# "WHAT IS DATA SCIENCE?" REVISITED

Jeff Goldsmith, PhD

Department of Biostatistics

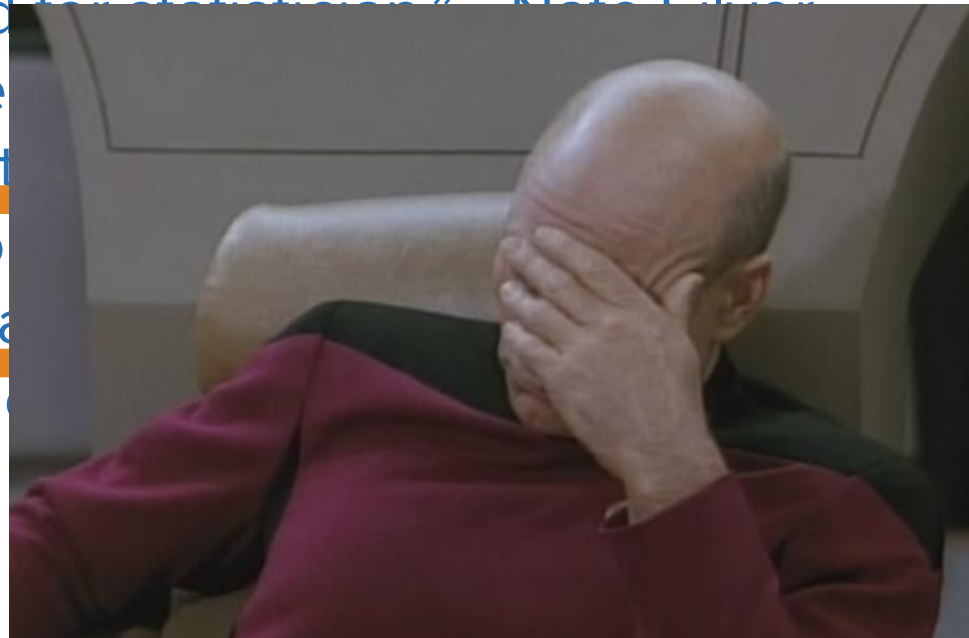# Some <sub>not great</sub> definitions

- Data science = statistics
- Data science = computer science
- Data science = machine learning
- Data science = statistics + computer science + machine learning
- Data scientists are big data wranglers
- "A data scientist is just a sexier word for statistician." –Nate Silver
- "A data scientist is a better computer scientist than a statistician and is a better statistician than a computer scientist."
- "A data scientist is a statistician who is useful" – Hadley Wickham
- A data scientist is a good statistical analyst
- A data scientist is a statistician who codes in python

