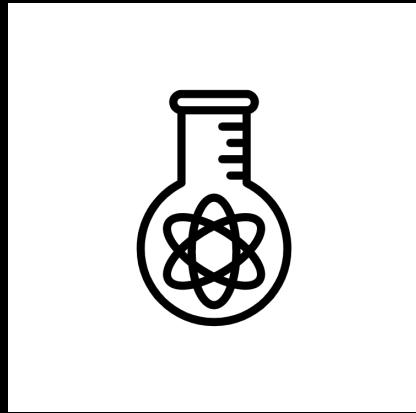


**tidylda**

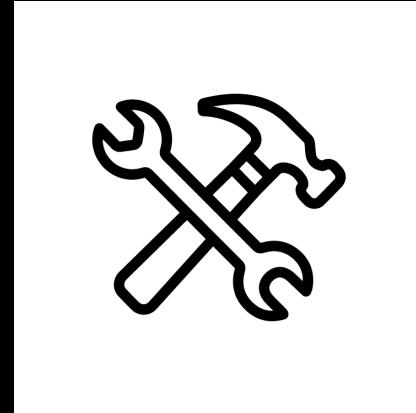
**Tommy Jones**  
Data Science DC  
June 14, 2022



# How to Take Over the World



Science!



Tools!

# Science = “Corpus Statistics”

Task-based  
vs.  
Inference on populations

- Sample composition
- Experimental design
- Model (mis)specification
- Estimate population parameters
- Uncertainty quantification
- etc...

Dictionary

🔍

 **statis·tics**

/stə'tistikəs/

*noun*

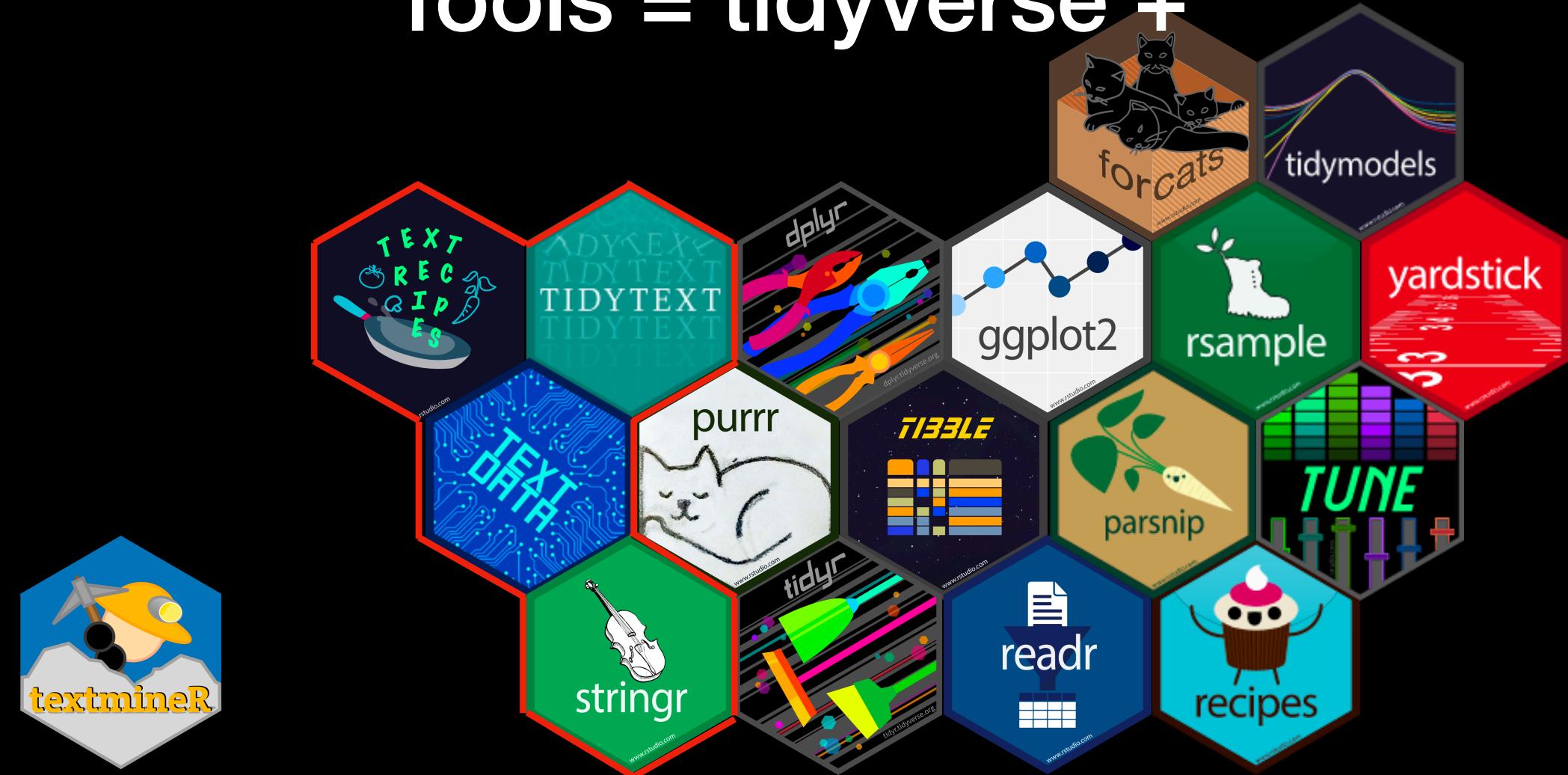
the practice or science of collecting and analyzing numerical data in large quantities, especially for the purpose of inferring proportions in a whole from those in a representative sample.

