dtype=dtype, copy=copy, errors=errors,
--> 4997                         **kwargs)
    4998         return self._constructor(new_data).__finalize__(self)
    4999

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in astype(self, dtype, **kwargs)
    3712
    3713     def astype(self, dtype, **kwargs):
--> 3714         return self.apply('astype', dtype=dtype, **kwargs)
    3715
    3716     def convert(self, **kwargs):

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in apply(self, f, axes, filter, do_integrity_check,
consolidate, **kwargs)
    3579
    3580         kwargs['mgr'] = self
--> 3581         applied = getattr(b, f)(**kwargs)
    3582         result_blocks = _extend_blocks(applied, result_blocks)
    3583

~/anaconda3/lib/python3.6/site-packages/pandas/core/internals.py in astype(self, dtype, copy, errors, values, **kwargs)
    573     def astype(self, dtype, copy=False, errors='raise', values=None, **kwargs):
    574         return self._astype(dtype, copy=copy, errors=errors, values=values,
--> 575                         **kwargs)
    576
    577     def _astype(self, dtype, copy=False, errors='raise', values=None,
```

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
(e.g.: assert, raise Error)
- Regular expressions.
- Entity matching.

e.g: match Thomas Stoeger to Thomas Stoger

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
(e.g.: assert, raise Error)
- Regular expressions.
- Entity matching.

e.g: match Thomas Stoeger to Thomas Stoger

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
(e.g.: assert, raise Error)
- Regular expressions.
- Entity matching.

e.g: match Thomas Stoeger to Thomas Stoger

Main tactics for data cleaning.

- Codify anticipation of unanticipated.
(e.g.: assert, raise Error)
- Regular expressions.
- Entity matching.

e.g: match Thomas Stoeger to Thomas Stoger

Missing data can be tricky.

Advance Access publication August 22, 2016

Political Analysis (2016) 24:414–433
doi:10.1093/pan/mpw020

How Multiple Imputation Makes a Difference

Ranjit Lall

Department of Government, Harvard University, 1737 Cambridge Street, Cambridge, MA 02138
e-mail: ranjitlall@fas.harvard.edu (corresponding author)

Edited by R. Michael Alvarez

Political scientists increasingly recognize that multiple imputation represents a superior strategy for analyzing missing data to the widely used method of listwise deletion. However, there has been little systematic investigation of how multiple imputation affects existing empirical knowledge in the discipline. This article presents the first large-scale examination of the empirical effects of substituting multiple imputation for listwise