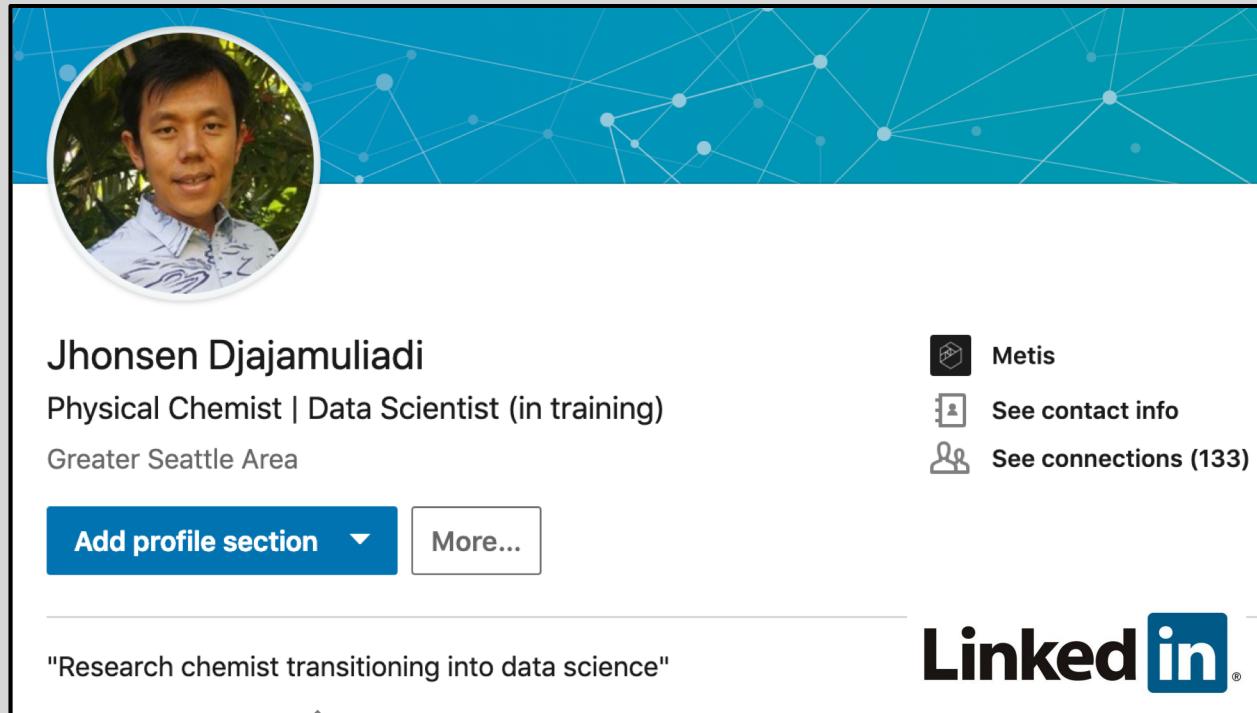


# *“Am I a Data Scientist?”*

## Finding My Identity Through Unsupervised Learning

Project-4 @ **METIS**  
Jhonsen Djajamuliadi

# *My Journey in Becoming a Data Scientist*



Jhonsen Djajamuliadi  
Physical Chemist | Data Scientist (in training)  
Greater Seattle Area

Add profile section ▾ More... See contact info See connections (133)

