# Data Science Bootcamps with 

# Business Problem

- Imagine yourself as a marketer for a data science boot camp.

   What might you care about?


- How do you know what people are honestly interested in since reviews are public.

# Reddit

- Site that hosts any topic imaginable

- High utility since there are numerous subreddits discussing data science, boot camps, and data science boot camps.

- Users are anonymous so opinions shared are more likely to be honest (increase in the validity of comments).

- Good site for asking questions and getting answers. StackOverflow for everything else.

- https://www.reddit.com/r/datascience/search?q=data+science+bootcamp&restrict_sr=on

# The Tools of the Trade

- Scrapy
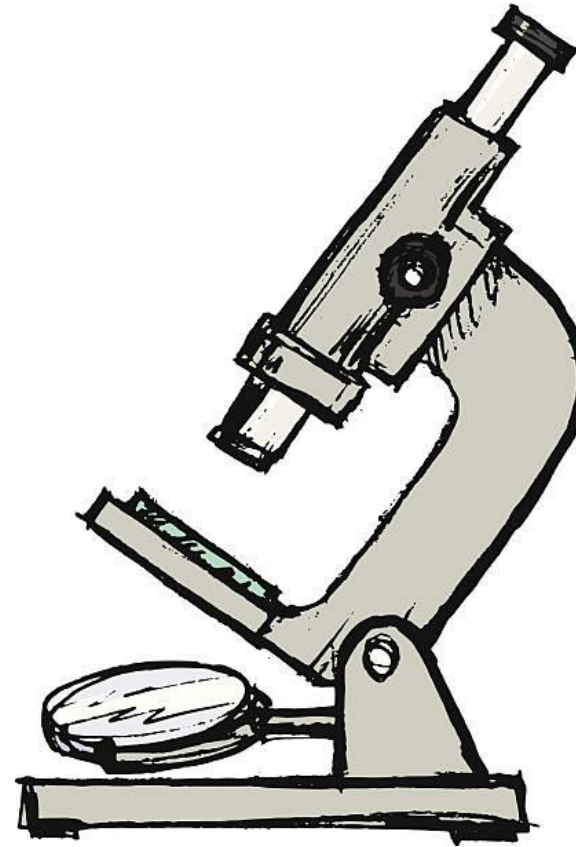
A general framework you can install on python.

- Latent Dirichlet Allocation(NLP)
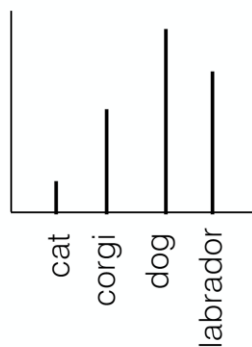- Unsupervised Learning Algorithm

- Probabilistic Model
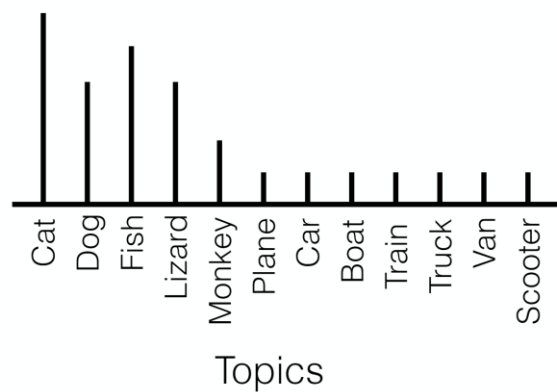- The algorithm discovers patterns in the text based upon variation.

