# N(eur)IPS

And an analysis of bleeding edge
ML research through Tweets

Barrett Williams — November 16, 2020

# Where do top-notch ML researchers congregate?

GitHub, obviously, but more formally:

- NeurIPS
- CVPR (computer vision)
- ICML (more statistical and theoretical)
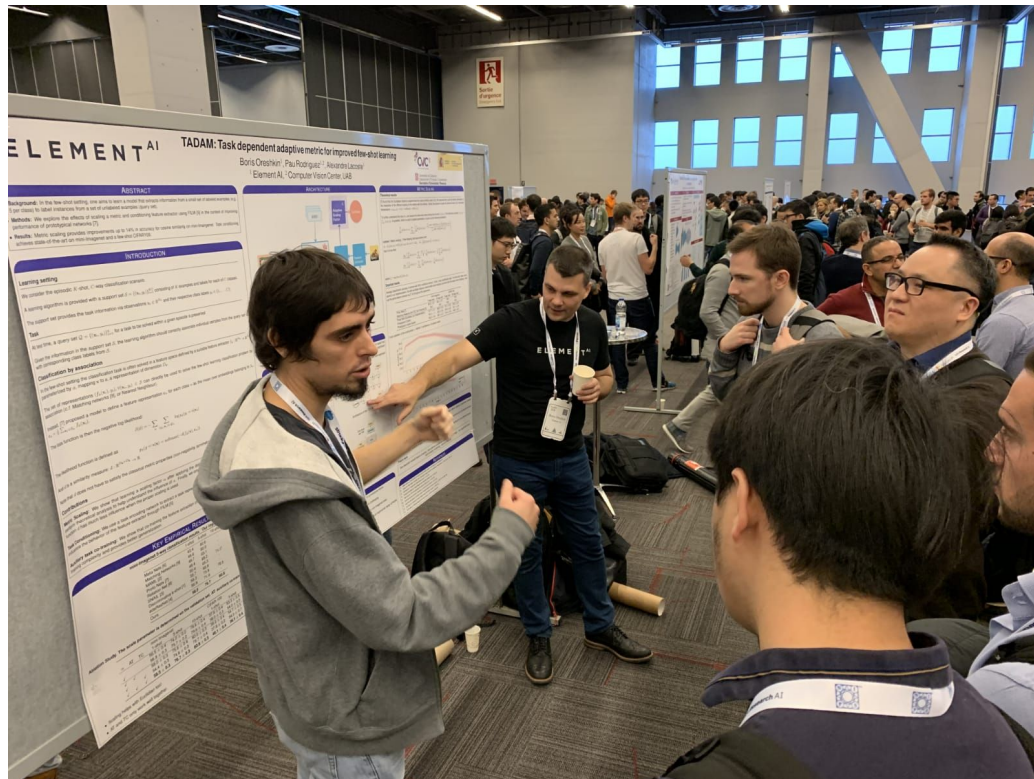
NEURAL INFORMATION
PROCESSING SYSTEMS

# Where do top-notch ML researchers congregate?

Also, Twitter!