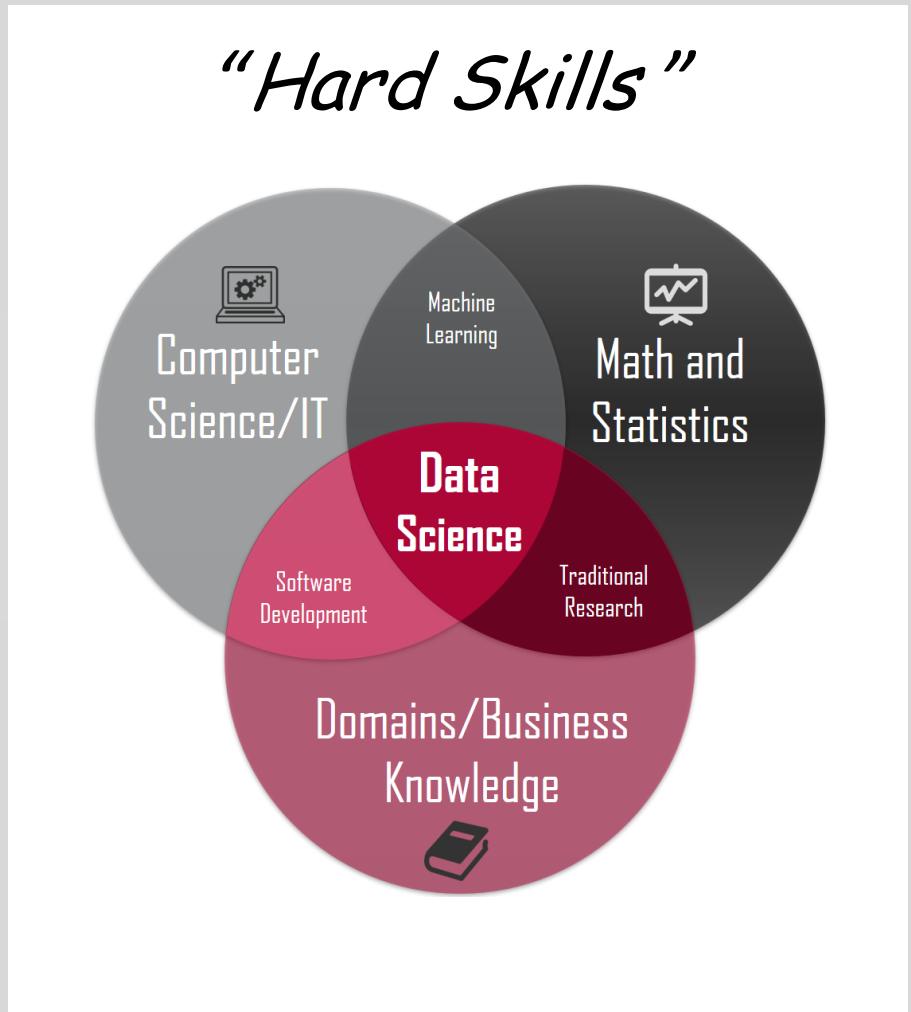
"Research chemist transitioning into data science"