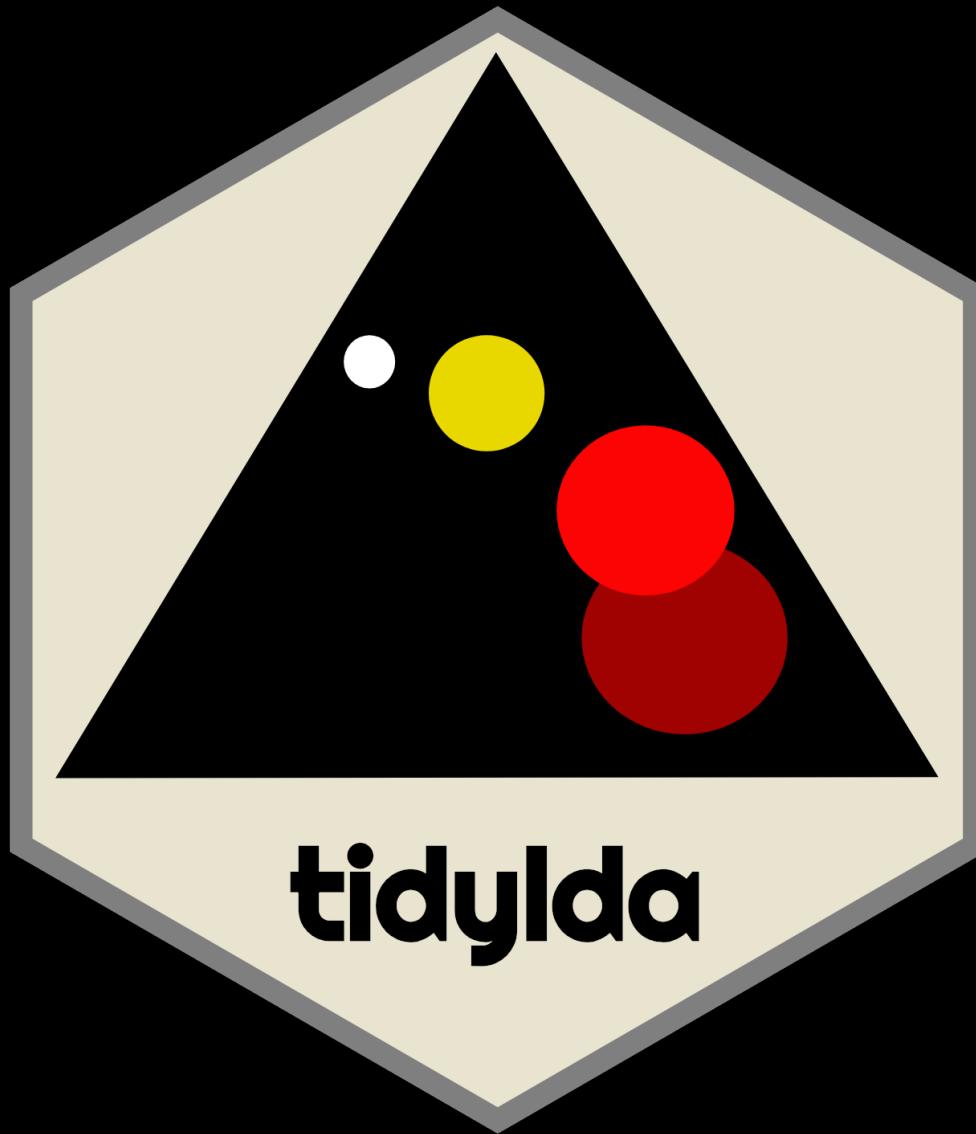
Definitions from Oxford Languages Feedback