Professional experience has taught me that most lab heads, PhDs, and post-docs will stay active on Twitter, especially around academic conferences.
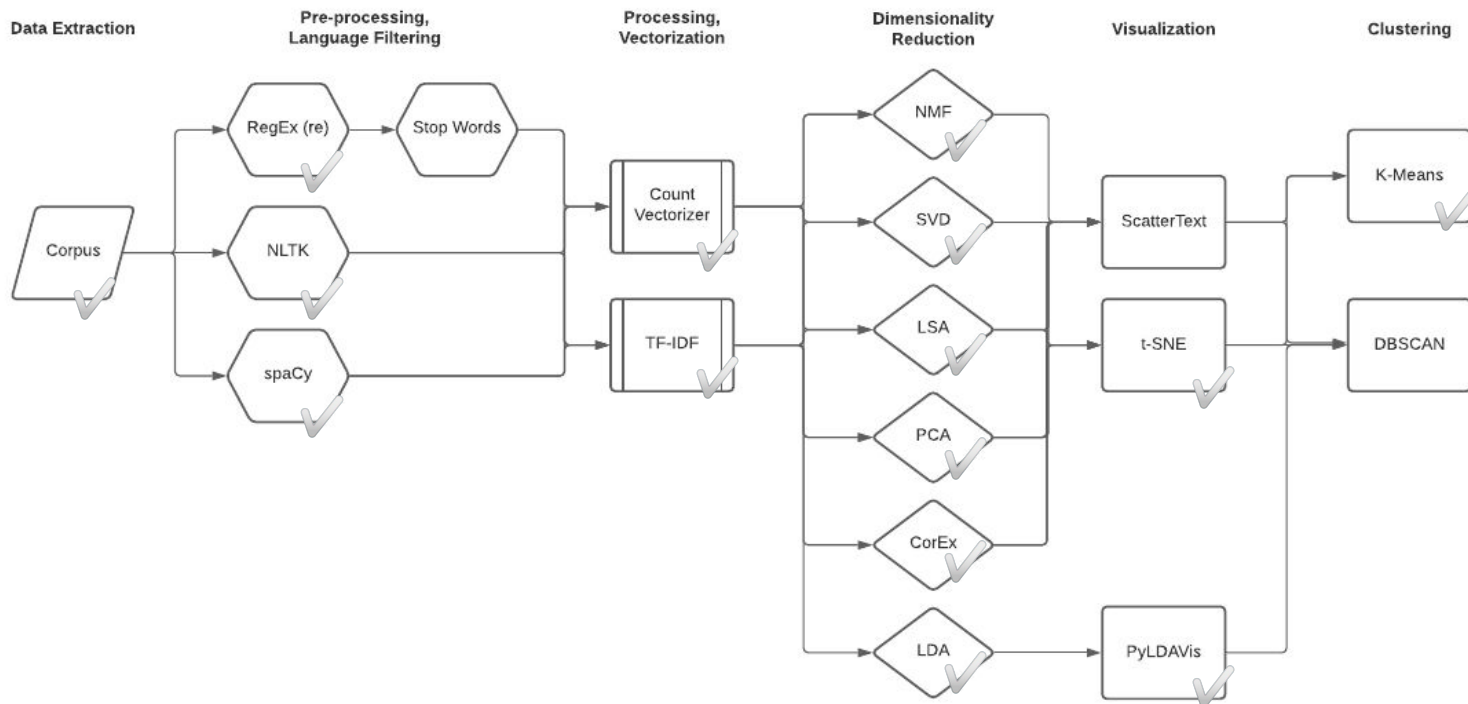




hardmaru @hardmaru · Dec 2, 2017
You know AI is a bubble when Intel invites Flo Rida to perform at their NIPS party

Intel Registration <Intel_Registration@regsvc.com>
to me

intel AI

Let the Gradient Flo
Celebrate NIPS 2017 with Intel AI

You're confirmed for the Intel AI party on December 5!

We look forward to seeing you at The Loft on Pine (230 Pine Ave, Long Beach, CA) at 9:00 PM.

Your conference pass will serve as your ticket to this event. Please ensure you are wearing it at all