LinkedIn



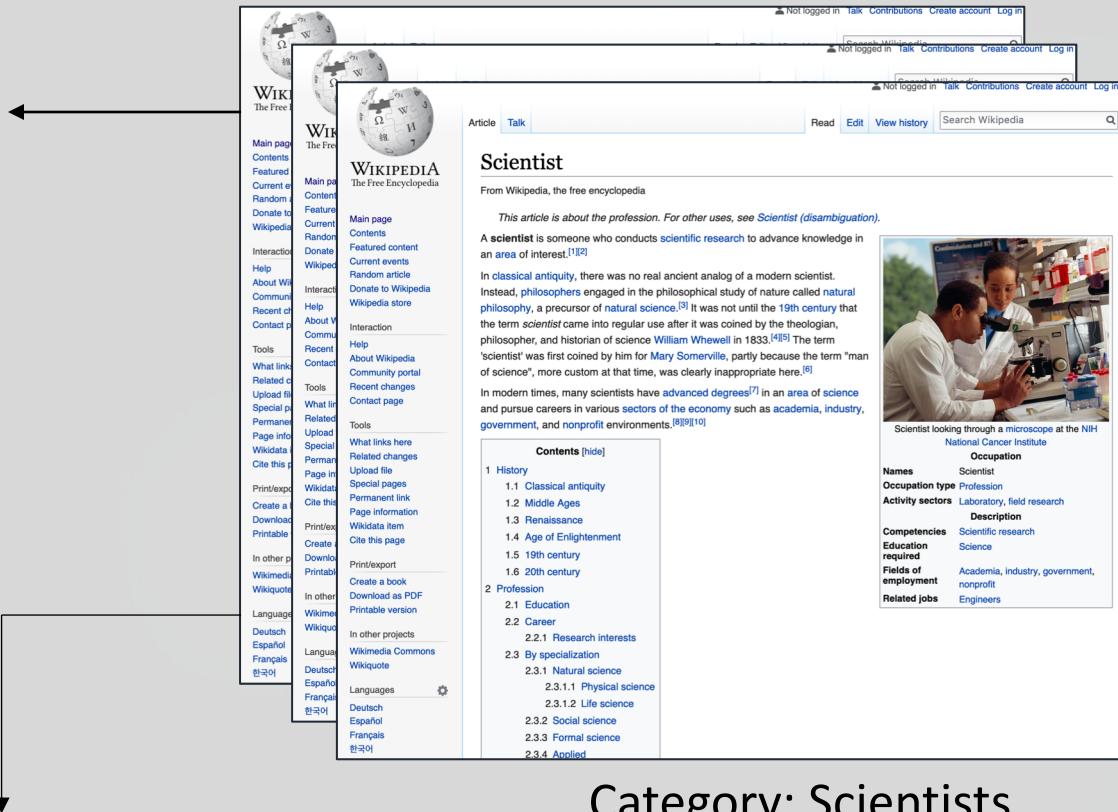
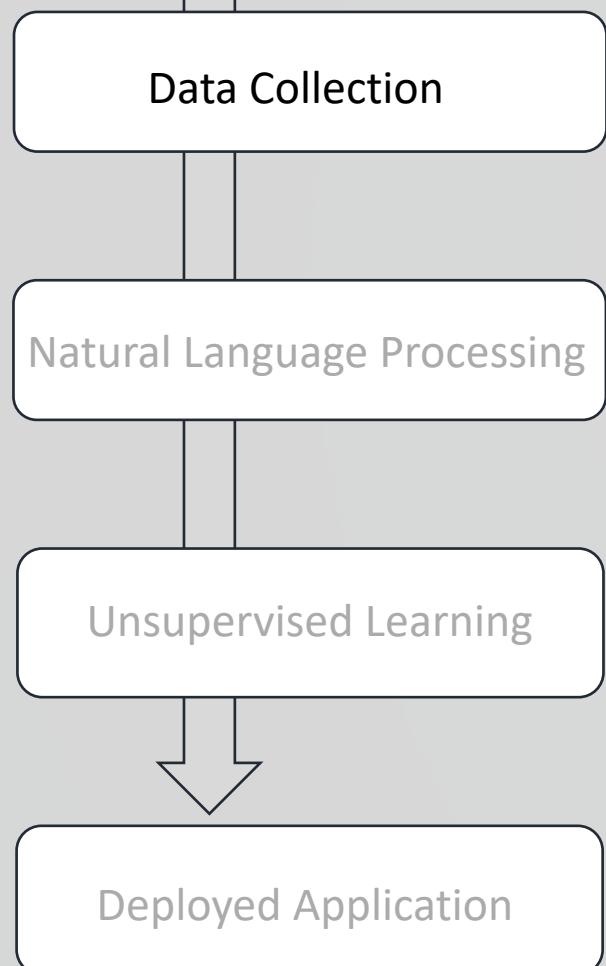
## *"Soft Skills"*

- How should I describe myself?
- How should I network?
- How should I negotiate?



# Seeking My Identity in Wikipedia Articles

- What is a **data scientist**?
  - Is it similar to other **scientists**?