# LDA continued

## Topics (Z)

- A topic is a probability distribution over words



cat · corgi · dog · labrador

## Topic Distribution for a Document



Topics: Cat, Dog, Fish, Lizard, Monkey, Plane, Car, Boat, Train, Truck, Van, Scooter

A document can be described by a recipe of topics and "how much" of each topic it contains

# LDA continued

## Logic of Process

**Document**

word word word
word word word
word word word
word word word
word word word
word word word
word word word

Topic 1
Topic 2
Topic 3

Basic Idea

- Documents are made up of words that belong (with some probability) to topics
- So…We can just reverse engineer these words to learn what a document is about
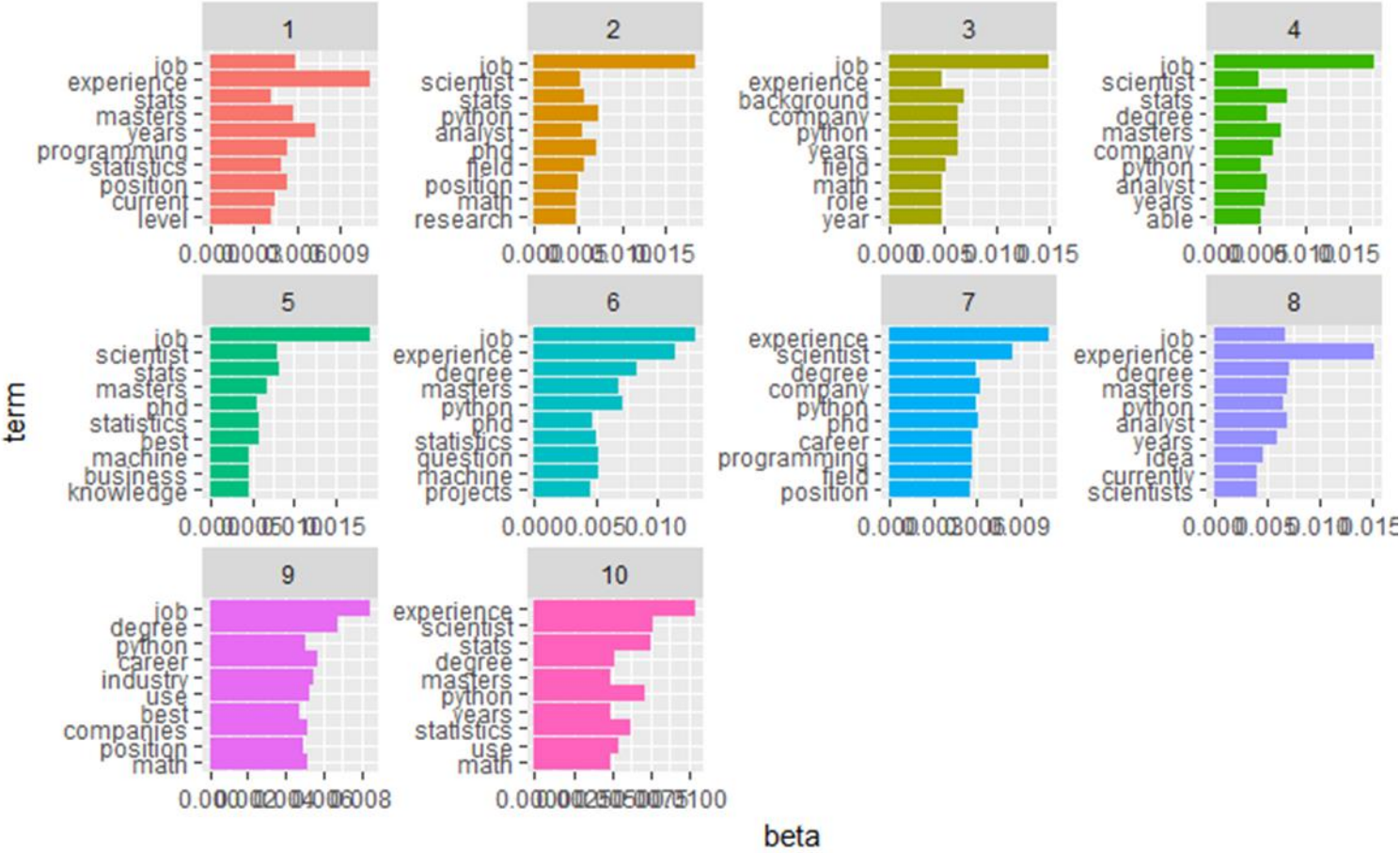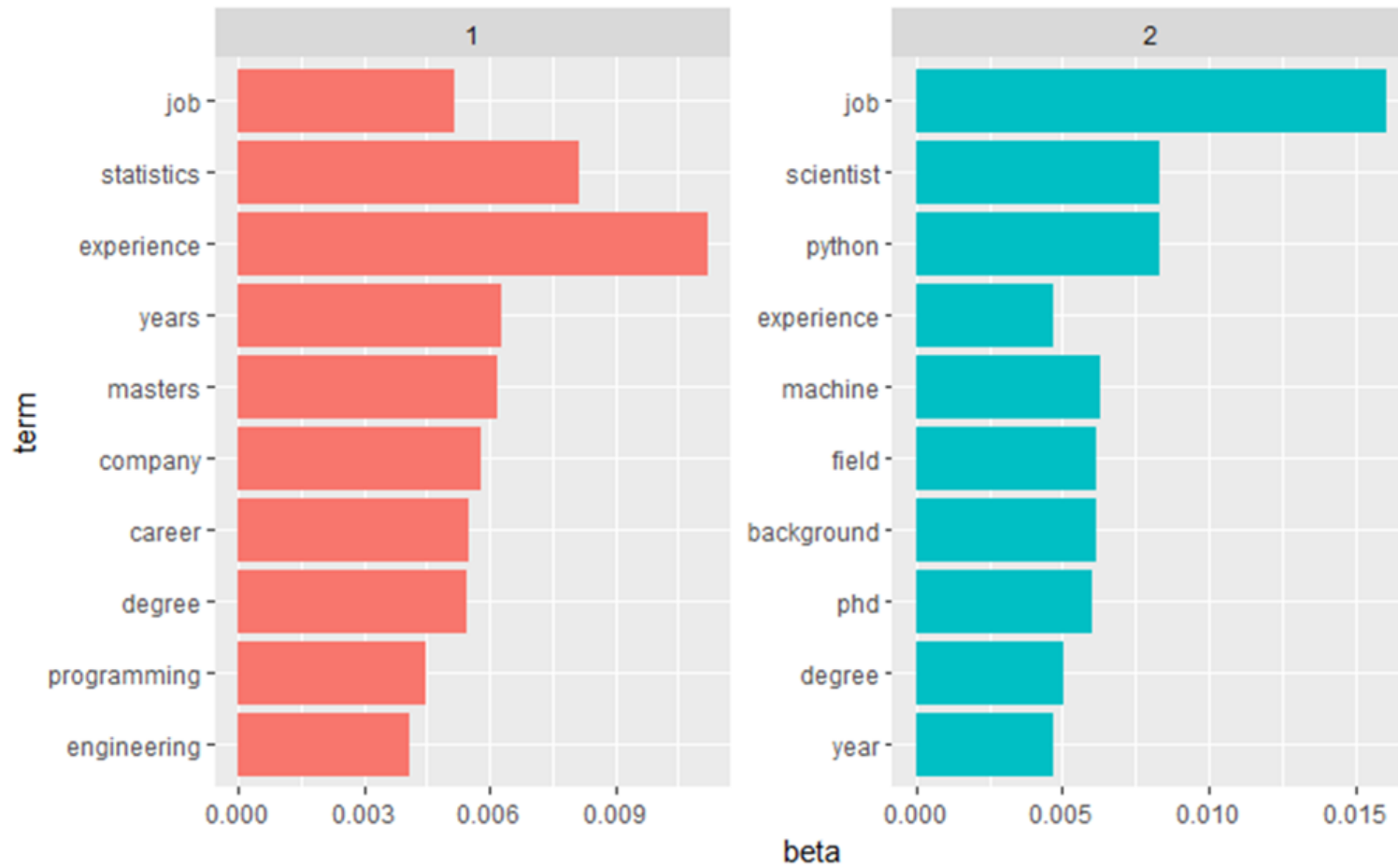
# Word Cloud