# Latent Dirichlet Allocation for Corpus Statistics

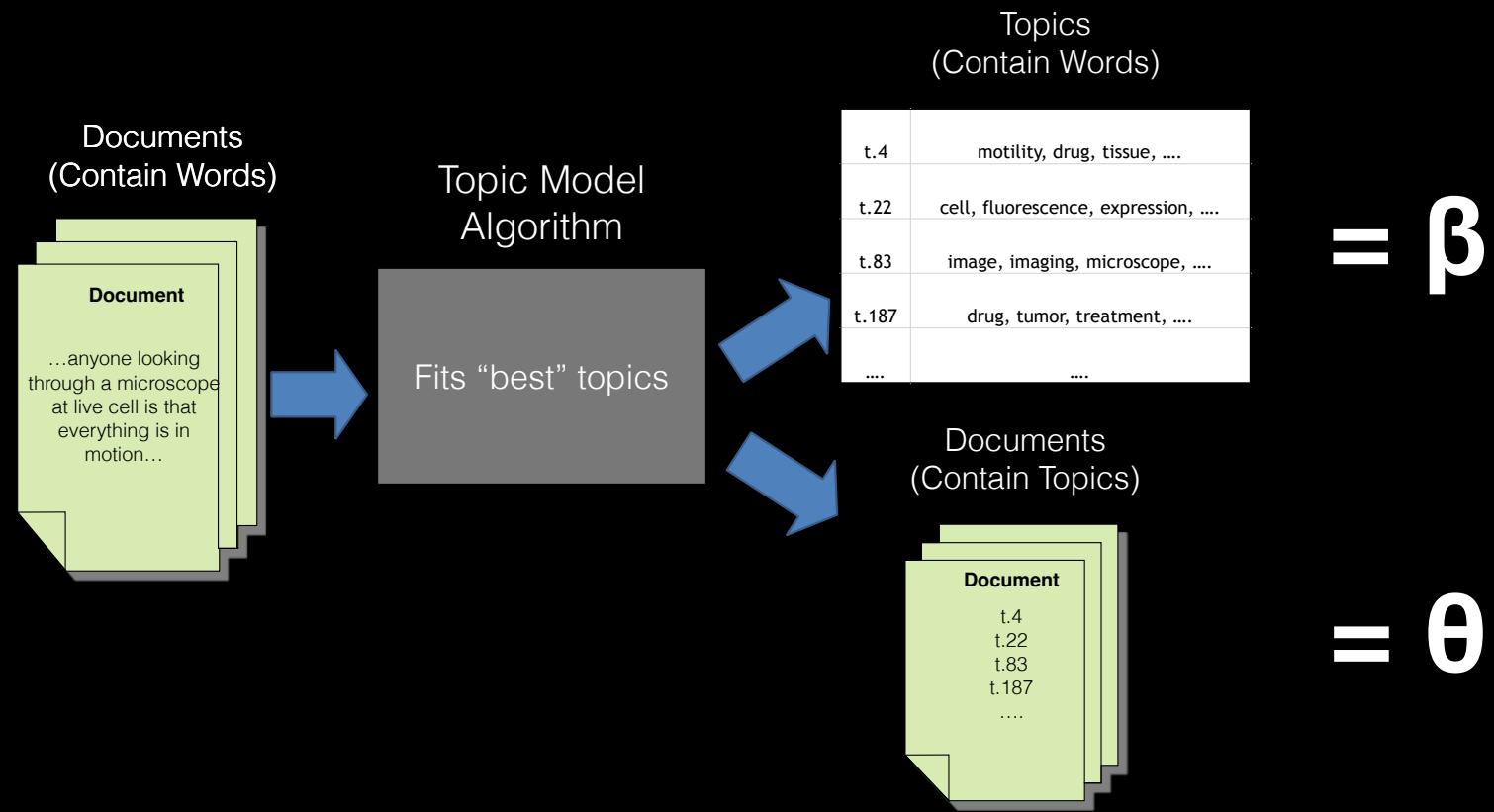
- **LDA embeds text into a probability space**  
Well-defined and interpretable relationships
- **LDA is a Bayesian parametric model**  
Uncertainty quantification, diagnostics, etc. in line with established statistical practice
- **LDA is a generative model of language**  
Study the LDA-DGP to build better models and diagnose pathological misspecification