## Category: Scientists

~15,000 articles x 1.2 million words

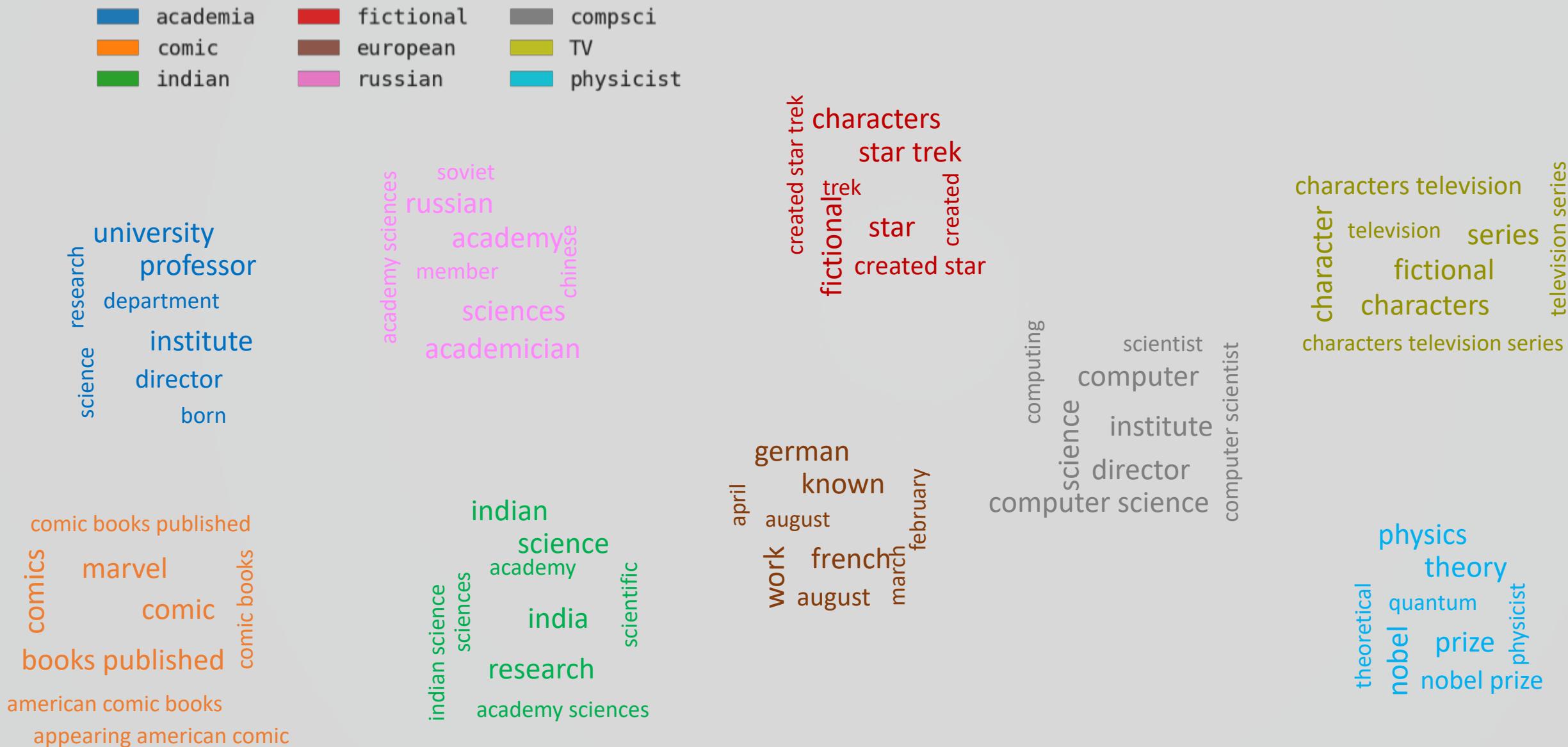
- ✓ Dimensionality reduction
- ✓ “Topics” in documents

## Non-negative Matrix Factorization

$$\text{scientists} \left\{ \begin{matrix} W \\ \times \\ H \end{matrix} \right\} \approx \begin{matrix} V \\ \approx \\ \text{words} \end{matrix}$$

