

Using Embeddings to Understand the Variance and Evolution of Data Science Skill Sets

Maryam Jahanshahi Ph.D.

Research Scientist

TapRecruit.co

Jobs are hard to categorize

Research Analyst

Entry-Level

DRG - New York City

We are seeking a Research Analyst to join our team which serves our pharmaceutical clients with actionable data. As a Research Analyst, you will provide data support for client-facing platforms, presentations, and client requests...

Research Analyst

Senior

MMP - New York City

We have an exciting opportunity for an individual to join MMP's Cyber Risk Group. The successful candidate will have the ability to shape our investors service strategy, analytic, research and outreach framework for cyber risk and its relationship to credit and the financial markets...

Research at TapRecruit

Helping companies make fairer and more efficient recruiting decisions

NLP and Data Science:

- What are distinguishing characteristics of successful career documents?
- What skills are increasingly important for different industries?

Decision Science:

- How do candidates make decisions about which jobs to apply to?
- How do hiring teams make decisions about candidate qualifications?

Job ▾

🔄 Sync

Similar Jobs ▾

Open

Large Candidate Pool

📈 Applicants: 202 ▾

3850 Characters

Notify ▾

Last edit: **System** ▾

28

Job will perform
poorlyThis job scores **lower than 95%** of **Junior Accounting** jobs in **Los Angeles, CA**

- Add preferred qualifications
- Add more "you" statements
- Perks included
- Equal opportunity statement is included

Neutral

Gendered



Senior Finance Analyst

TapRecruit - Los Angeles

\$76,300 ^{BETA}

\$65,200

\$98,600

TapRecruit is looking for a smart, detail-oriented person to serve as a senior financial analyst. This person will be responsible for supporting the company's FP&A requirements. Responsibilities will include working on TapRecruit Entertainment Group's FP&A model, supporting analysis for long term plan, as options, tracking key business operational metrics and producing monthly financial/operation role will require strong organizational skills to help manage the senior managers across the department and evaluate/implement management. This is a dynamic role that serves the finance department of Finance and will routinely interface with TapRecruit's top management. In addition to FP&A needs, this person will provide discussions with projects for top management and report to a Senior Manager.

Language that emphasizes an "intense" or "confusing" environment is known to deter qualified candidates.