deletion in political science. The examination focuses on research in the major subfield of comparative and international political economy (CIPE) as an illustrative example. Specifically, I use multiple imputation to reanalyze the results of almost every quantitative CIPE study published during a recent five-year period in *International Organization* and *World Politics*, two of the leading subfield journals in CIPE. The outcome is striking: in almost half of the studies, key results “disappear” (by conventional statistical standards) when reanalyzed.

Missing data can be tricky.

Advance Access publication August 22, 2016

Political Analysis (2016) 24:414–433
doi:10.1093/pan/mpw020

How Multiple Imputation Makes a Difference

Ranjit Lall

Department of Government, Harvard University
e-mail: ranjitlall@fas.harvard.edu

Edited by R.

Political scientists increasingly recognize that multiple imputation is a better way to handle missing data than the widely used method of listwise deletion. This article provides the first large-scale investigation of how multiple imputation affects estimates in political science. The examination focuses on comparative and international political economy (CIP) studies. The article uses multiple imputation to reanalyze the results of almost every study published in CIP journals over a 20-year period in *International Organization* and *World Politics*. The main finding is striking: in almost half of the studies, the results change when multiple imputation is used instead of listwise deletion. In some cases, the results change dramatically (ards) when reanalyzed.

Sometimes, there are smart people who actively work on misleading others.

A rough guide to missing data:

Random

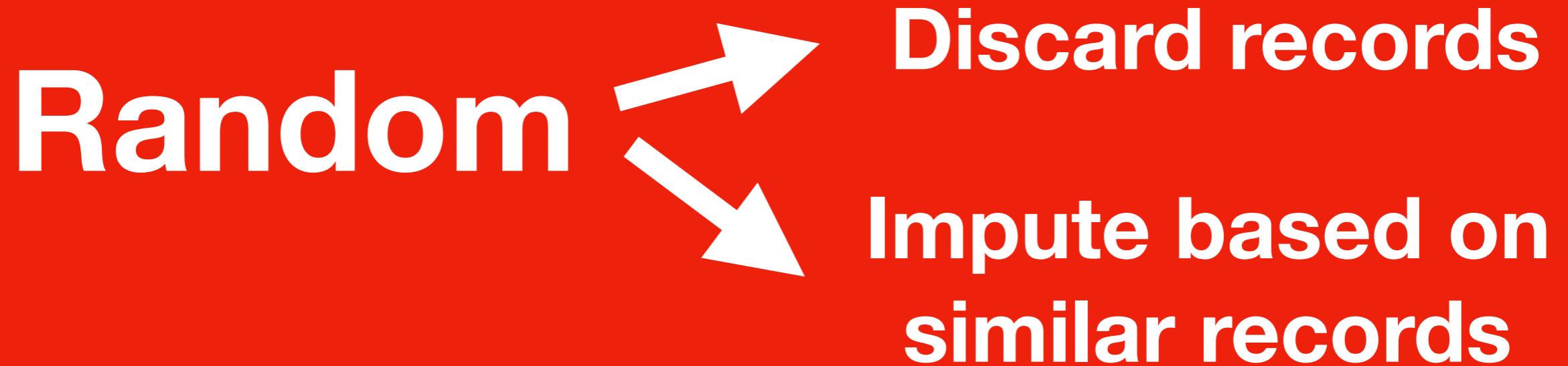
Non-random

A rough guide to missing data:

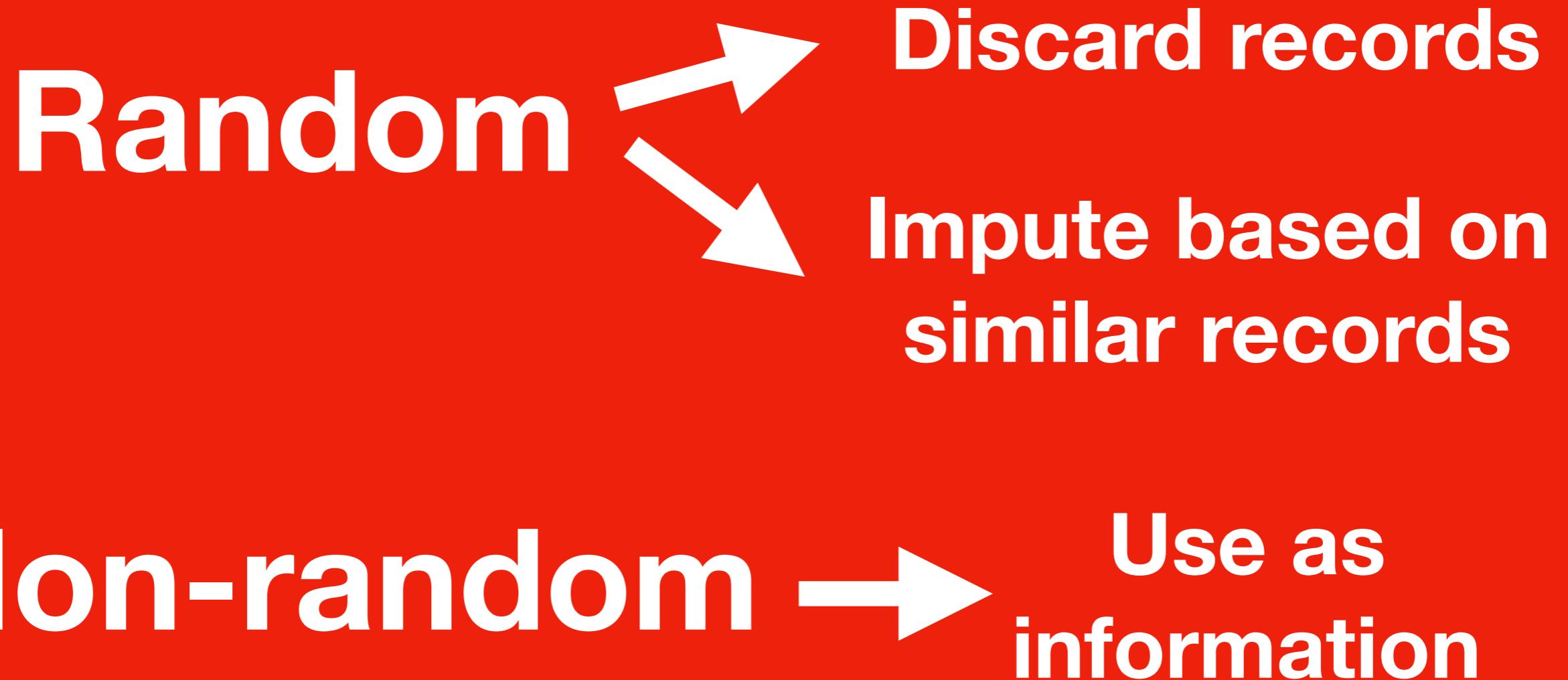
Random → Discard records

Non-random

A rough guide to missing data:



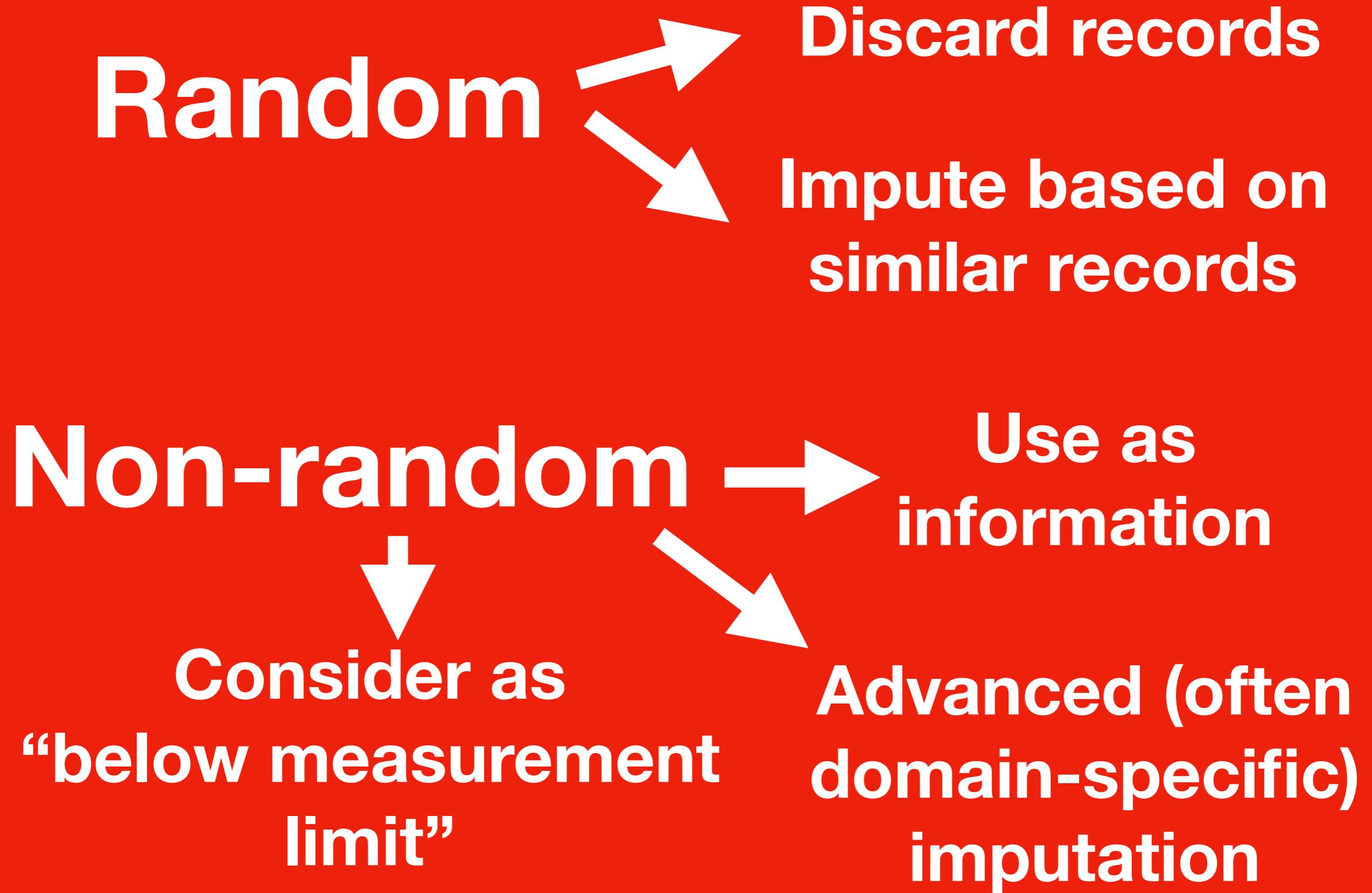
A rough guide to missing data:

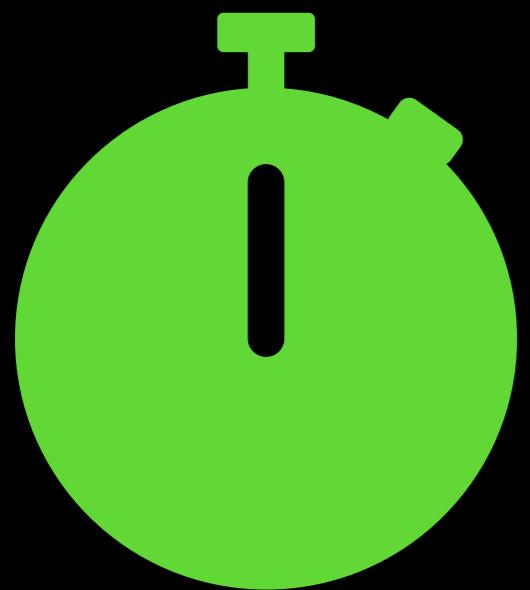


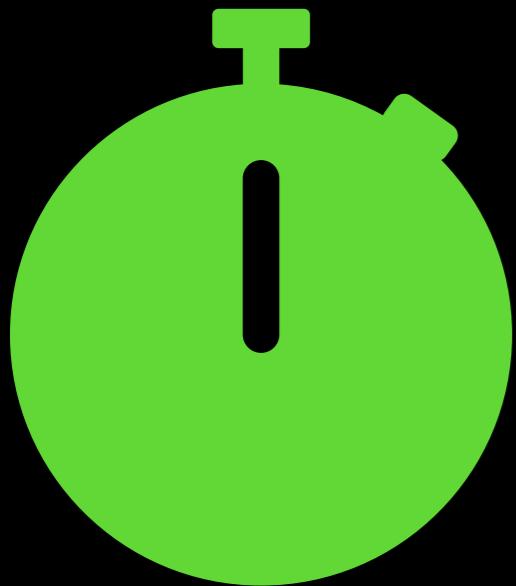
A rough guide to missing data:



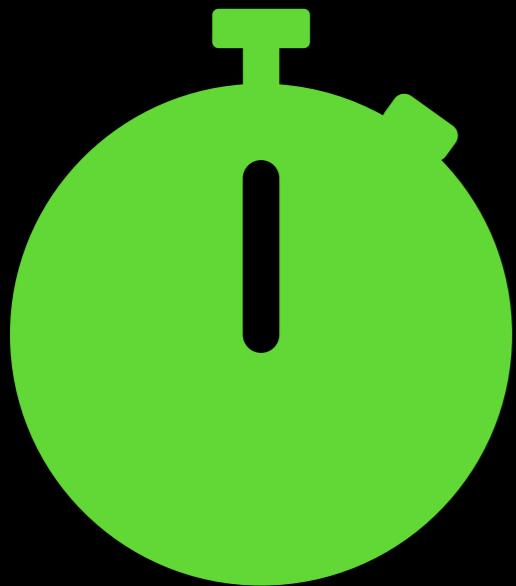
A rough guide to missing data:







**Remember your worst
encounter with data.**



**Remember your worst
encounter with data.**

**How should the data have been cleaned
to simplify your work?**

Workflow

- Anticipate problems
- Get overview
- Clean data
- **Find connections**

break

Setup

- Computer
- Organizing data
- Coding

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

Is there a gender bias in the salary of data scientists?

Is there a gender bias in the salary of data scientists?