# Tools = tidyverse +

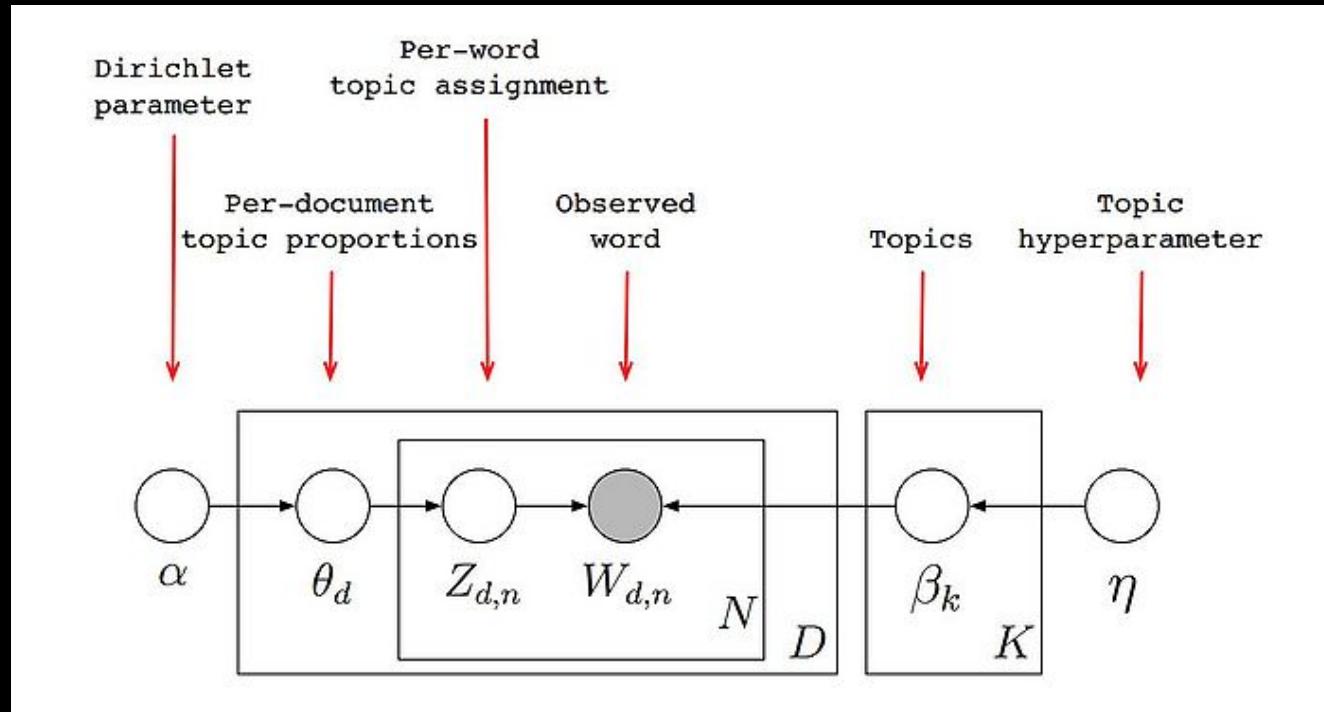




# Quick Review of LDA I



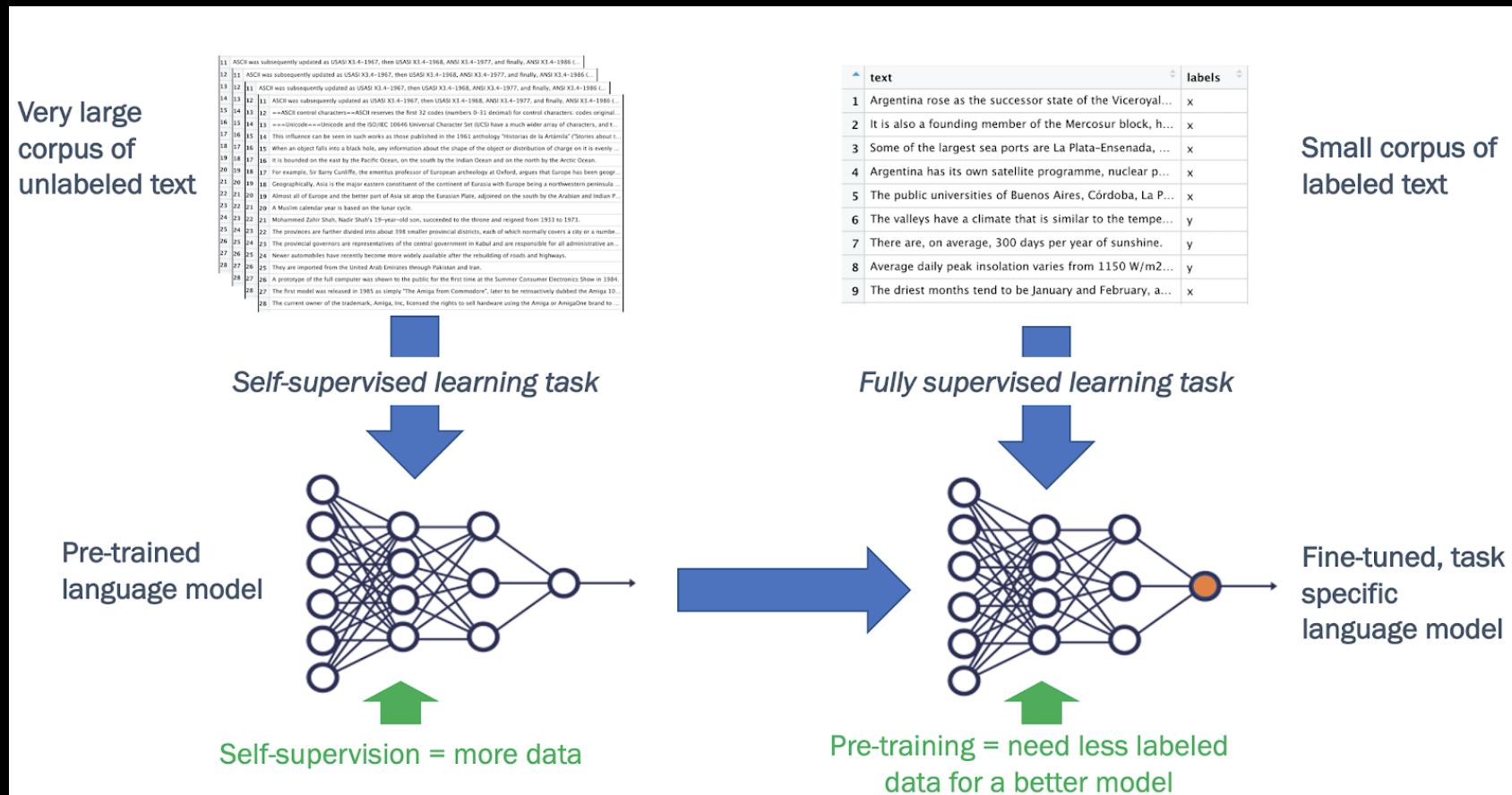
# Quick Review of LDA II



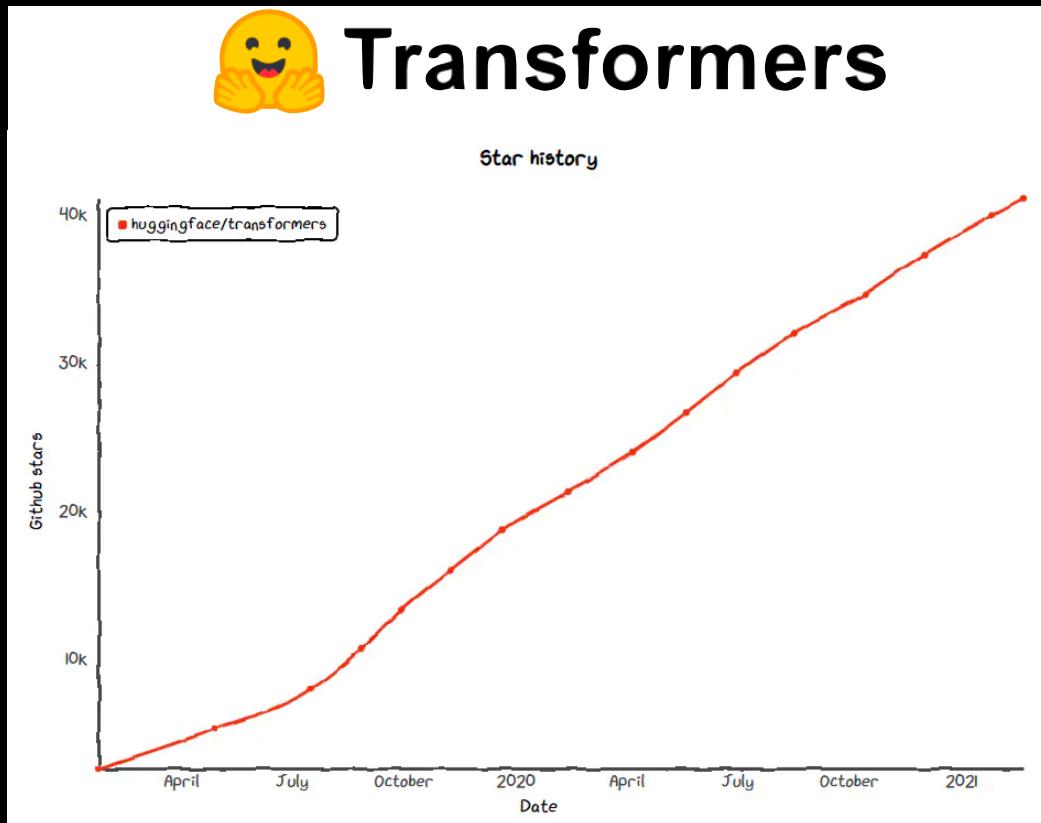
- $\eta, \alpha$ : priors
- $\beta, \theta$ : parameters of interest
- $w$ : observed word frequencies
- $z$ : latent topic frequencies

# Demo Time

# A Paradigm Shift in NLP



# Reign of the Pre-Trained Language Model



# Fine Tuning LDA for Transfer Learning

$$\vec{\beta}_k^{(t)} \sim \text{Dirichlet} \left( Cv_k^{(t)} + a^{(t)} \cdot \omega_k^{(t)} \cdot \mathbb{E} \left[ \vec{\beta}_k^{(t-1)} \right] \right)$$

Posterior

Data

Weight

Base Posterior



Prior

# Demo Time

# When does tLDA converge?

- tLDA has one tuning parameter:  $a^{(t)}$
- Do certain values of  $a^{(t)}$  allow/not allow tLDA to converge at all?
- Do certain values of  $a^{(t)}$  allow/not allow tLDA to converge towards the correct distribution?

# Revisiting the key equation

$$\vec{\beta}_k^{(t)} \sim \text{Dirichlet} \left( Cv_k^{(t)} + a^{(t)} \cdot \omega_k^{(t)} \cdot \mathbb{E} \left[ \vec{\beta}_k^{(t-1)} \right] \right)$$

Posterior

Data

Weight

Base Posterior



Prior

# Simulation Experiment: Setup

- 128 data generating distributions
- 100 corpora of 10,000 documents sampled from each
- For each corpus:
  - Train a model on the first 100 documents
  - Iteratively add 100 documents to update the model
- Do the above for  $\alpha$  in  $\{0.2, 0.4, \dots, 1, 1.2, 1.4, \dots, 2\}$
- Does the model converge? Does it converge towards the data generating distribution?