Delete

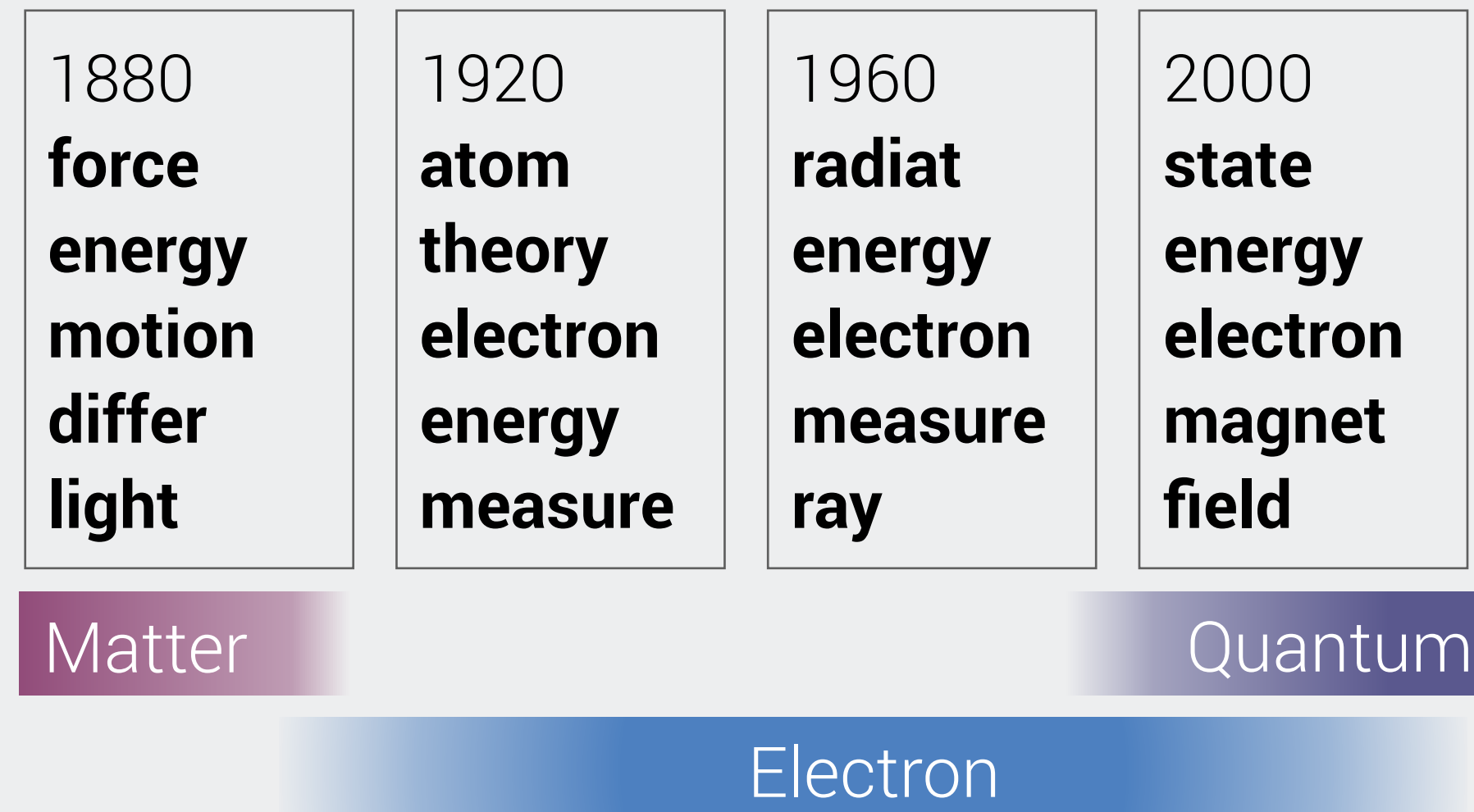
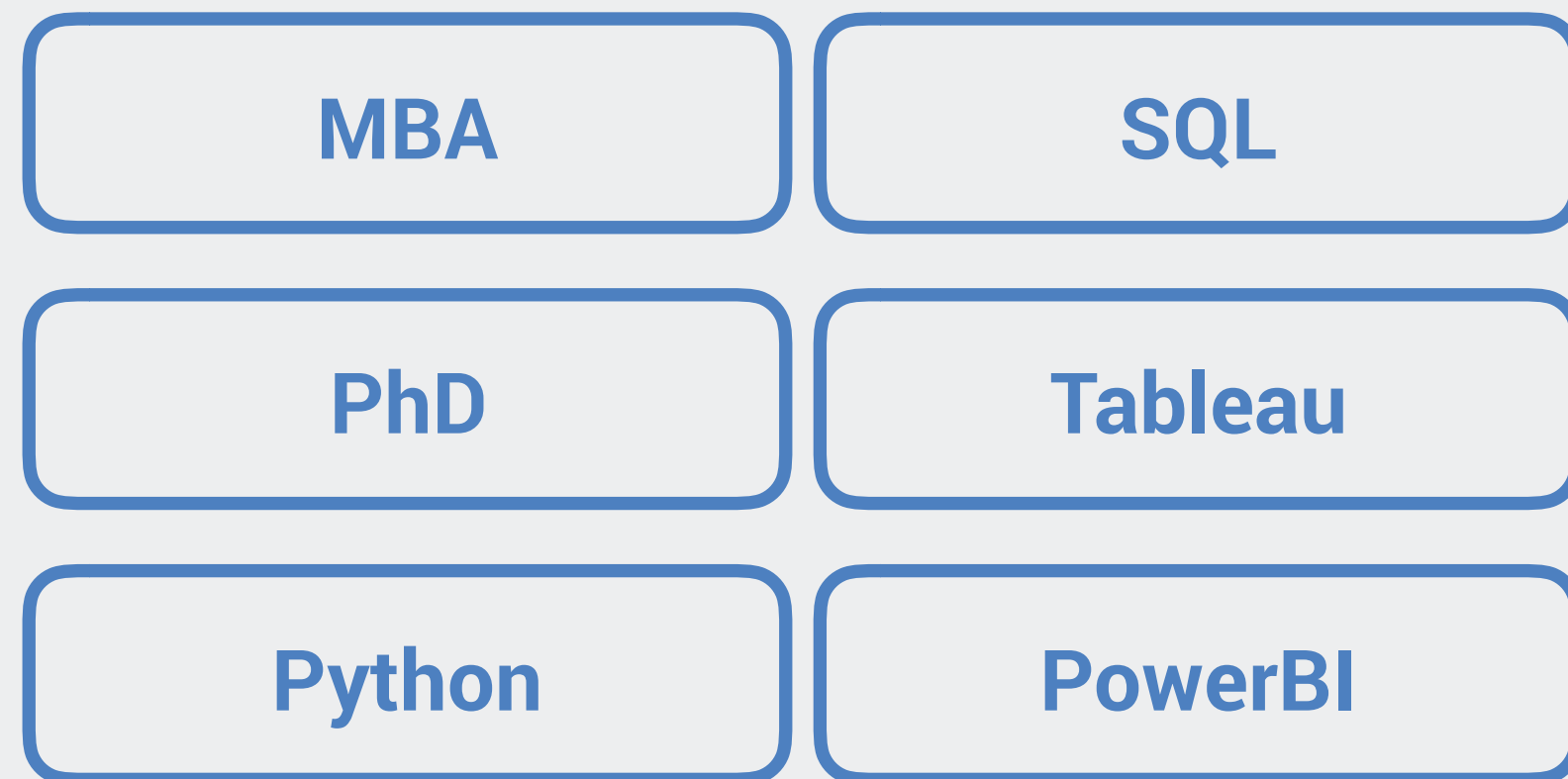
This is an ideal position for an individual who has gained strong experience at an investment bank or accounting firm and now seeks to apply those skills to a fast-growing entrepreneurial company. Strong quantitative and excel financial modeling skills are a must. The ideal candidate must be comfortable in a dynamic start-up environment, will bring energy and passion to everything he/she does, and will not be afraid to roll up his/her sleeves to tackle challenging analytical assignments.