122,336 words scraped
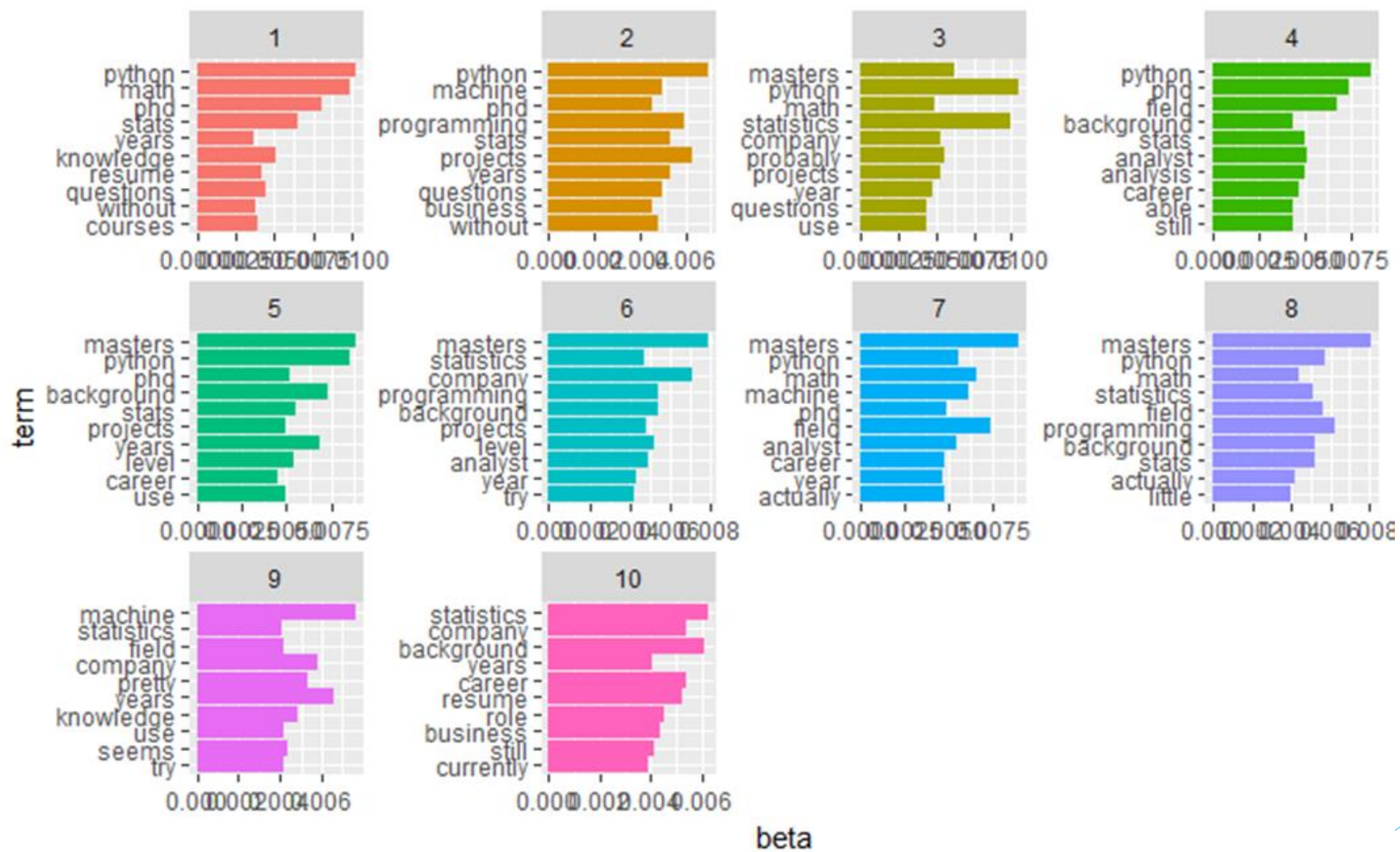