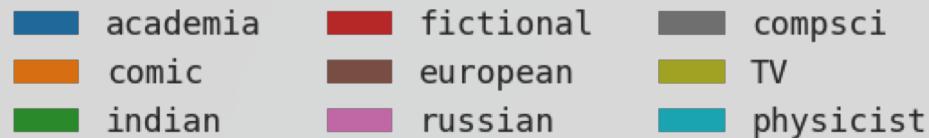
# Finding Topics in Summaries

## Topic Modeling



# Finding Structure Among Scientists

## Topic Modeling and 2D-Embedding with T-SNE



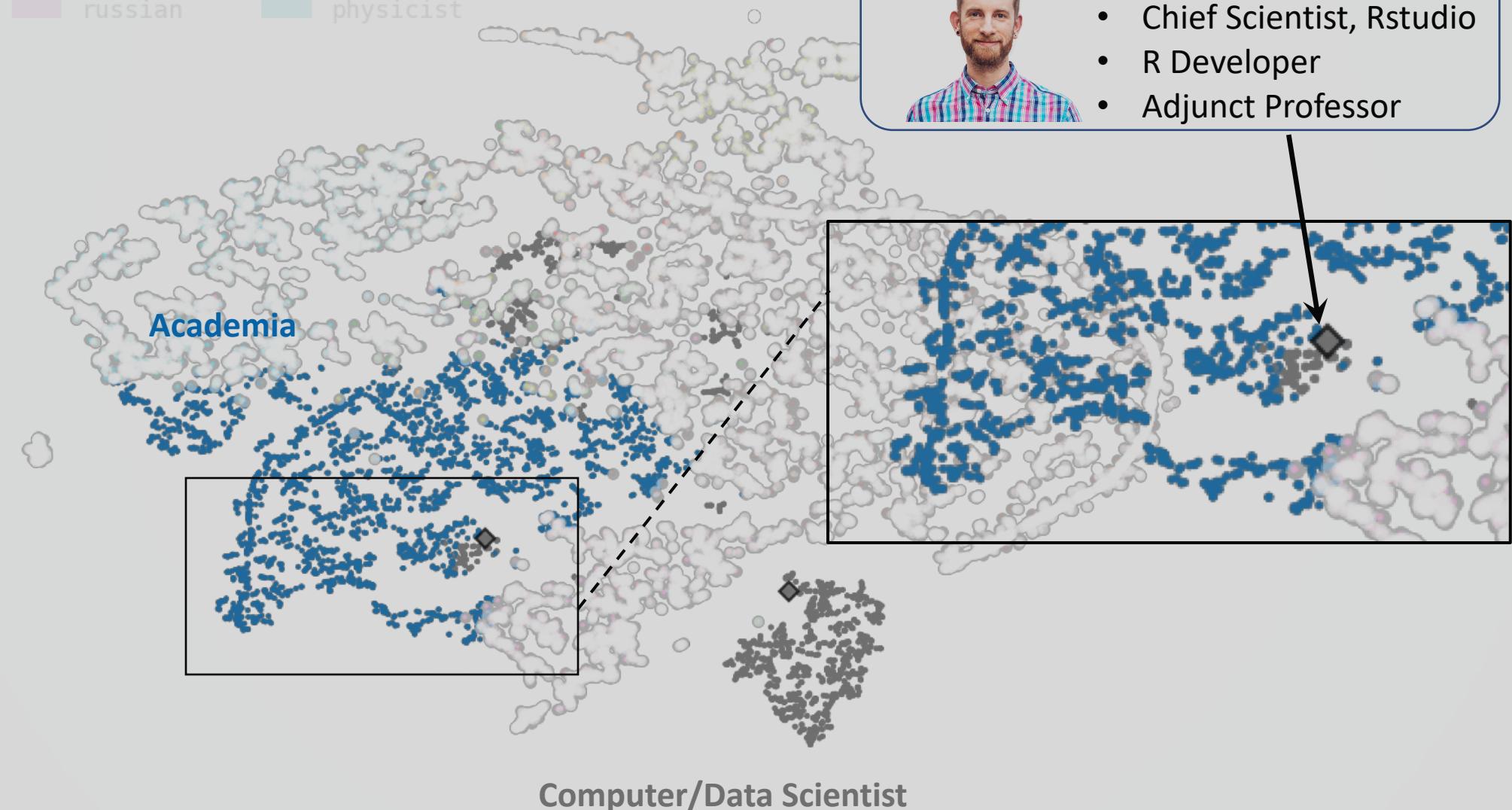
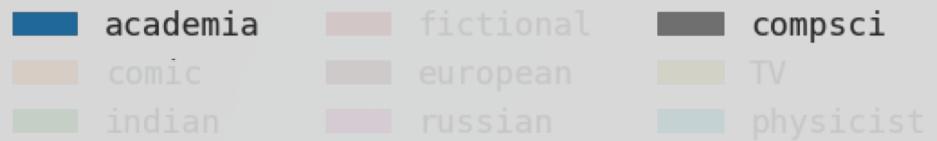
t-Distributed Stochastic  
Neighbor Embedding (t-SNE)<sup>[1]</sup>

- Each point – a scientist
- Colors – heaviest topic-weights
- Distances  $\propto$  similarity

[1] T-SNE, <https://lvdmaaten.github.io/tsne/>

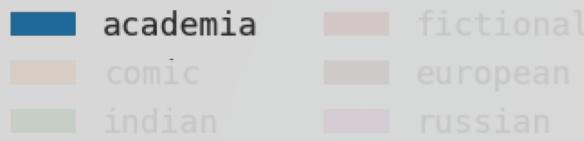
# Finding Structure Among Scientists

## Topic Modeling and 2D-Embedding with T-SNE



# Finding (Cosine)-Similarity with Others

## Topic Modeling and 2D-Embedding with T-SNE



*"Am I a Data Scientist?"*

**Jhonsen Djajamuliadi**

- Research Chemist
- Data Scientist
- Machine Learning

**Ricardo Bianchini**

- Microsoft Research
- American Computing Fellow
- Data & Server Management

vs.

Similarity  $\sim 0.95$

Academia

Computer/Data Scientist

We say: YES!!!