This job is full-time, based in Los Angeles. We offer competitive compensation and stock option program.

**How have data science skills
changed over time?**

Strategies to identify changes in texts

Traditional approaches do not capture syntactic and semantic shifts



Manual Feature Extraction

Require selection of key attributes, therefore difficult to discover new attributes

Dynamic Topic Models

Require experimentation with topic number

Adapted from Blei and Lafferty, ICML 2006.

Embeddings use context to extract meaning

Window sizes capture semantic similarity vs semantic relatedness

Statistical modeling through software (e.g. SPSS) or programming language (e.g. **Python**)

Context

Word

Experience in **Python**, Java or other object-oriented programming languages

Context

Word

Context

Proficiency programming in **Python**, Java or C++.

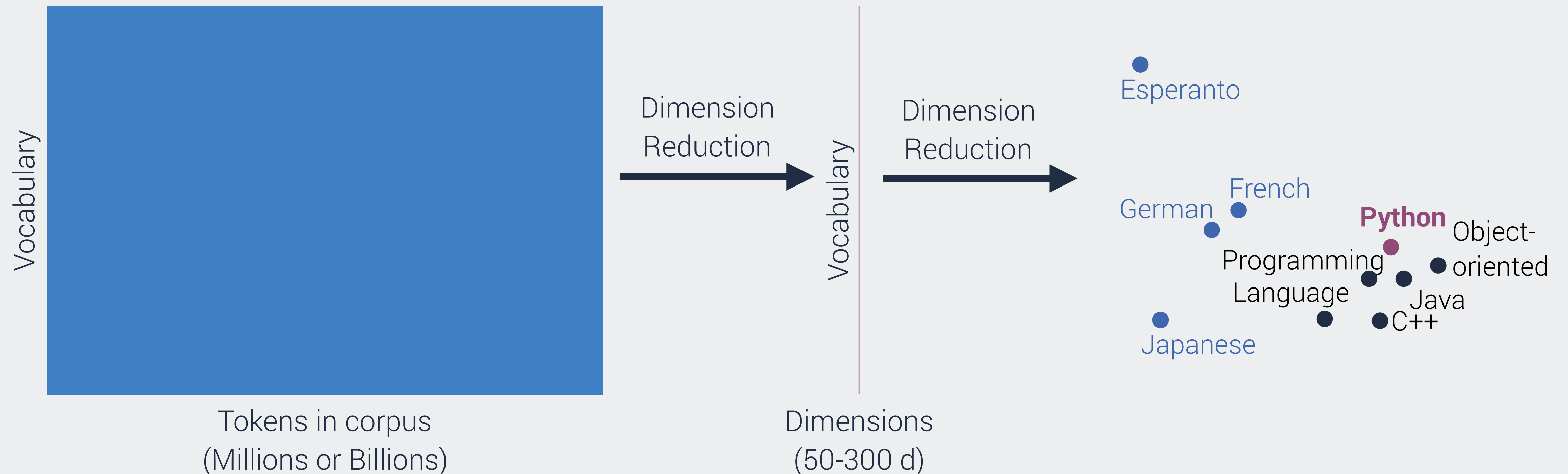
Context

Word

Context

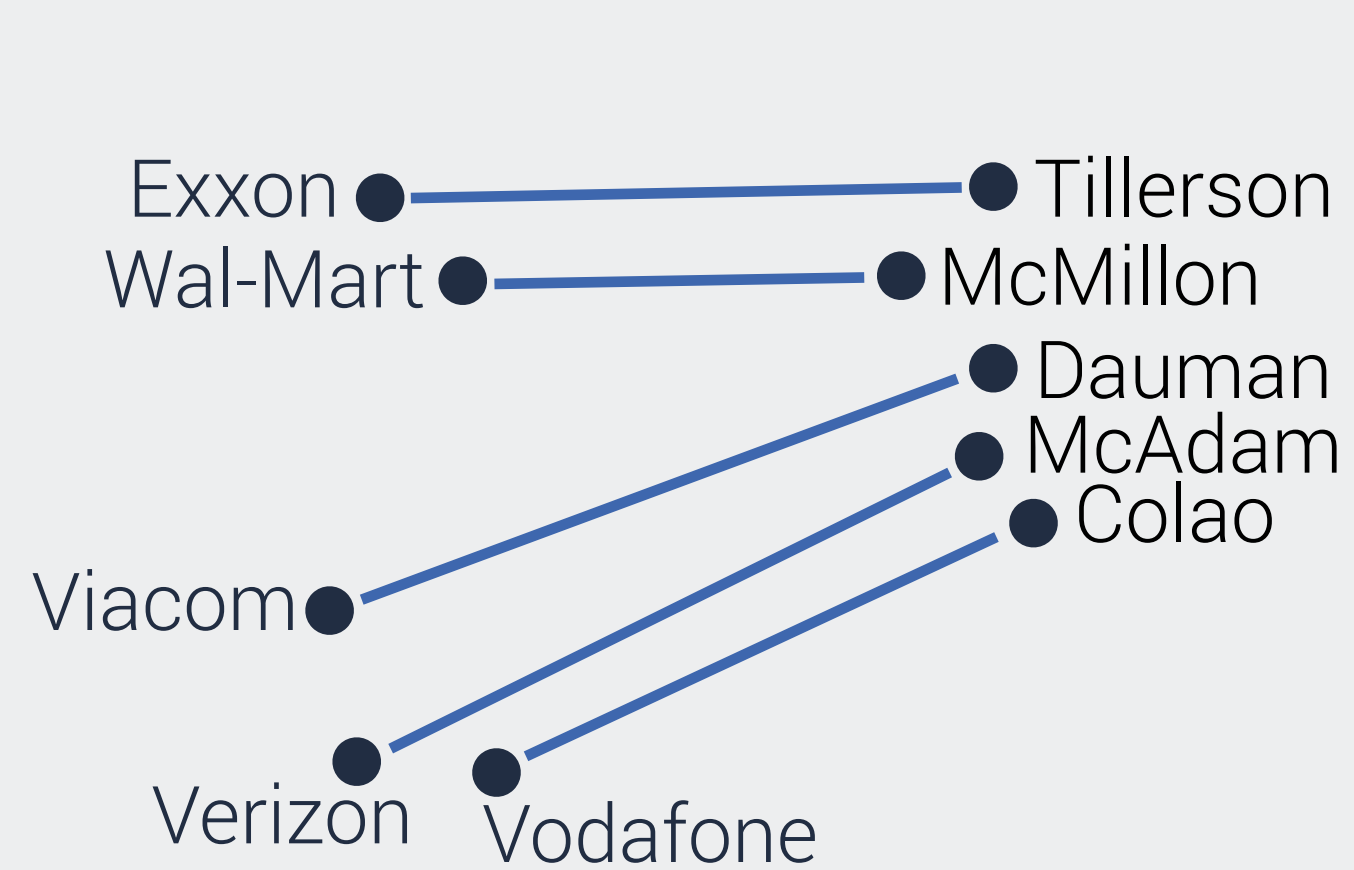
A simplified representation of word vectors