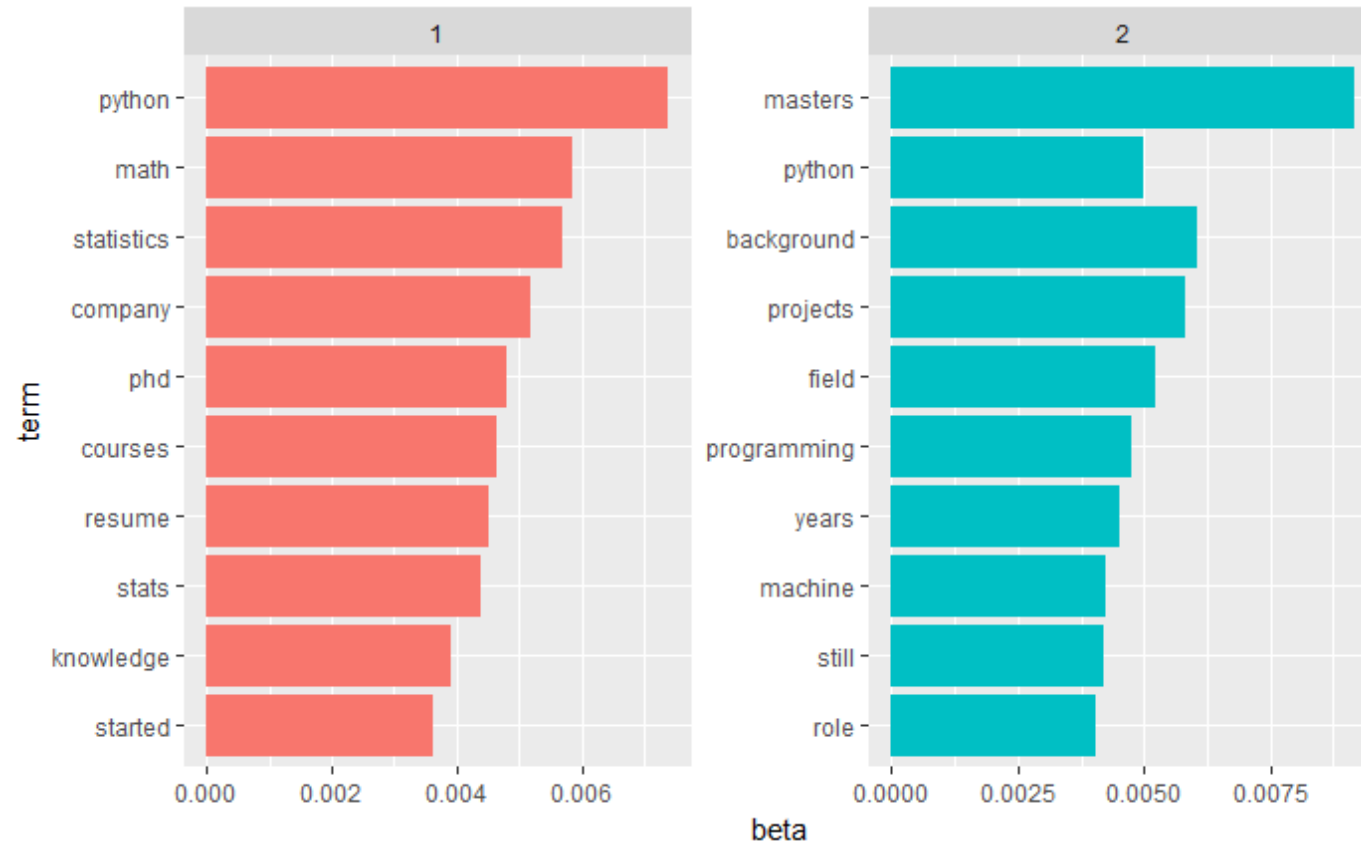