```
In [108]: f = (
    df_multiple_choice['GenderSelect'].isin(['Male', 'Female'])) & (
    df_multiple_choice['Country'].isin(['United States'])) & (
    df_multiple_choice['EmploymentStatus'].isin(['Employed full-time'])) & (
    df_multiple_choice['Age'].isin(range(20, 35))) & (
    df_multiple_choice['CurrentJobTitleSelect'].isin(['Data Scientist']))
)
```

Is there a gender bias in the salary of data scientists?

```
In [108]: f = (
    df_multiple_choice['GenderSelect'].isin(['Male', 'Female'])) & (
    df_multiple_choice['Country'].isin(['United States'])) & (
    df_multiple_choice['EmploymentStatus'].isin(['Employed full-time'])) & (
    df_multiple_choice['Age'].isin(range(20, 35))) & (
    df_multiple_choice['CurrentJobTitleSelect'].isin(['Data Scientist']))
)
```

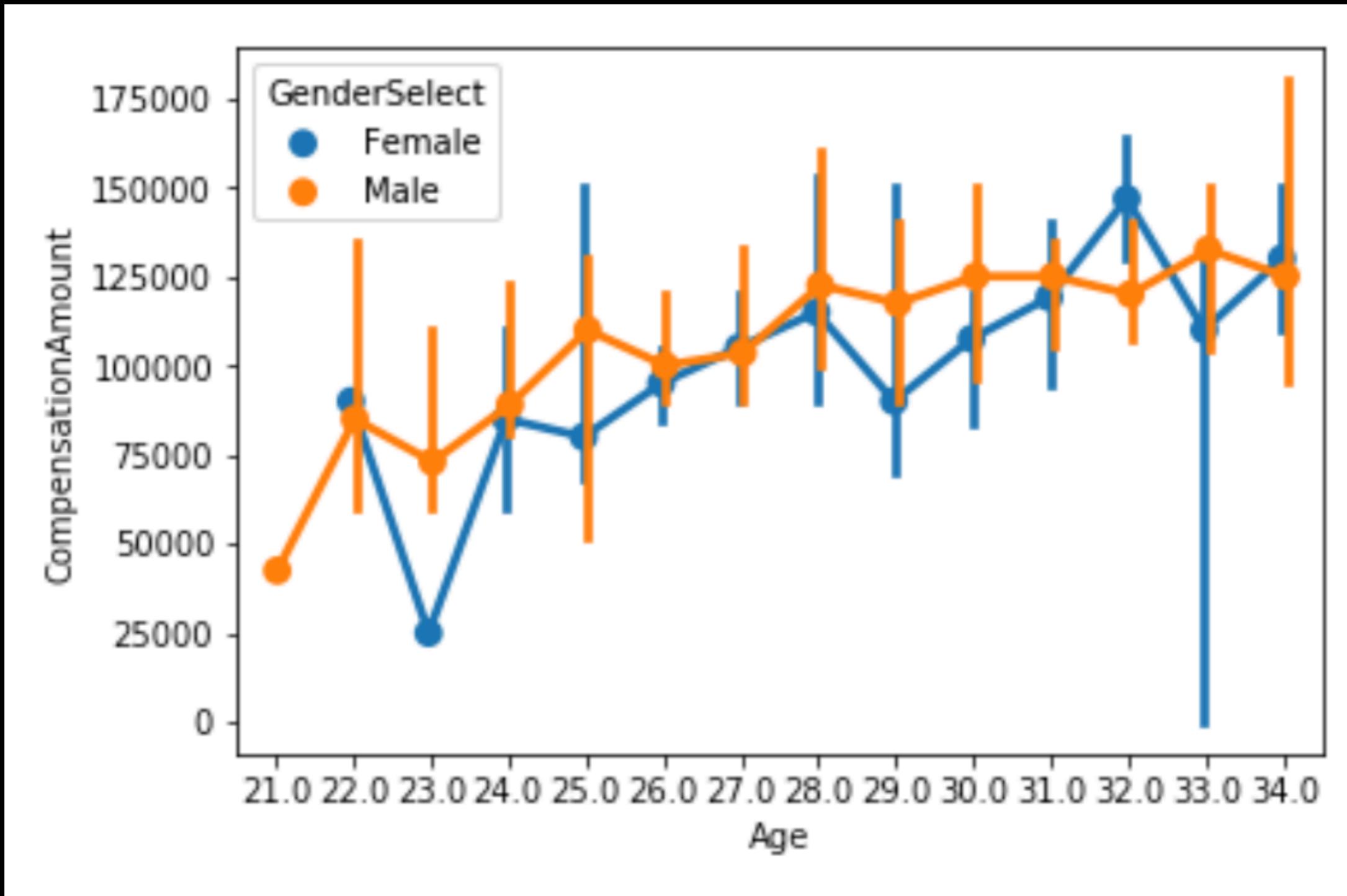