Dimension reduction is key to all types of embeddings models

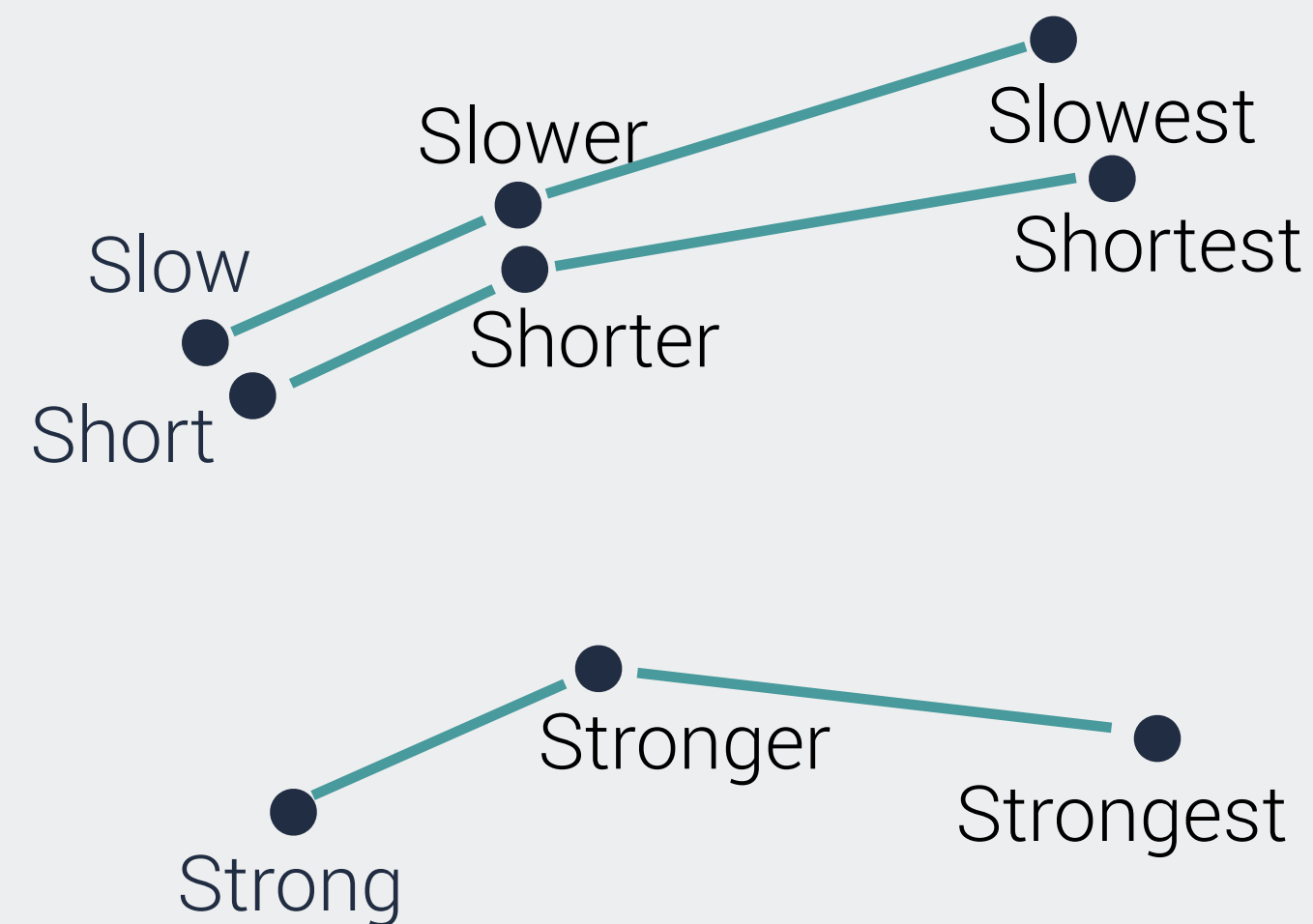


Embeddings capture entity relationships

Dimensionality enables comparison between word pairs along many axes



Hierarchies



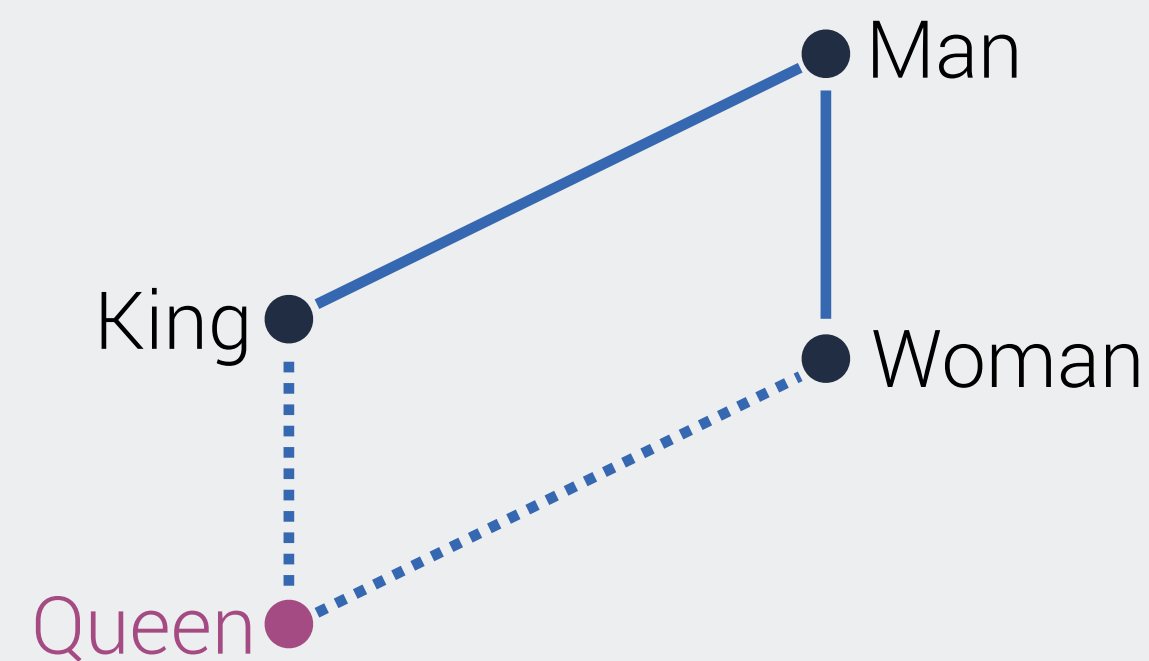
Comparatives and Superlatives



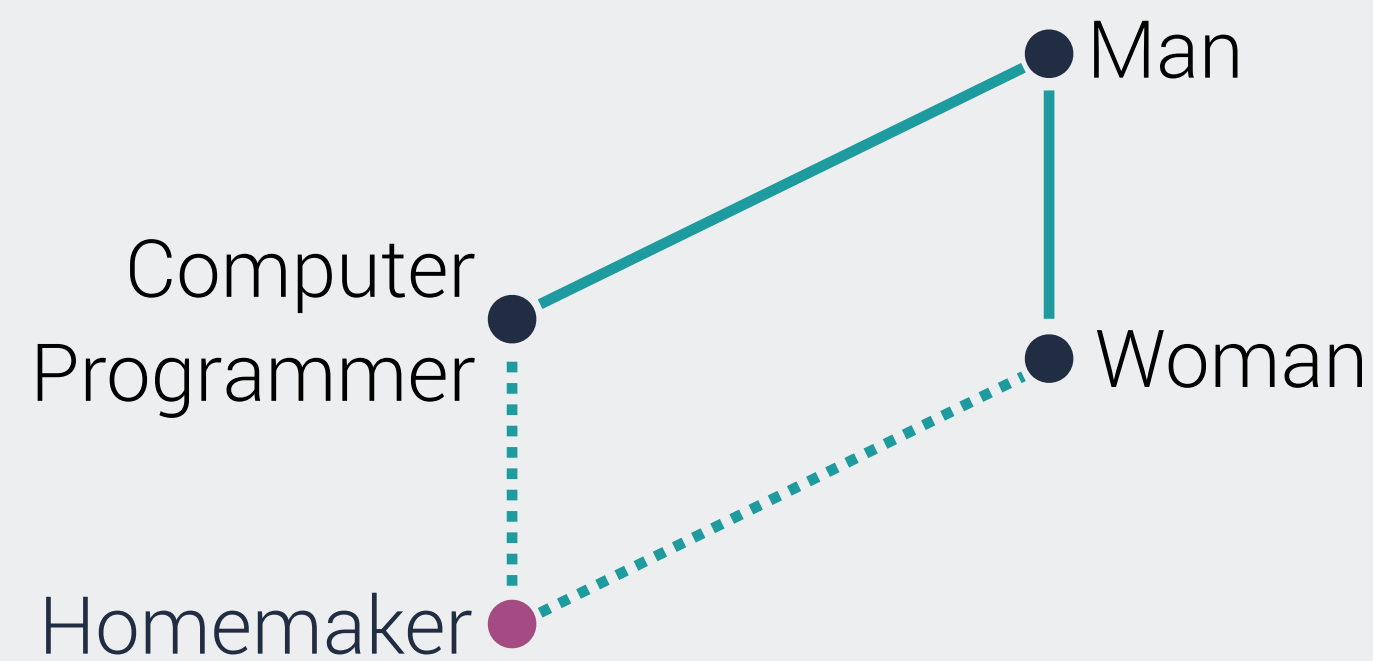
Man :: King as Woman :: ?

Embeddings reflect cultural bias in corpora

High dimensionality enables some bias reduction



Man :: King as Woman :: ?



Man :: Programmer as Woman :: ?

Adapted from Bolukbasi et al., [arXiv: 1607:06520](https://arxiv.org/abs/1607.06520).