40          330          780

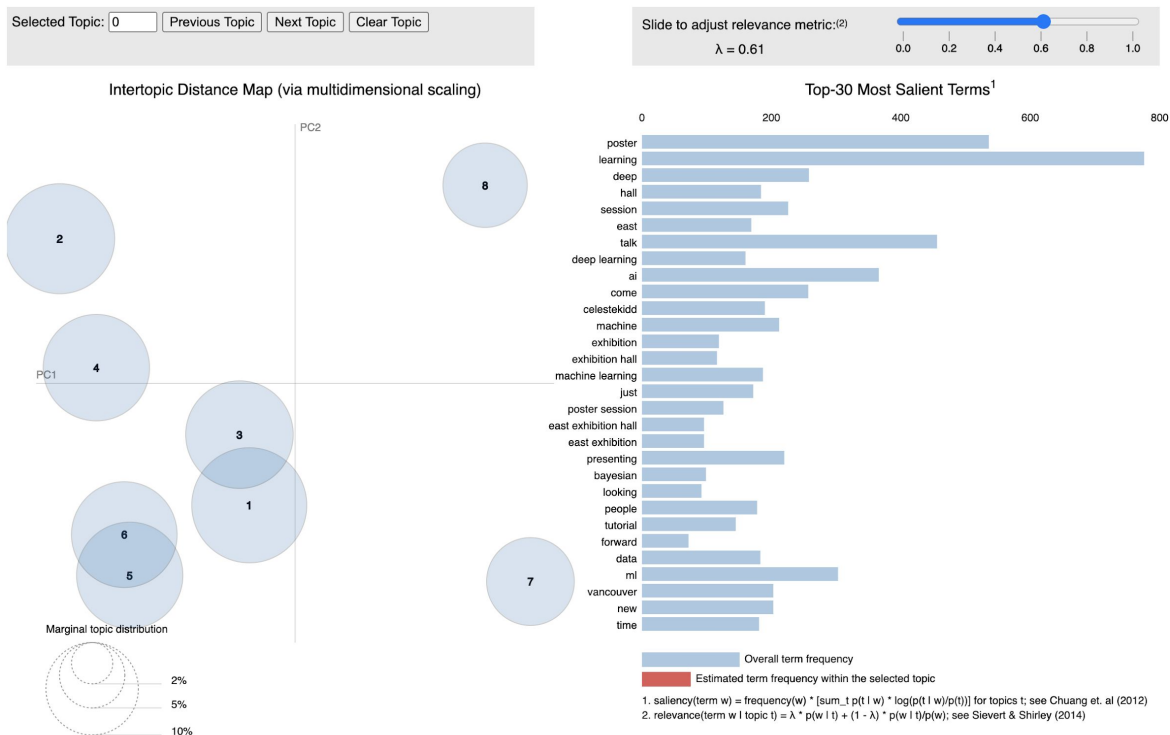# Poster sessions do indeed fill up (pre-COVID)

# Analysis and modeling workflow

I retrieved Tweet IDs with `snscrape`, then collected full text from `TweePy` and the Twitter API.

# LDA Visualization

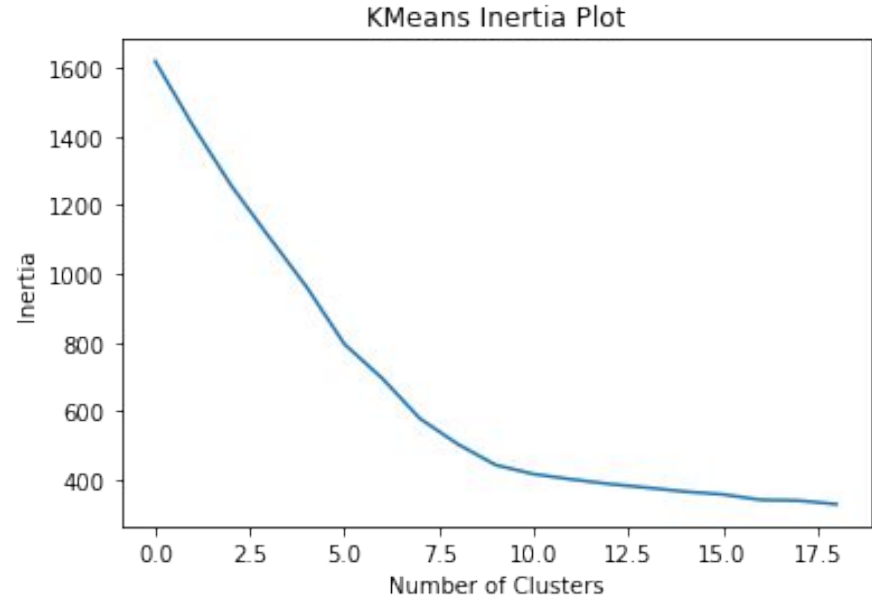8 topics marked the elbow of the inertias plot from K-means clustering:

# Topic categories: a stratified attendee experience

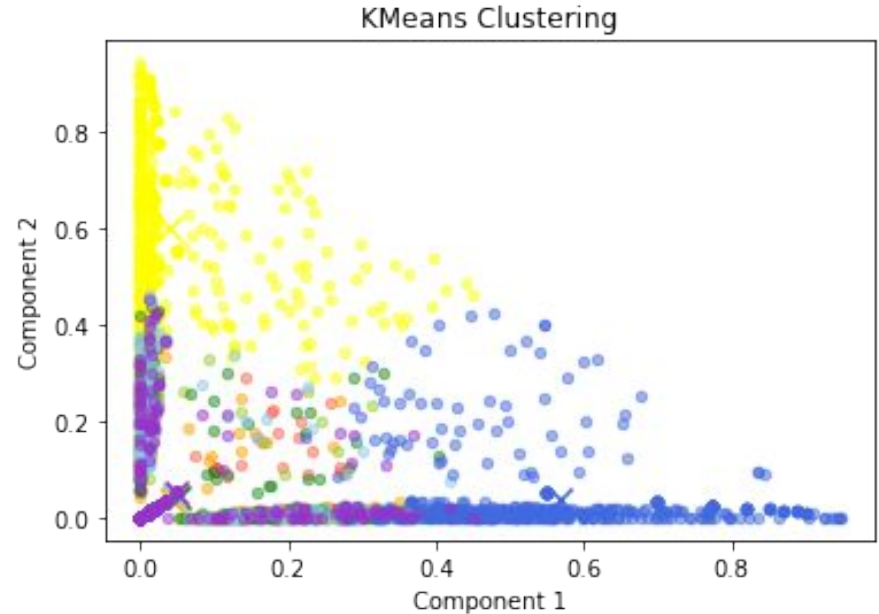| Num | Topic | Example |
|-----|-------|---------|
| 1 | Poster session | The poster session (east hall) was indeed packed |
| 2 | Climate Change Workshop | Workshop to use AI to tackle climate issues |
| 3 | Diversity Groups | LatinX in AI, Black in AI meetings |
| 4 | Sexual Harassment, Bayesian RL (how to separate?) | Celeste Kidd's presentation on avoiding, mitigating sexual harassment in the workplace |
| 5 | Yoshua Bengio and Google | Gaps in ML talk: "ML I to ML II" |
| 6 | Interpretable AI, avoiding bias | Interpretability in image recognition |
| 7 | Celeste Kidd | Highlight talk on the psychology of learning |
| 8 | NVIDIA DIB-R research | In conjunction with UToronto, interpolation rendering |

# Scree plot to determine dimensionality

Using K-means clustering, I was able
to find an "elbow" of 8 clusters:

# K-means clustering with 8 topics

Even with dimensionality reduction,
hard to establish separable clusters

Note centroids are indicated
by an "X" marker



KMeans Clustering

# Future work

- Topic modeling with arXiv abstracts might provide more precise insights on where ML research is heading, on a year-by-year basis
- CorEx modeling was interesting, but ultimately provided no useful output data; perhaps with better anchoring terms this might be different
- Establishing certain papers on arXiv as the centers of topic clusters might be illustrative
- Vader sentiment analysis to determine which areas of research are more highly criticized
- Focus on workshops, papers, or talks: these keywords might reveal a more topical representation of the event, even if used as CorEx anchors

# Thank you!

You can find me at:

https://linkedin.com/in/barwi

https://github.com/brwillia

And you can find the code for this project at:

https://github.com/brwillia/metis-neurips-twitter-proj-4

Questions?

# Appendix

# Things to know going in:

1. There are too many processes to apply in NLP: you can't try it all
2. Labels are still helpful; I considered using "has arXiv link" as a label
3. Don't forget to standardize/de-mean when running PCA
4. t-SNE may not work, and that's OK (see next slide)
5. PyLDAVis will break your Jupyter Lab, but not VS Code
6. Don't download your own Word2Vec, the package will do it for you
7. There isn't a "best practice" for cleaning tweets, but consider the following:
   a. Hashtags
   b. @ mentions
   c. Links
   d. Emojis
   e. Language detection (somewhat time-consuming to run on a large corpus)