# Web Application in Development

## Input:

- **LinkedIn** Summary

"Research chemist transitioning into data science"

I have an academic background in spectroscopy and molecular modeling, which are the science of "extracting signals out of the noise" and "computationally recreating or simulating their interactions".

I love integrating experimental techniques with computational approaches, to find actionable insights and most consistent answer to research questions. "Hands-on" laboratory work is fun, but "in silico" computational projects have always peaked my interest. The latter has led me to this exciting field of data science. I'm currently taking a deeper dive into machine learning and AI for industrial applications, especially, in the biomedical and healthcare fields.

I'm also passionate about science communication, which is a way of "storytelling using data". I find it enjoyable to decompose technical concepts, and convey them to a mixed audience, e.g., scientists in different disciplines or others without scientific backgrounds.

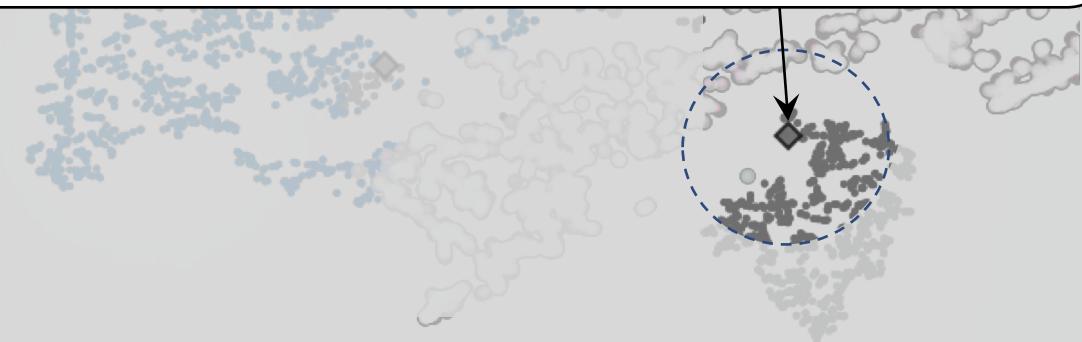
Submit

## Output:

- **Top 10** (or n-) similar scientists

### Ricardo Bianchini

Ricardo Bianchini from [Rutgers University](#) & [Microsoft Research](#), Bellevue, WA was named [Fellow of the Institute of Electrical and Electronics Engineers \(IEEE\)](#) in 2015<sup>[1]</sup> for contributions to server and data center energy management. He was named an [Association for Computing Machinery \(ACM\) Fellow](#) in 2016<sup>[2]</sup> for contributions to power, energy and thermal management of servers and datacenters.



## Future Work:

- "Cleaner data", **LinkedIn**
- Target data science domain

*"Am I a Data Scientist?"*

# Web Application in Development

## Input:

- **LinkedIn** Summary

"Research chemist transitioning into data science"

I have an academic background in spectroscopy and molecular modeling, which are the science of "extracting signals out of the noise" and "computationally recreating or simulating their interactions".

I love integrating experimental techniques with computational approaches, to find actionable insights and most consistent answer to research questions. "Hands-on" laboratory work is fun, but "in silico" computational projects have always peaked my interest. The latter has led me to this exciting field of data science. I'm currently taking a deeper dive into machine learning and AI for industrial applications, especially, in the biomedical and healthcare fields.

I'm also passionate about science communication, which is a way of "storytelling using data". I find it enjoyable to decompose technical concepts, and convey them to a mixed audience, e.g., scientists in different disciplines or others without scientific backgrounds.

Submit

## Output:

- **Top 10** (or n-) similar scientists

### Ricardo Bianchini

Ricardo Bianchini from [Rutgers University](#) & [Microsoft Research](#), Bellevue, WA was named [Fellow of the Institute of Electrical and Electronics Engineers \(IEEE\)](#) in 2015<sup>[1]</sup> for contributions to server and data center energy management. He was named an [Association for Computing Machinery \(ACM\) Fellow](#) in 2016<sup>[2]</sup> for contributions to power, energy and thermal management of servers and datacenters.



## Future Work:

- "Cleaner data", **LinkedIn**
- Target data science domain

*"Are you a good fit?"*

## Application to Industry:

- Tool for company recruiters
- "New feature" on **LinkedIn**

Candidates

Employers



# Thank You

---

---

Physical Chemist

| Research Scientist

| Data Analyst

| Data Scientist