```
In [109]: import numpy as np
```

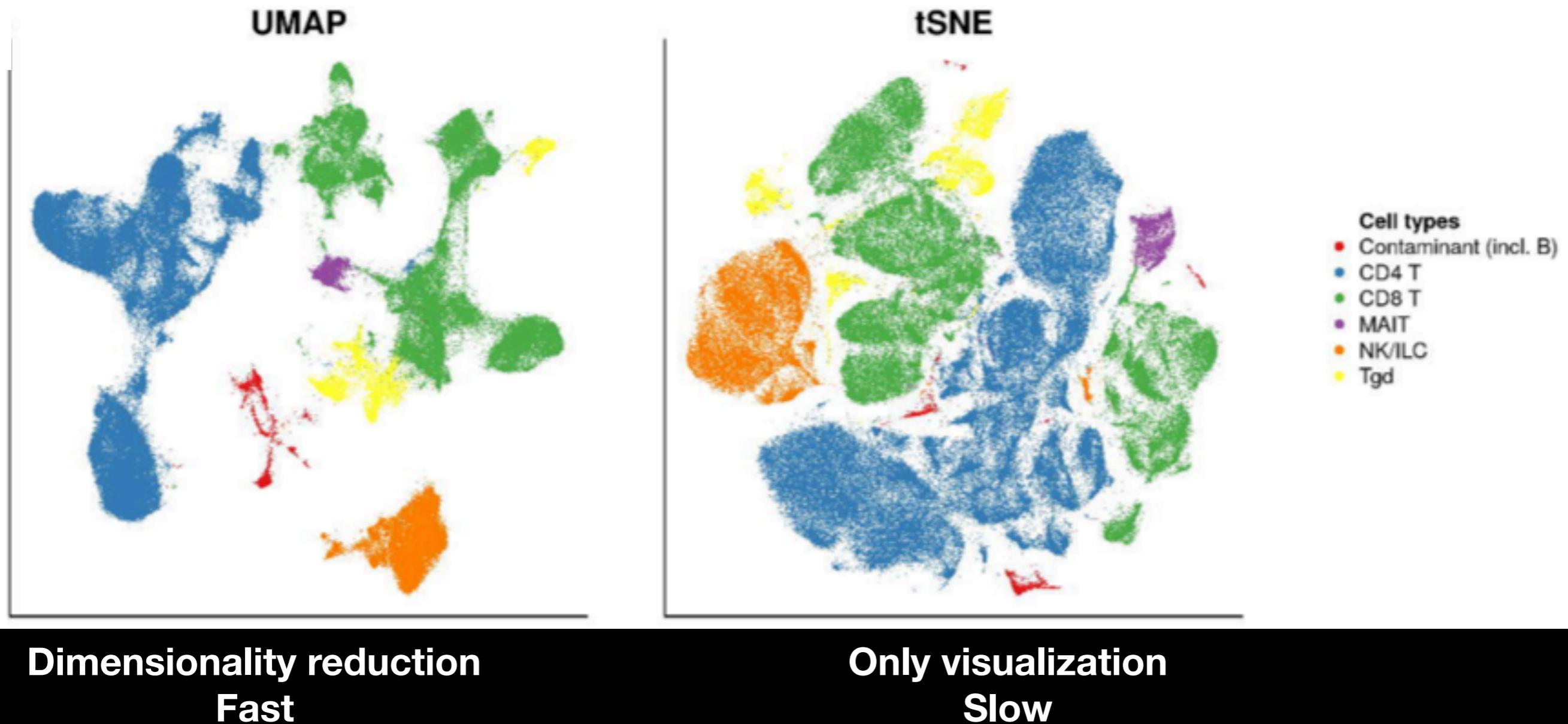


```
In [110]: sns.pointplot(
    x='Age',
    y='CompensationAmount',
    data=df_multiple_choice[f],
    estimator=np.median,
    hue='GenderSelect',
    dodge=True
)
```

Median salary of data scientists at university ages.



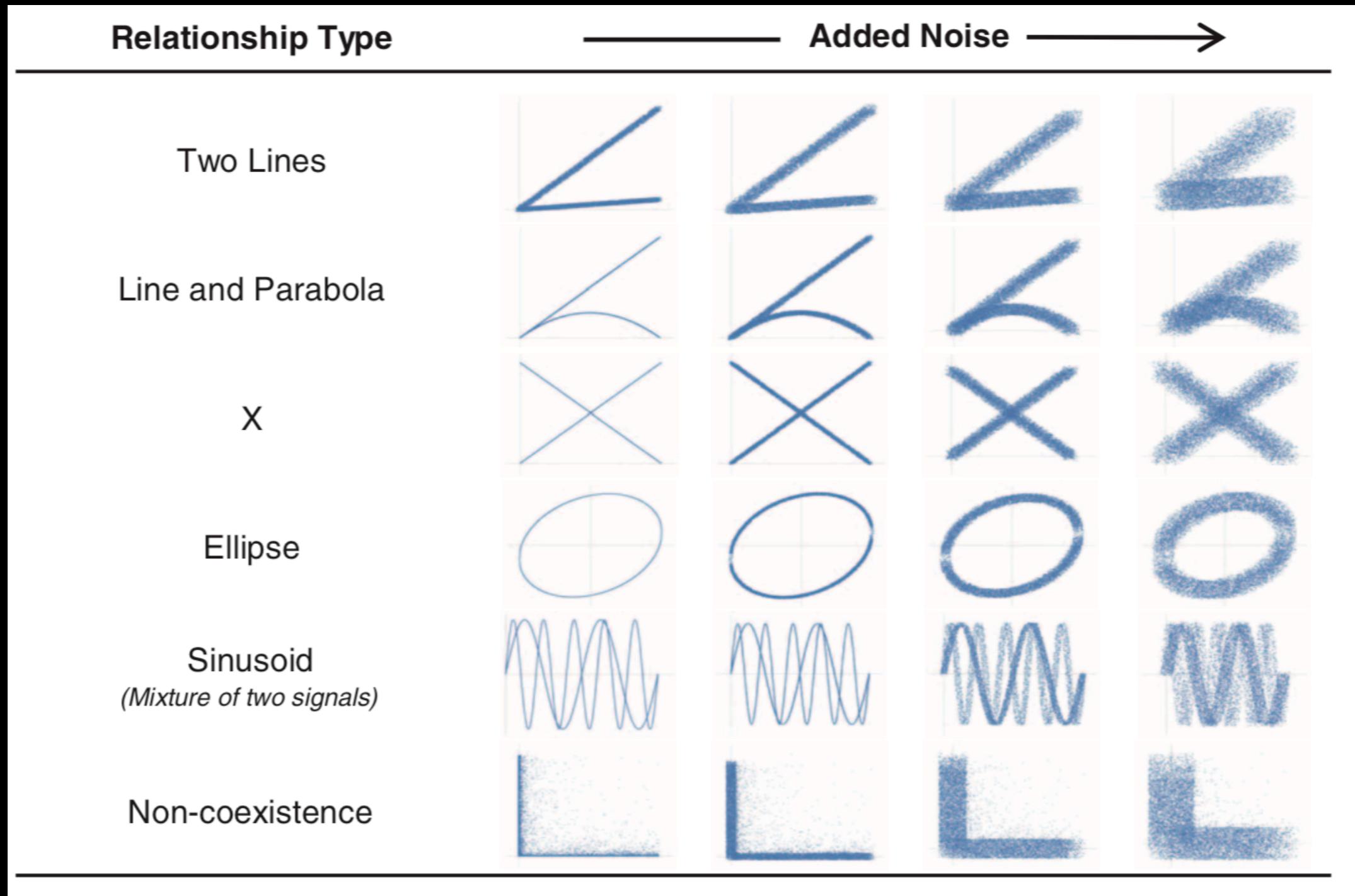
Further default tactic: visualize upon dimensionality reduction.



Becht et al., bioRxiv 2018

Distinct types of correlations can embody distinct assumptions on the data.

Distinct types of correlations can embody distinct assumptions on the data.



Reshef et al. 2011

Distinct types of correlations can embody distinct assumptions on the data.

Relationship Type	MIC	Pearson	Spearman	Mutual Information (KDE)	Mutual Information (Kraskov)	CorGC (Principal Curve-Based)	Maximal Correlation
Random	0.18	-0.02	-0.02	0.01	0.03	0.19	0.01
Linear	1.00	1.00	1.00	5.03	3.89	1.00	1.00
Cubic	1.00	0.61	0.69	3.09	3.12	0.98	1.00
Exponential	1.00	0.70	1.00	2.09	3.62	0.94	1.00
Sinusoidal <i>(Fourier frequency)</i>	1.00	-0.09	-0.09	0.01	-0.11	0.36	0.64
Categorical	1.00	0.53	0.49	2.22	1.65	1.00	1.00
Periodic/Linear	1.00	0.33	0.31	0.69	0.45	0.49	0.91
Parabolic	1.00	-0.01	-0.01	3.33	3.15	1.00	1.00
Sinusoidal <i>(non-Fourier frequency)</i>	1.00	0.00	0.00	0.01	0.20	0.40	0.80
Sinusoidal <i>(varying frequency)</i>	1.00	-0.11	-0.11	0.02	0.06	0.38	0.76

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

Main tactics for finding connections.

- Exploratory visualization
- Correlations
- Machine learning

```
In [113]: df_multiple_choice[ 'MLMethodNextYearSelect' ].value_counts()
```

```
In [113]: df_multiple_choice['MLMethodNextYearSelect'].value_counts()
```

```
Out[113]: Deep learning                                535
          Neural Nets                                 207
          Bayesian Methods                            101
          Time Series Analysis                      90
          Genetic & Evolutionary Algorithms        66
          Anomaly Detection                           61
          Text Mining                               52
          Social Network Analysis                  48
          Other                                     46
          Ensemble Methods (e.g. boosting, bagging) 42
          I don't plan on learning a new ML/DS method 33
          Cluster Analysis                           31
          Monte Carlo Methods                      28
          Support Vector Machines (SVM)            26
          Proprietary Algorithms                   21
          Random Forests                            18
          Survival Analysis                         18
          Regression                                15
          Rule Induction                            8
          Link Analysis                             8
          Decision Trees                            8
          Uplift Modeling                           6
          MARS                                      5
          Factor Analysis                           5
          Association Rules                        3
Name: MLMethodNextYearSelect, dtype: int64
```

Reminder