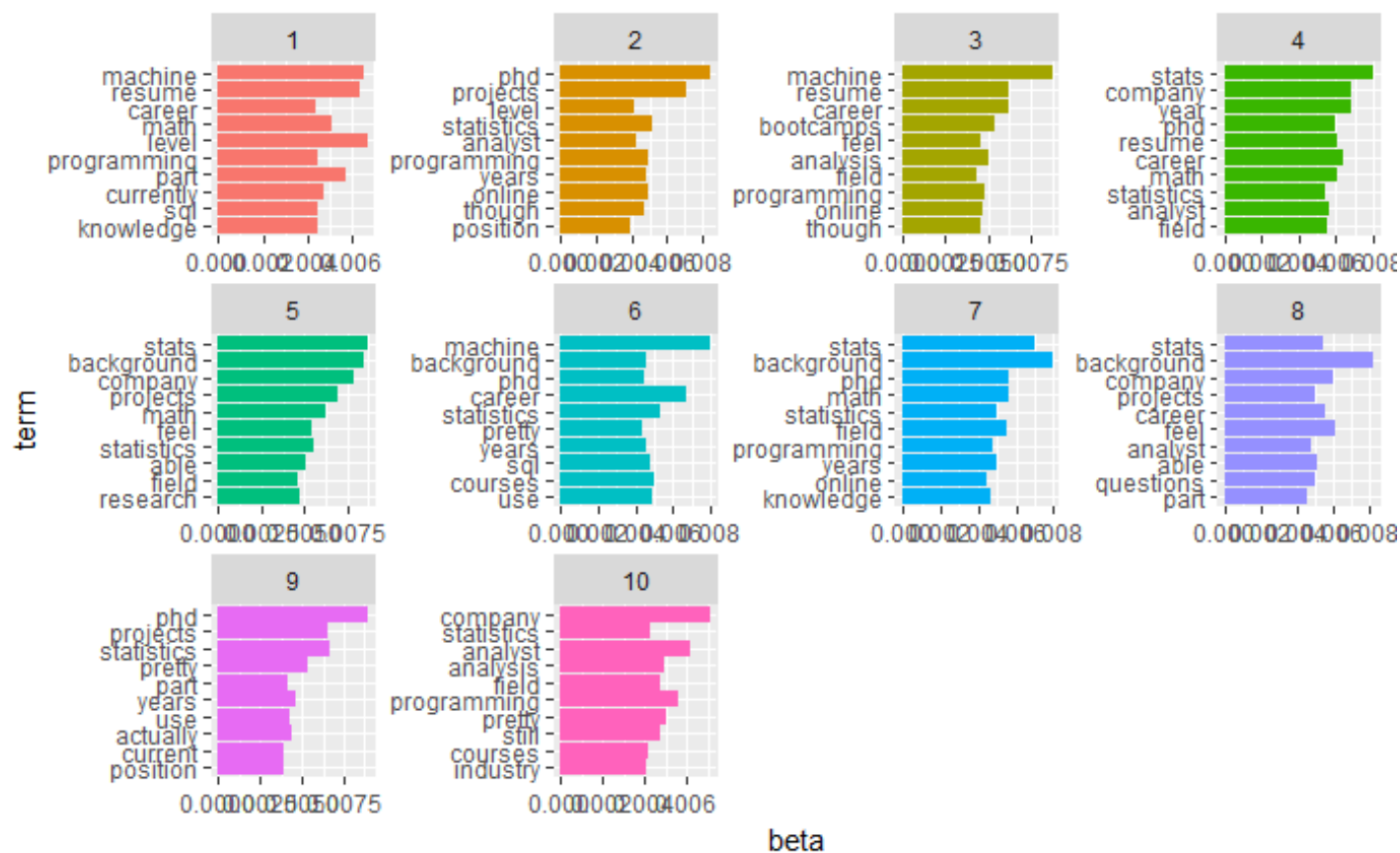
Total of 5 page lists