# Some <sub>not great</sub> definitions

- Data science = statistics
- Data science = computer science
- Data science = machine learning
- Data science = statistics + computer science + machine learning
- Data scientists are big data wranglers
- "A data scientist is just a sexier word for statistician" – Nate Silver
- "A data scientist is a better compute[r ...] a better statistician than a computer scientist
- "A data scientist is a statistician who [...]
- A data scientist is a good statistical a[...]
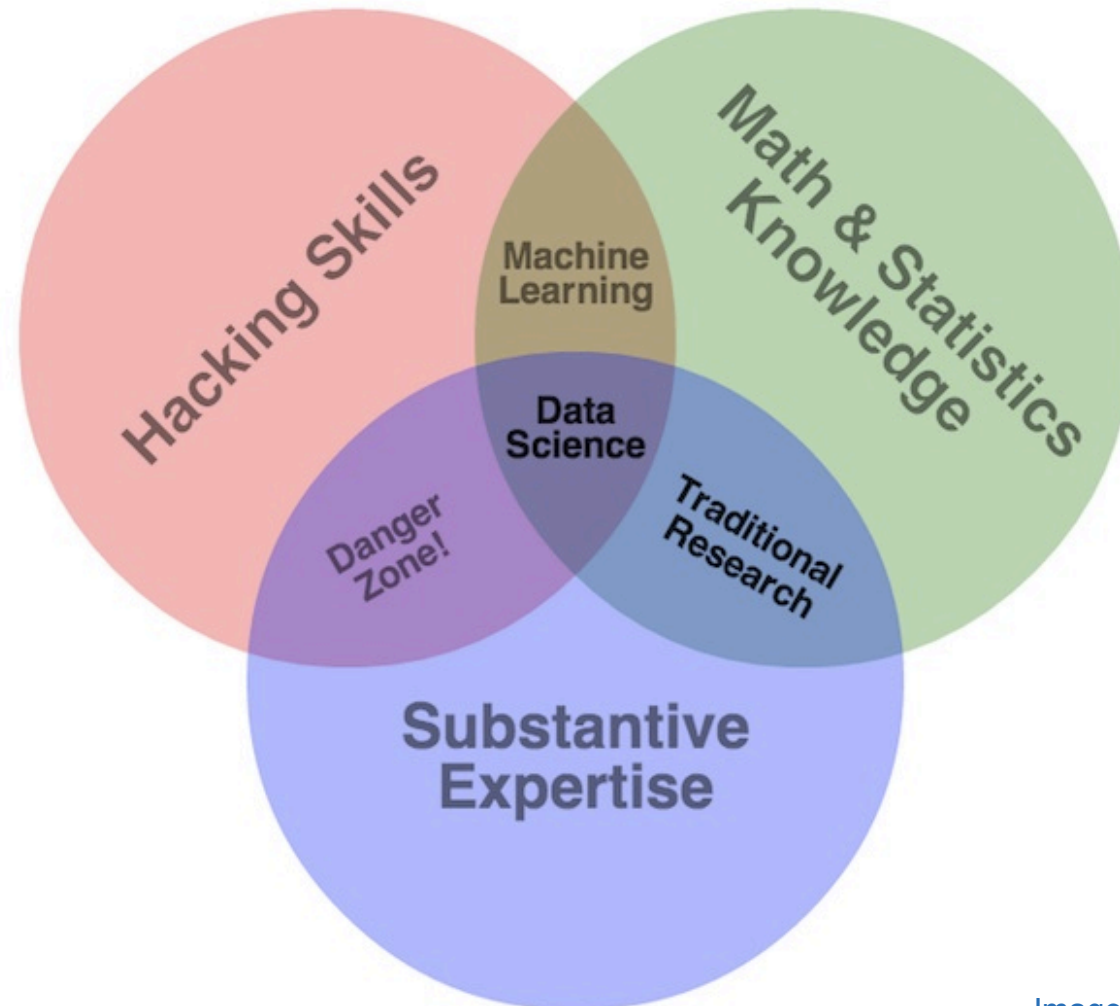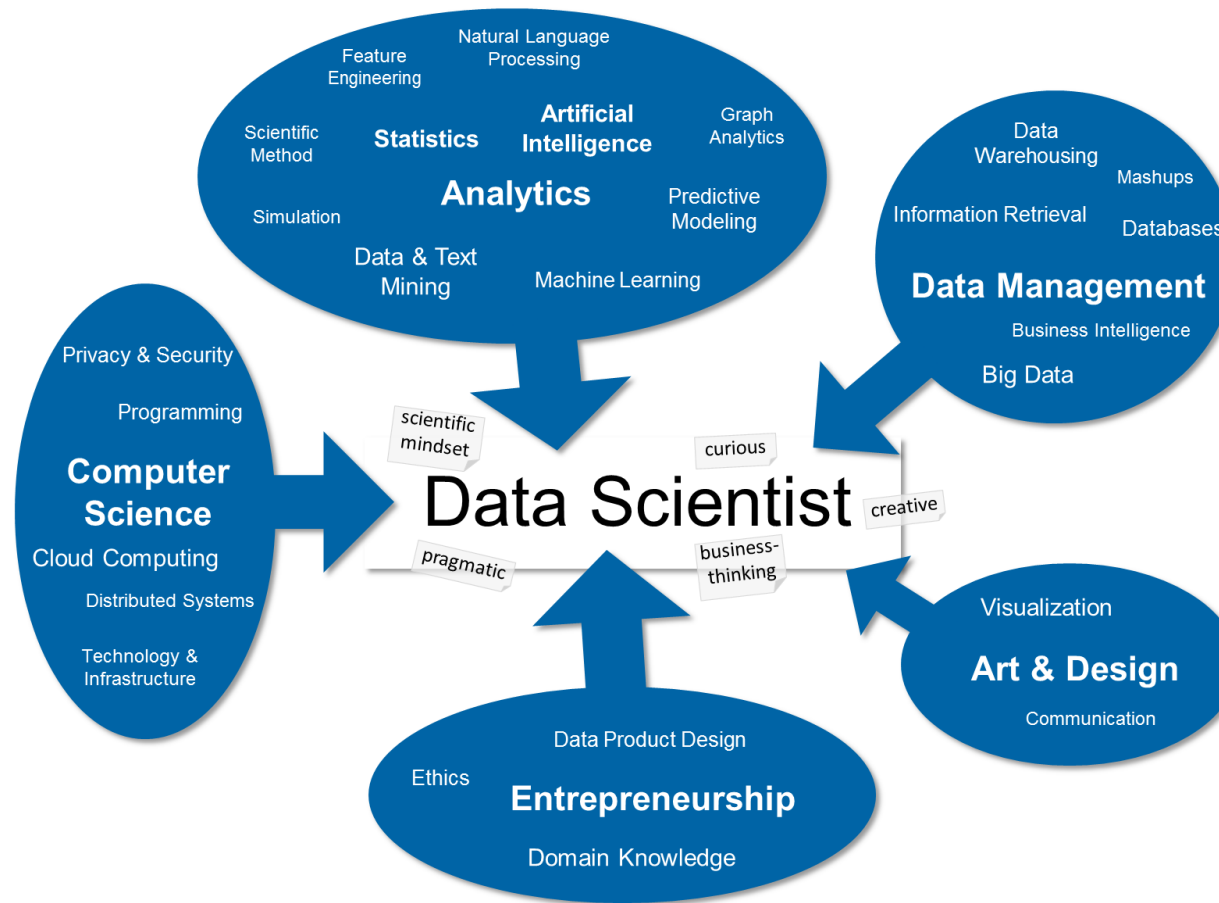- A data scientist is a statistician who [...]