```
In [113]: df_multiple_choice['MLMethodNextYearSelect'].value_counts()
```

Out[113]:	Deep learning	535
	Neural Nets	207
	Bayesian Methods	101
	Time Series Analysis	90
	Genetic & Evolutionary Algorithms	66
	Anomaly Detection	61
	Text Mining	52
	Social Network Analysis	48
	Other	46
	Ensemble Methods (e.g. boosting, bagging)	42
	I don't plan on learning a new ML/DS method	33
	Cluster Analysis	31
	Monte Carlo Methods	28
	Support Vector Machines (SVM)	26
	Proprietary Algorithms	21
	Random Forests	18
	Survival Analysis	18
	Regression	15
	Rule Induction	8
	Link Analysis	8
	Decision Trees	8
	Uplift Modeling	6
	MARS	5
	Factor Analysis	5
	Association Rules	3
	Name: MLMethodNextYearSelect, dtype: int64	

Any question?

Note: Some discipline-specific approaches have very interesting properties.

TIME SERIES ANALYSIS

Information leverage in interconnected ecosystems: Overcoming the curse of dimensionality

Hao Ye and George Sugihara*

Turn high dimensionality from a curse to a blessing.

In ecological analysis, complexity has been regarded as an obstacle to overcome. Here we present a straightforward approach for addressing complexity in dynamic interconnected systems. We show that complexity, in the form of multiple interacting components, can actually be an asset for studying natural systems from temporal data. The central idea is that multidimensional time series enable system dynamics to be reconstructed from multiple viewpoints, and these viewpoints can be combined into a single model. We show how our approach, multiview embedding (MVE), can improve forecasts for simulated ecosystems and a mesocosm experiment. By leveraging complexity, MVE is particularly effective for overcoming the limitations of short and noisy time series and should be highly relevant for many areas of science.

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

Workflow

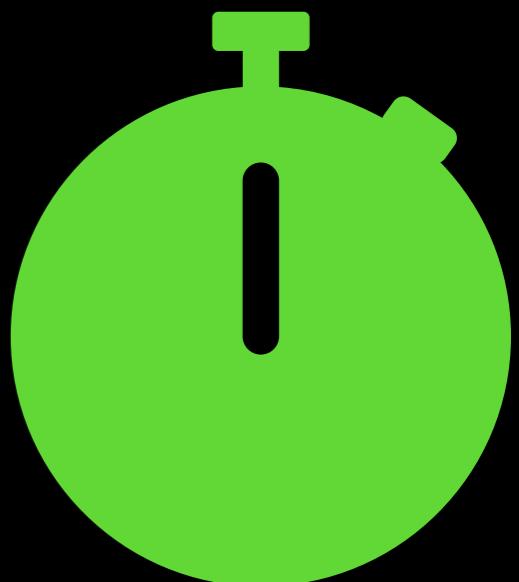
- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

break

**How does your setup (computer,
code, libraries, data formats, ...)
reduce the efficiency of your work?**



Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- **Environment**
- Organizing data
- Coding

break

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)

Constant Learning

- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

Constant Learning

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)
- Meetups
- www.data-science-nights.org

Workstation

- much RAM
- Fast-single core
- Don't use workstation for computations!

Space

- Avoid isolated rooms.

Constant Learning

- Books, scientific literature; e.g.: note that NU has license for huge ebook collection on data science through **Safari Online** (<https://www-safaribooksonline-com.turing.library.northwestern.edu/library/view/temporary-access/>)
- Meetups
- www.data-science-nights.org

Please object. As some of this is scientific field specific.

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

break

Setup

- Computer
- **Organizing data**
- Coding



**Data integration
is an entire
scientific field.**

**PRINCIPLES OF
DATA INTEGRATION**

ANHAI DOAN ALON HALEVY ZACHARY IVES

e.g.: dealing with
ambiguity in matching
of entries, or speeding
processing up,...

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA



OPEN ACCESS

Check for updates

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

Be consistent.



OPEN ACCESS

Check for updates

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

Be consistent.

Just put one thing in a cell.



OPEN ACCESS

Check for updates

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

Be consistent.

Just put one thing in a cell.

Create a data dictionary.

Data Organization in Spreadsheets

Karl W. Broman^a and Kara H. Woo^b

^aDepartment of Biostatistics & Medical Informatics, University of Wisconsin-Madison, Madison, WI; ^bInformation School, University of Washington, Seattle, WA

Be consistent.

Just put one thing in a cell.

Create a data dictionary.

Do not use color as data.



Journal of Statistical Software

August 2014, Volume 59, Issue 10.

<http://www.jstats.org>

**popular in some
academic fields**

Tidy Data

Hadley Wickham
RStudio

Abstract

A huge amount of effort is spent cleaning data to get it ready for analysis, but there has been little research on how to make data cleaning as easy and effective as possible. This paper tackles a small, but important, component of data cleaning: data tidying. Tidy datasets are easy to manipulate, model and visualize, and have a specific structure: each variable is a column, each observation is a row, and each type of observational unit is a table. This framework makes it easy to tidy messy datasets because only a small set of tools are needed to deal with a wide range of un-tidy datasets. This structure also makes it easier to develop tidy tools for data analysis, tools that both input and output tidy datasets. The advantages of a consistent data structure and matching tools are demonstrated with a case study free from mundane data manipulation chores.

Portability and access.

Central location

Database

Folder with datasets

Domain-specific (e.g.: NIH's Synapse)

Local copies

Single reference!

Portability and access.

Central location
Database

Folder with datasets
Domain-specific (e.g.: NIH's Synapse)

Local copies
Single reference!

Portability and access.

Central location

Database

Folder with datasets

Domain-specific (e.g.: NIH's Synapse)

Local copies

Single reference!

Portability and access.

Central location
Database
Folder with datasets
Domain-specific (e.g.: NIH's Synapse)

Local copies
Single reference!

Portability and access.

Central location

Database

Folder with datasets

Domain-specific (e.g.: NIH's Synapse)

Local copies

Single reference!

many technical choices

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_clipboard
In [ ]: df_multiple_choice.to_csv
In [ ]: df_multiple_choice.to_dense
In [ ]: df_multiple_choice.to_dict
In [ ]: df_multiple_choice.to_excel
In [ ]: df_multiple_choice.to_feather
In [ ]: df_multiple_choice.to_gbq
In [ ]: df_multiple_choice.to_hdf
In [ ]: df_multiple_choice.to_html
In [ ]: df_multiple_choice.to_json
```

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_latex
In [ ]: df_multiple_choice.to_msgpack
In [ ]: df_multiple_choice.to_panel
In [ ]: df_multiple_choice.to_parquet
In [ ]: df_multiple_choice.to_period
In [ ]: df_multiple_choice.to_pickle
In [ ]: df_multiple_choice.to_records
In [ ]: df_multiple_choice.to_sparse
In [ ]: df_multiple_choice.to_sql
In [ ]: df_multiple_choice.to_stata
In [ ]:
```

many technical choices

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_clipboard
In [ ]: df_multiple_choice.to_csv
In [ ]: df_multiple_choice.to_dense
In [ ]: df_multiple_choice.to_dict
In [ ]: df_multiple_choice.to_excel
In [ ]: df_multiple_choice.to_feather
In [ ]: df_multiple_choice.to_gbq
In [ ]: df_multiple_choice.to_hdf
In [ ]: df_multiple_choice.to_html
In [ ]: df_multiple_choice.to_json
```