# Results

# Continued

# Continued

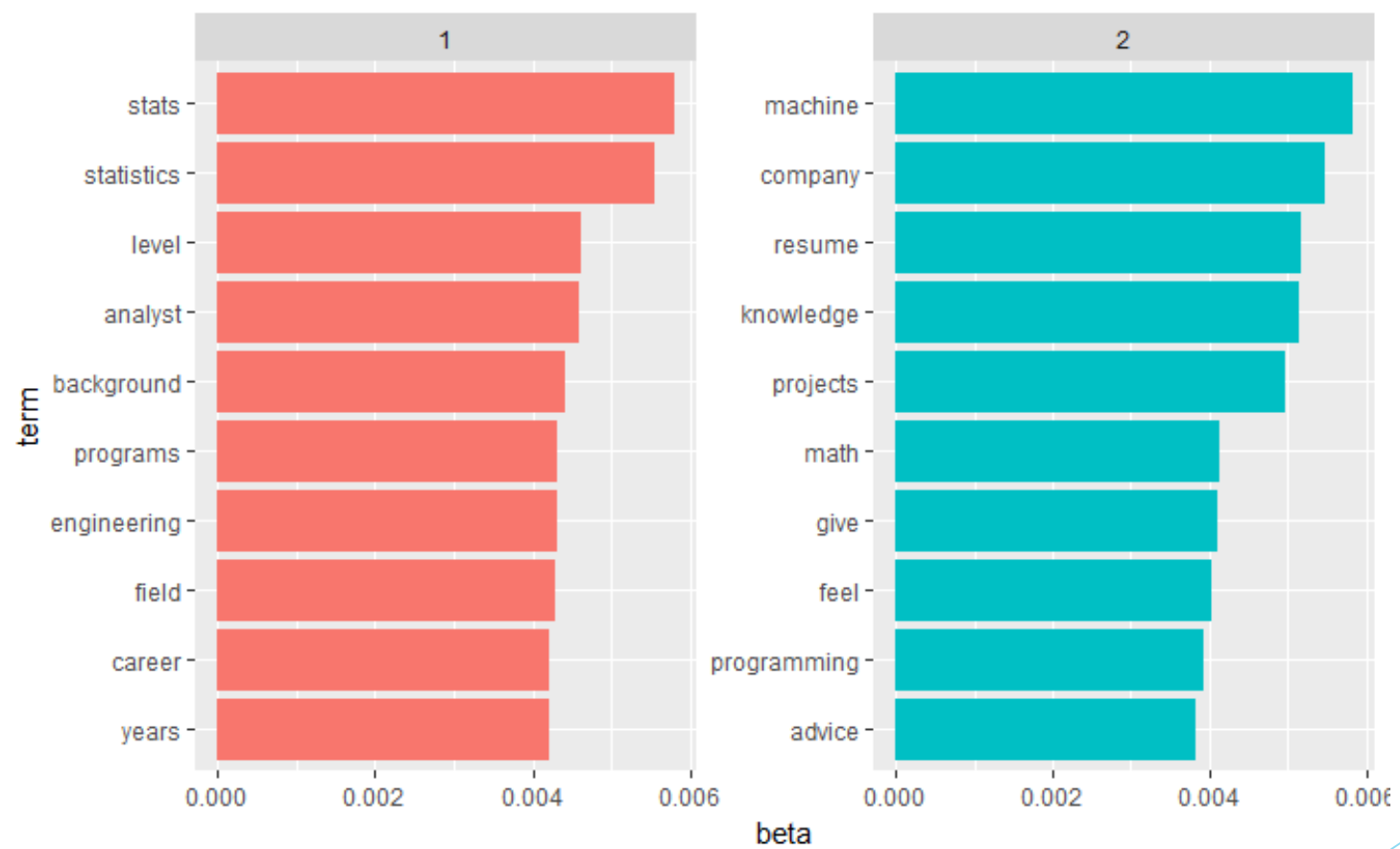# Continued

# Continued

# Continued

# Conclusion: Surprise!

- Major Topics of interest amongst reddit users: Job, Experience

- Python and masters and statistics seem to be other major topics of interest

- Limitations: number of words (algorithm clustering the same terms under different topics). Generalizability (sample may not reflect the interests of people interested in data science boot camps)

# The future

- Increase power of the topic modeling by adding more text and test even more numbers of topics.

- Scrape reviews based other discussion based sites like Quora and Hacker News
  And see how they relate to each other

- Compare to other programming/bootcamps

- Focus marketing strategies around related topics of interest
  Keep your marketing focused on helping people find a job and host sessions python