# Simulation Experiment: Results I

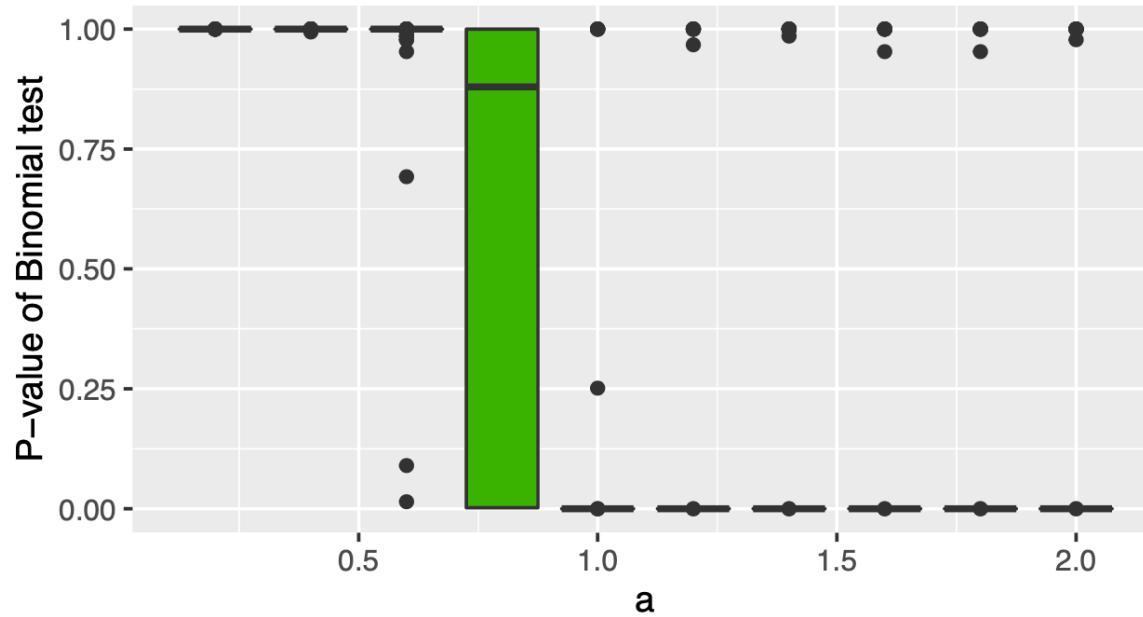


Figure 6: Boxplot of the p-values assessing convergence for each value of  $a$ . Below  $a = 0.8$  almost all p-values are near one. Above  $a = 0.8$  almost all p-values are near zero, with a few exceptions. But at  $a = 0.8$  there is high variance in the p-values assessing model convergence.

# Simulation Experiment: Results II

Table 1: Effects of  $a$  on convergence

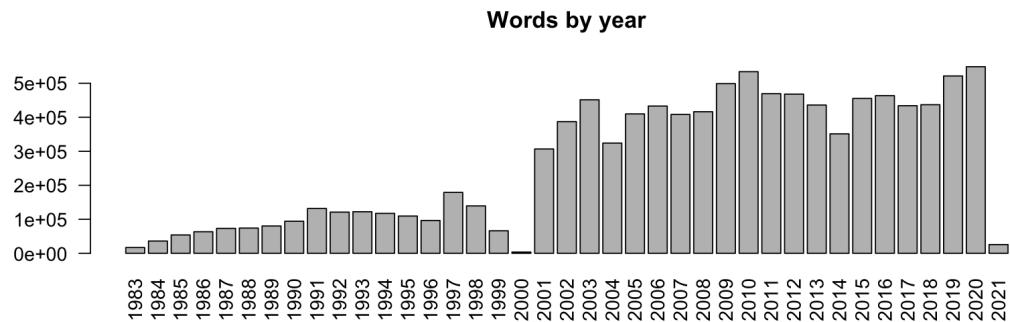
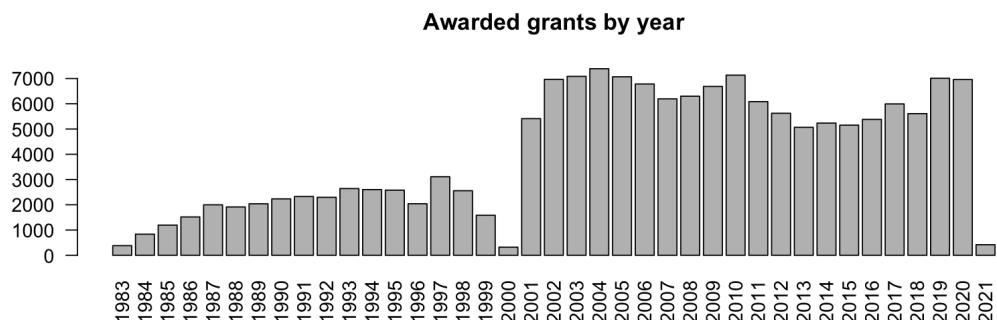
Variable	Convergence		Convergence		Direction		Direction	
Intercept	-6.01	***	-17.09	***	64.59	***	65.19	***
$\$a\$$	5.92	***	28.38	***	1.94	***	8.62	***
Avg. Doc. Length	0.00		-0.01	**	0.01	***	0.02	***
$\$\\sum(\\boldsymbol\\eta) \$$	0.00		0.00	**	-0.02	***	-0.06	***
$\$\\sum(\\boldsymbol\\alpha) \$$	0.32	***	0.40		-0.18	***	0.01	
$\$a^2 \$$			-9.20	***			-2.28	***
$\$(\\text{Avg. Doc. Length})^2 \$$			0.00	**			0.00	***
$\$\\sum(\\boldsymbol\\eta)^2 \$$			0.00	**			0.00	***
$\$\\sum(\\boldsymbol\\alpha)^2 \$$			0.02				-0.02	**