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_latex
In [ ]: df_multiple_choice.to_msgpack
In [ ]: df_multiple_choice.to_panel
In [ ]: df_multiple_choice.to_parquet
In [ ]: df_multiple_choice.to_period
In [ ]: df_multiple_choice.to_pickle
In [ ]: df_multiple_choice.to_records
In [ ]: df_multiple_choice.to_sparse
In [ ]: df_multiple_choice.to_sql
In [ ]: df_multiple_choice.to_stata
In [ ]:
```



CSV
Spreadsheets

many technical choices

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_clipboard
In [ ]: df_multiple_choice.to_csv
In [ ]: df_multiple_choice.to_dense
In [ ]: df_multiple_choice.to_dict
In [ ]: df_multiple_choice.to_excel
In [ ]: df_multiple_choice.to_feather
In [ ]: df_multiple_choice.to_gbq
In [ ]: df_multiple_choice.to_hdf
In [ ]: df_multiple_choice.to_html
In [ ]: df_multiple_choice.to_json
```

CSV
Spreadsheets

HDF5
Selectively load
Parts of data
(avoid breaking
cluster)

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_latex
In [ ]: df_multiple_choice.to_msgpack
In [ ]: df_multiple_choice.to_panel
In [ ]: df_multiple_choice.to_parquet
In [ ]: df_multiple_choice.to_period
In [ ]: df_multiple_choice.to_pickle
In [ ]: df_multiple_choice.to_records
In [ ]: df_multiple_choice.to_sparse
In [ ]: df_multiple_choice.to_sql
In [ ]: df_multiple_choice.to_stata
```

many technical choices

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_clipboard
In [ ]: df_multiple_choice.to_csv
In [ ]: df_multiple_choice.to_dense
In [ ]: df_multiple_choice.to_dict
In [ ]: df_multiple_choice.to_excel
In [ ]: df_multiple_choice.to_feather
In [ ]: df_multiple_choice.to_gbq
In [ ]: df_multiple_choice.to_hdf
In [ ]: df_multiple_choice.to_html
In [ ]: df_multiple_choice.to_json
```

CSV
Spreadsheets

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_latex
In [ ]: df_multiple_choice.to_msgpack
In [ ]: df_multiple_choice.to_panel
In [ ]: df_multiple_choice.to_parquet
In [ ]: df_multiple_choice.to_period
In [ ]: df_multiple_choice.to_pickle
In [ ]: df_multiple_choice.to_records
In [ ]: df_multiple_choice.to_sparse
In [ ]: df_multiple_choice.to_sql
In [ ]: df_multiple_choice.to_stata
```

HDF5
Selectively load
Parts of data
(avoid breaking
cluster)

PARQUET
Very fast reading
and writing of
entire tables

many technical choices

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_clipboard
In [ ]: df_multiple_choice.to_csv
In [ ]: df_multiple_choice.to_dense
In [ ]: df_multiple_choice.to_dict
In [ ]: df_multiple_choice.to_excel
In [ ]: df_multiple_choice.to_feather
In [ ]: df_multiple_choice.to_gbq
In [ ]: df_multiple_choice.to_hdf
In [ ]: df_multiple_choice.to_html
In [ ]: df_multiple_choice.to_json
```

CSV
Spreadsheets

```
In [ ]: df_multiple_choice.to_
In [ ]: df_multiple_choice.to_latex
In [ ]: df_multiple_choice.to_msgpack
In [ ]: df_multiple_choice.to_panel
In [ ]: df_multiple_choice.to_parquet
In [ ]: df_multiple_choice.to_period
In [ ]: df_multiple_choice.to_pickle
In [ ]: df_multiple_choice.to_records
In [ ]: df_multiple_choice.to_sparse
In [ ]: df_multiple_choice.to_sql
In [ ]: df_multiple_choice.to_stata
```

HDF5
Selectively load
Parts of data
(avoid breaking
cluster)

PARQUET
Very fast reading
and writing of
entire tables
SQL
Elegant query
language

```
In [172]: example_data = df_multiple_choice[['Country', 'Age', 'GenderSelect']].dropna()
```

```
In [173]: example_data.head()
```

Out[173]:

	Country	Age	GenderSelect
1	United States	30.0	Female
2	Canada	28.0	Male
3	United States	56.0	Male
4	Taiwan	38.0	Male
5	Brazil	46.0	Male

HDF5

```
In [174]: example_data.to_hdf(  
    './example.h5',  
    format='table',  
    key='searchable_data',  
    mode='w',  
    data_columns=True,  
    )
```

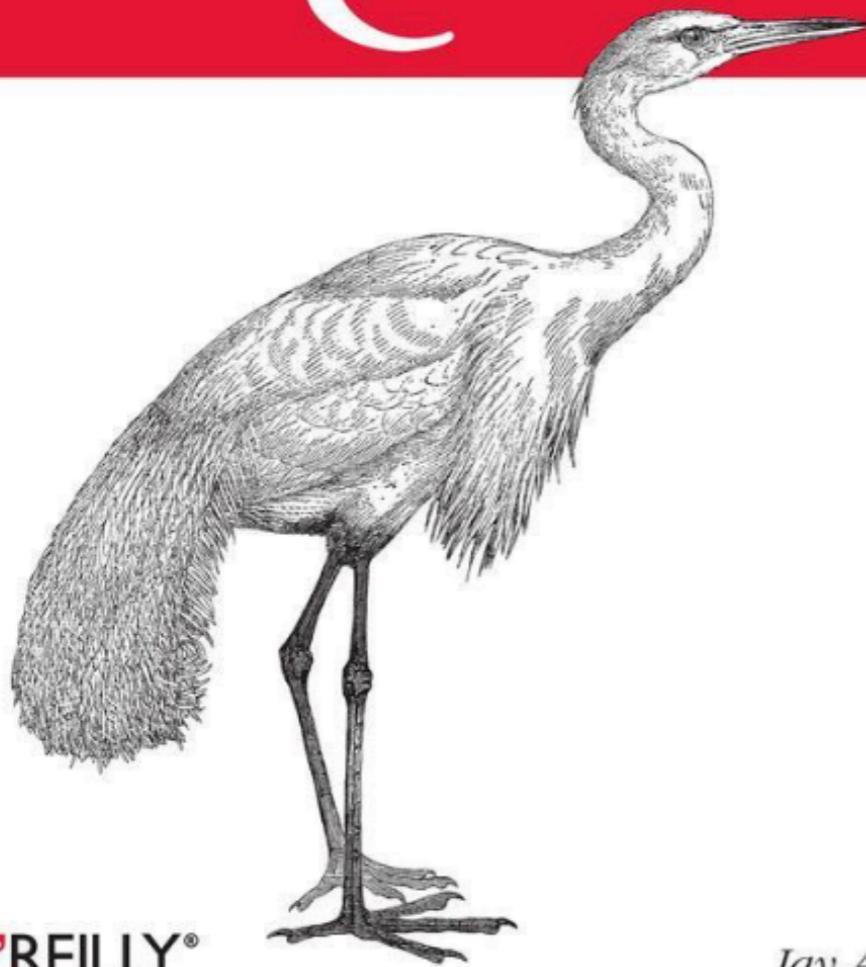
```
In [175]: pd.read_hdf('./example.h5', 'searchable_data', where=['Age=20'])
```

Out[175]:

	Country	Age	GenderSelect
10	Russia	20.0	Female
123	United States	20.0	Female
125	Turkey	20.0	Male
144	United States	20.0	Male
148	India	20.0	Male
172	China	20.0	Male