# Maybe pictures will help?
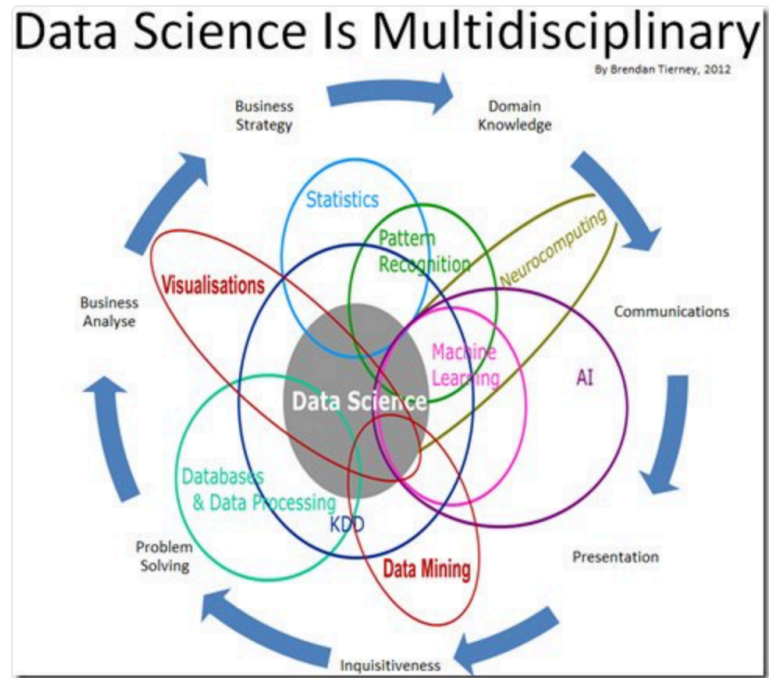


Image from Drew Conway

# Maybe pictures will help?

https://blog.zhaw.ch/datascience/the-data-science-skill-set/

# Maybe pictures will help?

From twitter

# Maybe pictures will help?

# Recurring themes

- You need "data skills"
  - Data wrangling
  - Reproducibility
  - Communication
  - Analytics and modeling
- These have been the focus of this course and others, and will continue to be the focus

- You also need a mindset
  - Intellectual curiosity
  - Ability to solve problems
  - Interest in domain, even empathy with collaborators

# Problem solving

"I've interviewed a lot of people over the years.... Recently, when people have an interview, I ask a single question that I think tries to get at the point of problem solving. The question I ask is along the lines of '[Imagine you had access to a database of 100 million mobile devices.] What questions would you ask? What types of things do you think you could learn, and how would you go about doing it?'"

From "How Industry Views Data Science Education in Statistics Departments", Chris Volinsky's JSM 2015 talk

# Practice problem solving

- You can (and should) practice having a mindset, or a style of thinking
  - Make a habit of asking yourself what you would like to do with a data resource
  - Think about how you would accomplish it

- Be on the lookout for cool projects, and learn from them
  - Pay attention to the thought process, not just the specific tools

- Many projects need overlapping skill sets
  - You don't have to be a domain expert yourself, but you may need to work with one
  - You'll also have to communicate effectively with that person, which means at least taking an interest