# The Data

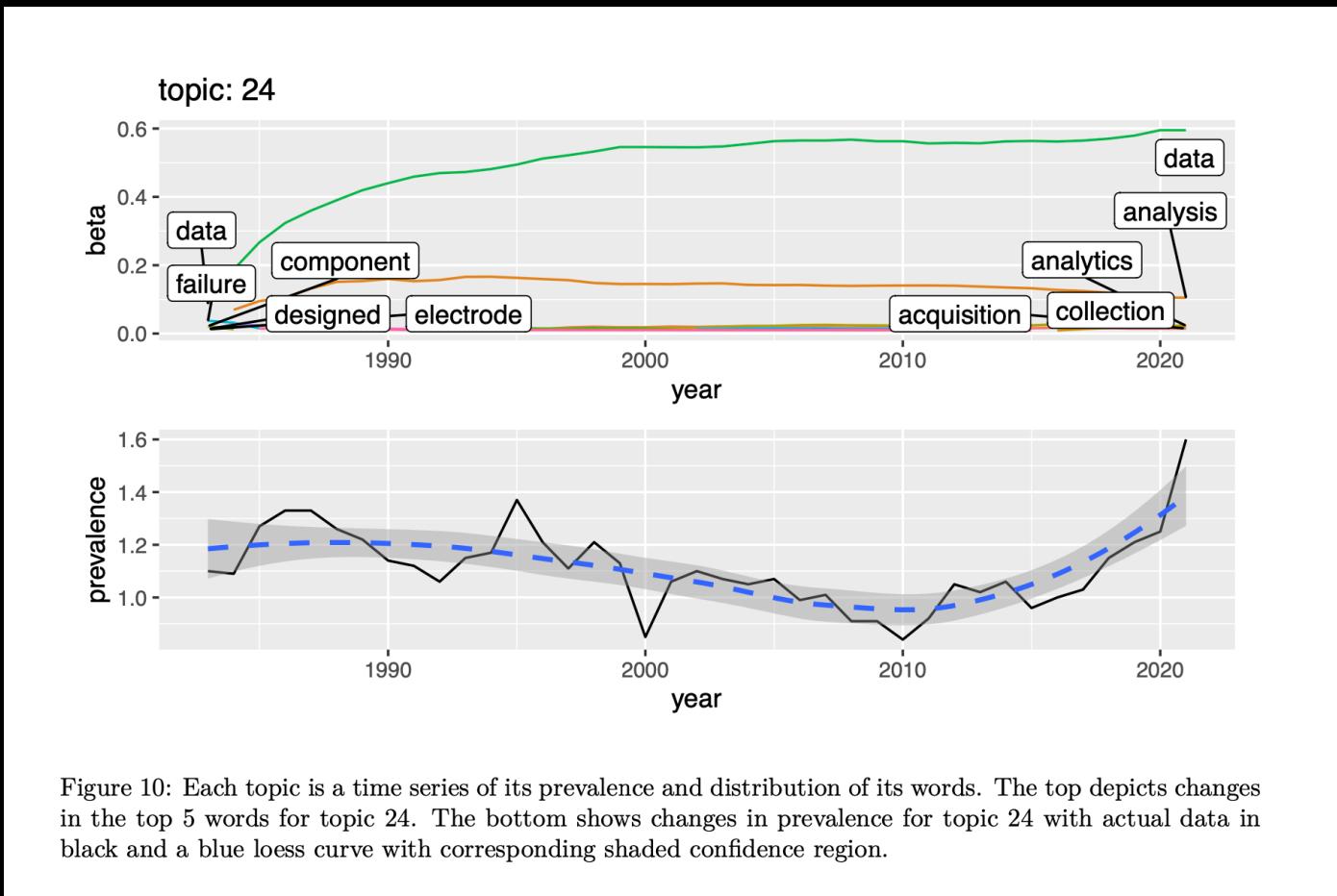
Small Business Innovation  
Research (SBIR) Program

Technology R&D grants from  
US Small Business Admin.