Using

SQLite



O'REILLY®

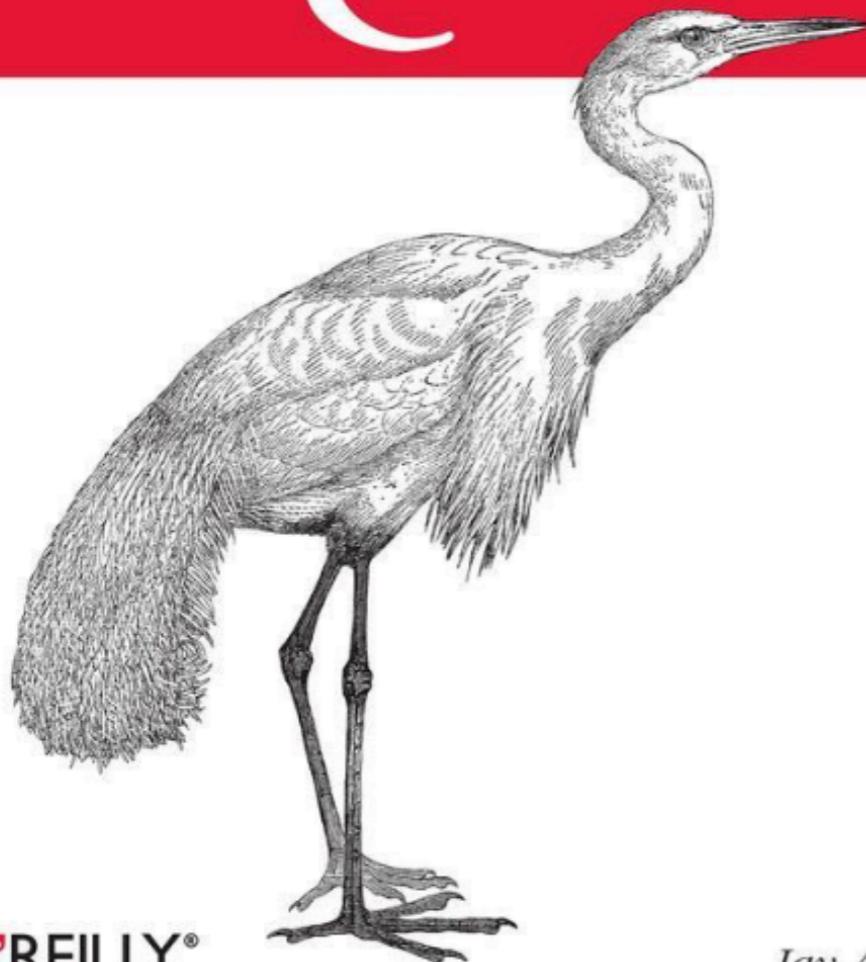
Jay A. Kreibich

Remember

You can have a server-less database.

Using

SQLite



O'REILLY®

Jay A. Kreibich

Remember

You can have a server-less database.

NU provides tons of free ebooks.

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- **Coding**

break

Save time!

Important to many.

Reduce time coding.

Specialist applications.

Reduce time running
code.

Save time!

Important to many.

Reduce time coding.

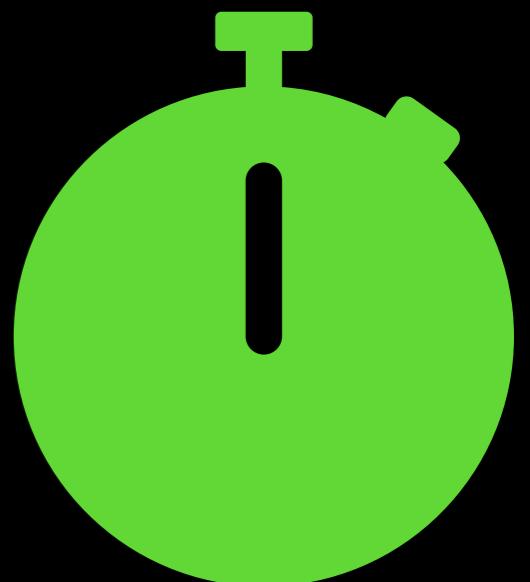
Important to even more

Reduce time others need to understand code.

Specialist applications.

Reduce time running
code.

**Which coding practices
could save time?**



Watch your language!

ISSN: 2158-2229

PLOS ONE

PUBLISH ABOUT BROWSE

OPEN ACCESS PEER-REVIEWED

RESEARCH ARTICLE

A large-scale analysis of bioinformatics code on GitHub

Pamela H. Russell , Rachel L. Johnson, Shreyas Ananthan, Benjamin Harnke, Nichole E. Carlson

Published: October 31, 2018 • <https://doi.org/10.1371/journal.pone.0205898>

Article	Authors	Metrics	Comments	Media Coverage
				

Abstract

Introduction
Results
Discussion
Methods
Supporting information
Acknowledgments
References

Reader Comments (0)

Abstract

In recent years, the explosion of genomic data and bioinformatic tools has been accompanied by a growing conversation around reproducibility of results and usability of software. However, the actual state of the body of bioinformatics software remains largely unknown. The purpose of this paper is to investigate the state of source code in the bioinformatics community, specifically looking at relationships between code properties, development activity, developer communities, and software impact. To investigate these issues, we curated a list of 1,720 bioinformatics repositories on GitHub through their mention in peer-reviewed bioinformatics articles. Additionally, we included 23 high-profile repositories identified by their popularity in an online bioinformatics forum. We analyzed repository metadata, source code, development activity, and team dynamics using data made available publicly through the GitHub API, as well as article metadata. We found key relationships within our dataset, including: certain scientific topics are associated with more active code development and higher community interest in the repository; most of the code in the main dataset is written in dynamically typed languages, while most of

GitHub Education

Students Teachers Schools Events

[Get benefits](#)

Real-world tools, engaged students

GitHub Education helps students, teachers, and schools access the tools and events they need to shape the next generation of software development.



[GitHub Student Developer Pack](#)

The best developer tools, free for students



[GitHub Campus Experts](#)

Training to enrich the technology community at your school



[GitHub Campus Program](#)

GitHub for your whole school, with everything you need to make it great



[GitHub Classroom](#)

The GitHub workflow, scaled for the needs of students



[GitHub Campus Advisors](#)

Teacher training to master Git and GitHub



Track and share your code!
e.g.: GitHub, which offers
free repositories for academics
and research groups
education.github.com

The art of naming variables

July 30th 2018

There are only two hard things in Computer Science: cache invalidation and naming things..

The art of naming variables

July 30th 2018

There are only two hard things in Computer Science: cache invalidation and naming things..

```
['apple', 'banana', 'cucumber']
```

The art of naming variables

July 30th 2018

There are only two hard things in Computer Science: cache invalidation and naming things..

`['apple', 'banana', 'cucumber']`

`fruits`

The art of naming variables

July 30th 2018

There are only two hard things in Computer Science: cache invalidation and naming things..

```
['apple', 'banana', 'cucumber']
```

fruits

```
function any(x.isin(fruits))
```

The art of naming variables

July 30th 2018

There are only two hard things in Computer Science: cache invalidation and naming things..

```
['apple', 'banana', 'cucumber']
```

fruits

```
function any(x.isin(fruits))
```

is_fruit

Principles

[edit]

Principles are listed as follows:

Beautiful is better than ugly.

Explicit is better than implicit.

Simple is better than complex.

Complex is better than complicated.

Flat is better than nested.

Sparse is better than dense.

Readability counts.

Special cases aren't special enough to break the rules.

Although practicality beats purity.

Errors should never pass silently.

Unless explicitly silenced.

In the face of ambiguity, refuse the temptation to guess.

There should be one—and preferably only one—obvious way to do it.

Although that way may not be obvious at first unless you're Dutch.

Now is better than never.

Although never is often better than *right* now.^[n 1]

If the implementation is hard to explain, it's a bad idea.

If the implementation is easy to explain, it may be a good idea.

Namespaces are one honking great idea—let's do more of those!

style conventions
e.g. Zen of Python
(and subsequently
PEP8)

Did we miss some useful?

the last lesson: Plan for **changing data.**

Saves time!
allows expansion of project
allows to only spend time on improving access when needed

the last lesson: Plan for **changing data.**

Saves time!
allows expansion of project
allows to only spend time on improving access when needed

the last lesson: Plan for **changing data.**

Saves time!
allows expansion of project
allows to only spend time on improving access when needed

the last lesson: Plan for **changing data.**

Saves time!
allows expansion of project
allows to only spend time on improving access when needed

Don't access data directly!

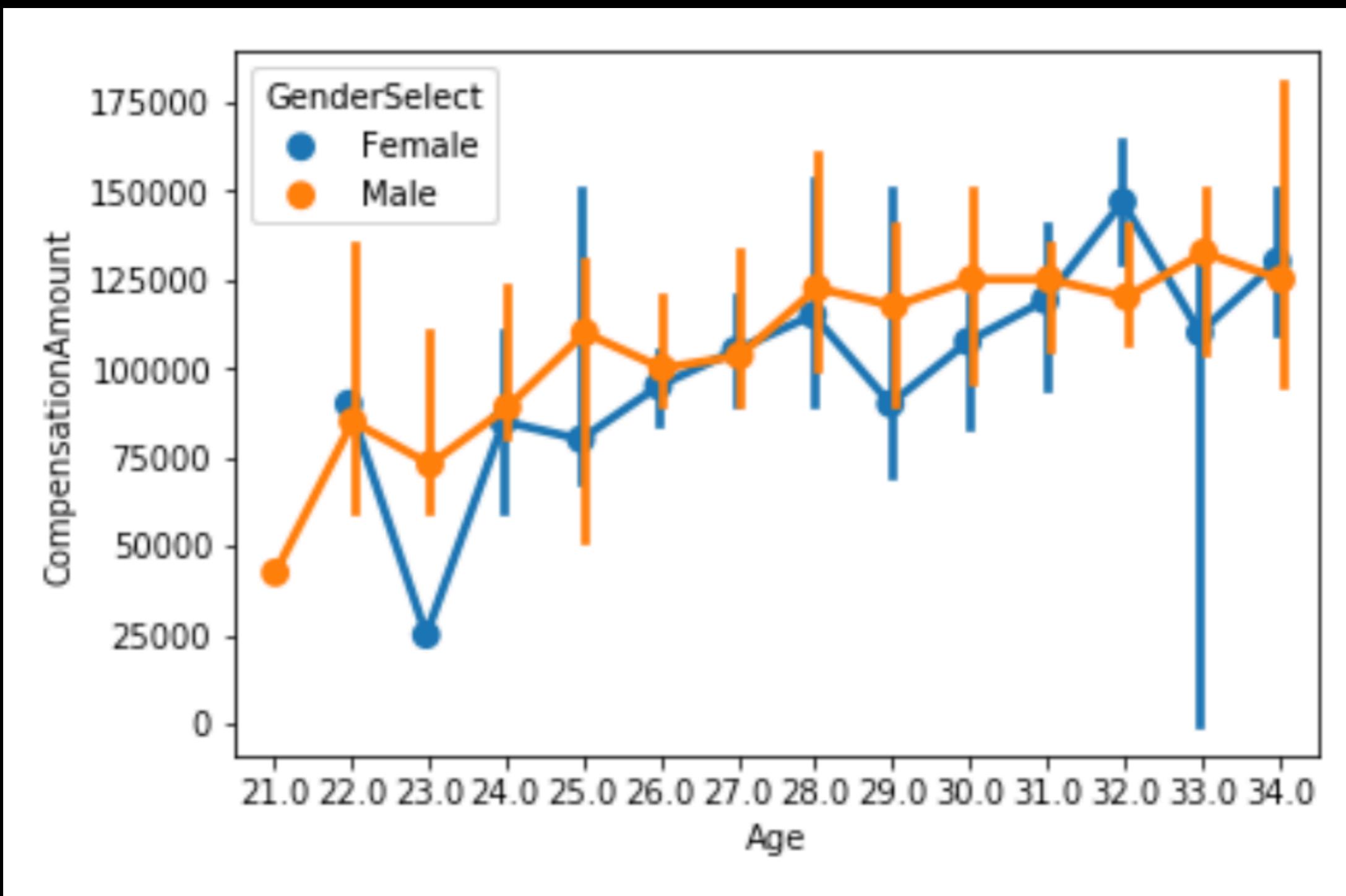
Don't access data directly!

Notebook → **Adaptor Function** → Data resource

Don't access data directly!



**Can change data resource
without worrying in ongoing research.**



Direct loading is acceptable for previewing new datasets.

```
In [1]: %matplotlib inline
```

```
In [2]: cd '/Users/tstoeger/Dropbox/Work/kaggle_survey'  
/Users/tstoeger/Dropbox/Work/kaggle_survey
```

```
In [3]: import pandas as pd  
import seaborn as sns
```

```
In [4]: df_multiple_choice = pd.read_csv(  
    './multipleChoiceResponses.csv',  
    encoding='latin-1', # files have some special characters  
    low_memory=False) # some columns have mixed type
```

```
In [5]: df_multiple_choice.shape # get idea of number of datasets  
Out[5]: (16716, 228)
```

... but move your loading into adaptor functions in the long run.

```
In [190]: def load_country_gender_at_given_age(age_of_interest):
    """
    Loads the country and gender of survey takes, which
    are of a given age

    Input:
        age_of_interest    int

    Output:
        dataframe

    """

    df = pd.read_csv(
        './multipleChoiceResponses.csv',
        encoding='latin-1', # files have some special characters
        low_memory=False)   # some columns have mixed type
    df = df[[ 'Country', 'Age', 'GenderSelect']].dropna()
    df = df[df[ 'Age']==age_of_interest]

    return df
```

Define function
In human digestible
Manner

... but move your loading into adaptor functions in the long run.

```
In [190]: def load_country_gender_at_given_age(age_of_interest):
    """
    Loads the country and gender of survey takes, which
    are of a given age

    Input:
        age_of_interest    int

    Output:
        dataframe

    """

    df = pd.read_csv(
        './multipleChoiceResponses.csv',
        encoding='latin-1', # files have some special characters
        low_memory=False)  # some columns have mixed type
    df = df[['Country', 'Age', 'GenderSelect']].dropna()
    df = df[df['Age']==age_of_interest]

    return df
```

```
In [191]: d = load_country_gender_at_given_age(20)
```

```
In [192]: d.head()
```

Out[192]:

	Country	Age	GenderSelect
10	Russia	20.0	Female
123	United States	20.0	Female

```
In [193]: def load_country_gender_at_given_age(age_of_interest):
    """
    Loads the country and gender of survey takes, which
    are of a given age

    Input:
        age_of_interest    int

    Output:
        dataframe

    """
    df = pd.read_hdf('./example.h5', 'searchable_data', where=['Age={}'.format(
        int(age_of_interest))])

    return df
```

Changed to HDF5



```
In [193]: def load_country_gender_at_given_age(age_of_interest):
    """
    Loads the country and gender of survey takes, which
    are of a given age

    Input:
        age_of_interest    int

    Output:
        dataframe

    """
    df = pd.read_hdf('./example.h5', 'searchable_data', where=['Age={}'.format(
        int(age_of_interest))])

    return df
```

```
In [194]: d = load_country_gender_at_given_age(20)
```

```
In [195]: d.head()
```

```
Out[195]:
```

	Country	Age	GenderSelect
10	Russia	20.0	Female
123	United States	20.0	Female
125	Turkey	20.0	Male
144	United States	20.0	Male
148	India	20.0	Male

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding

2 minutes:
Which encountered
strategy could
facilitate your work?

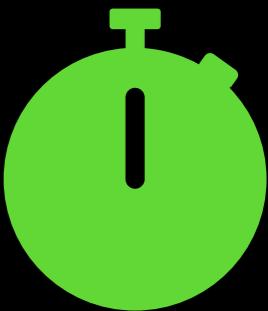


Workflow

- Anticipate problems
- Get overview
- Clean data
- Find connections

Setup

- Computer
- Organizing data
- Coding



**2 minutes:
Which encountered
strategy could
facilitate your work?**

Questions / Comments?

thank you