159,734 abstracts of  
awarded grants from  
1983 - 2021



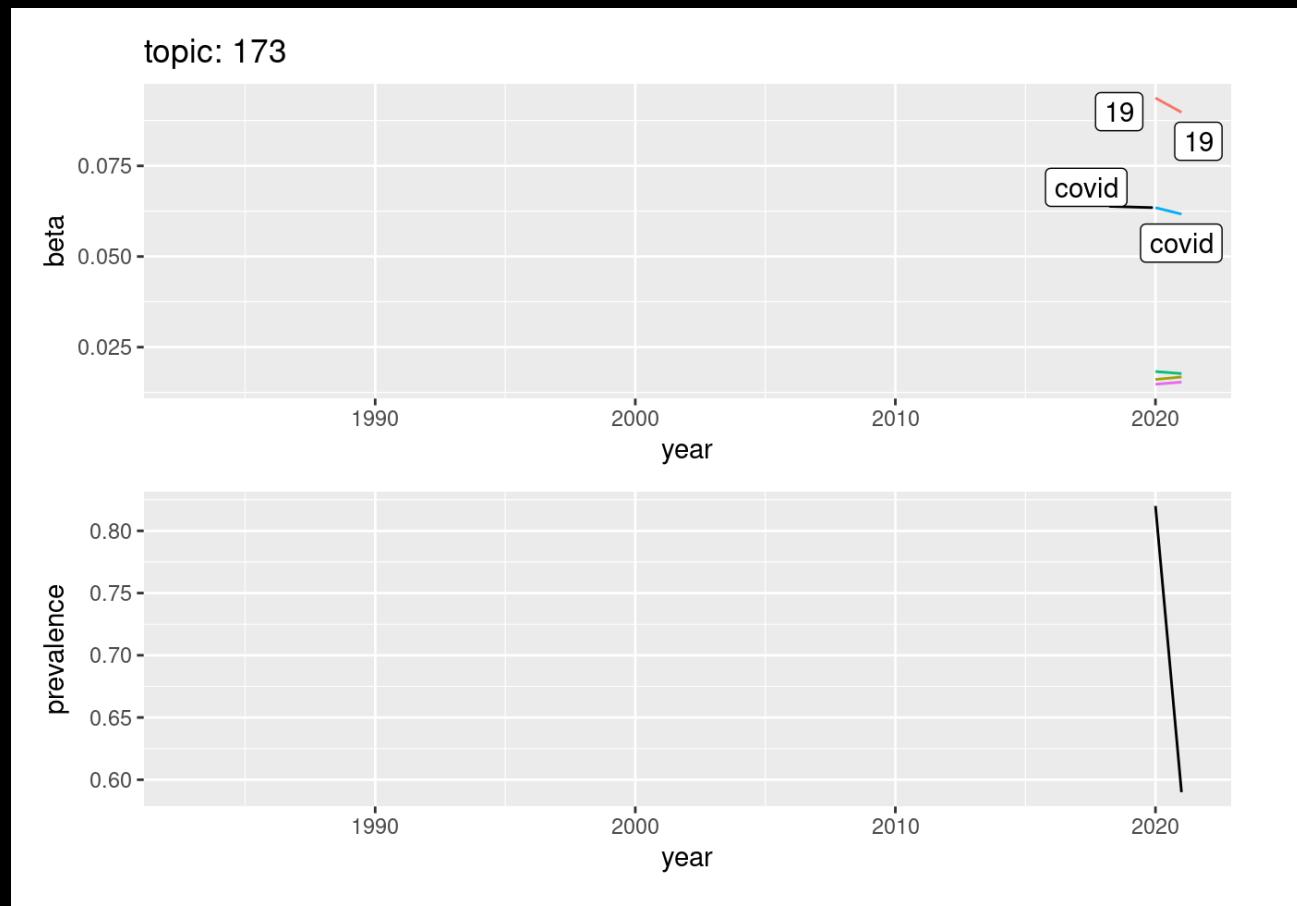
# SBIR Time Series



# SBIR Results

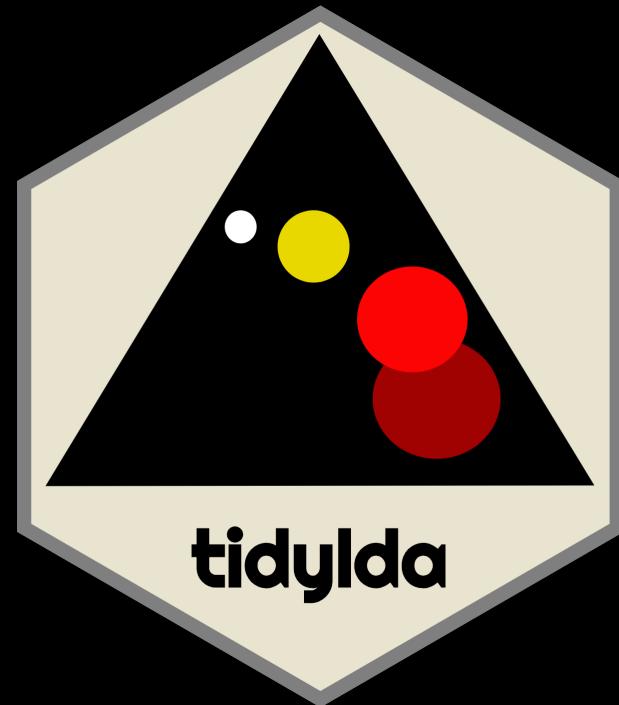


# SBIR Time Series II



# Next Steps for tidylda

- Speed and scalability
- UX for seeding topics
- Helper functions for common tasks
- Fold tidylda into textmineR



# Thank You!



[jones.thos.w@gmail.com](mailto:jones.thos.w@gmail.com)



@thos\_jones



<https://github.com/tommyjones>