Pretrained embeddings facilitate fast prototyping

Embeddings training should match corpus that is being tested on

| | | | | | |
|---------------------------|----------------------------|---------------------|--------------------------|---------------------|---------------------|
| Corpus Generation | Corpus Tokens | Twitter 27 B | Common Crawl 42-840 B | GoogleNews 100 B | Wikipedia 6 B |
| Corpus Processing | Vocabulary Size | 1.2 M | 1.9-2.2 M | 3 M | 400 k |
| Language Model Generation | Algorithm Vector Length | GLoVE 25 - 200 d | GLoVE 300 d | word2vec 300 d | GLoVE 50 - 300 d |
| Language Model Tuning | | | | | |
| Final Application | | | | | |

Problems with pretrained embedding models

| | |
|--------------------------------|---|
| Casing | Abbreviations vs Words e.g. IT vs it |
| Out of Vocabulary Words | Domain Specific Words & Acronyms |
| Polysemy | Words with multiple meanings e.g. drive (a car) vs drive (results) e.g. Chef (the job) vs Chef (the language) |
| Multi-word Expressions | Phrases that have new meanings e.g. Front-end vs front + end |

Custom language models tools

Modularized for different data and modeling requirements

spaCy

OPEN NLP™

gensim

PYTORCH

CoreNLP

SyntaxNet

TensorFlow



Amazon SageMaker

Corpus Processing

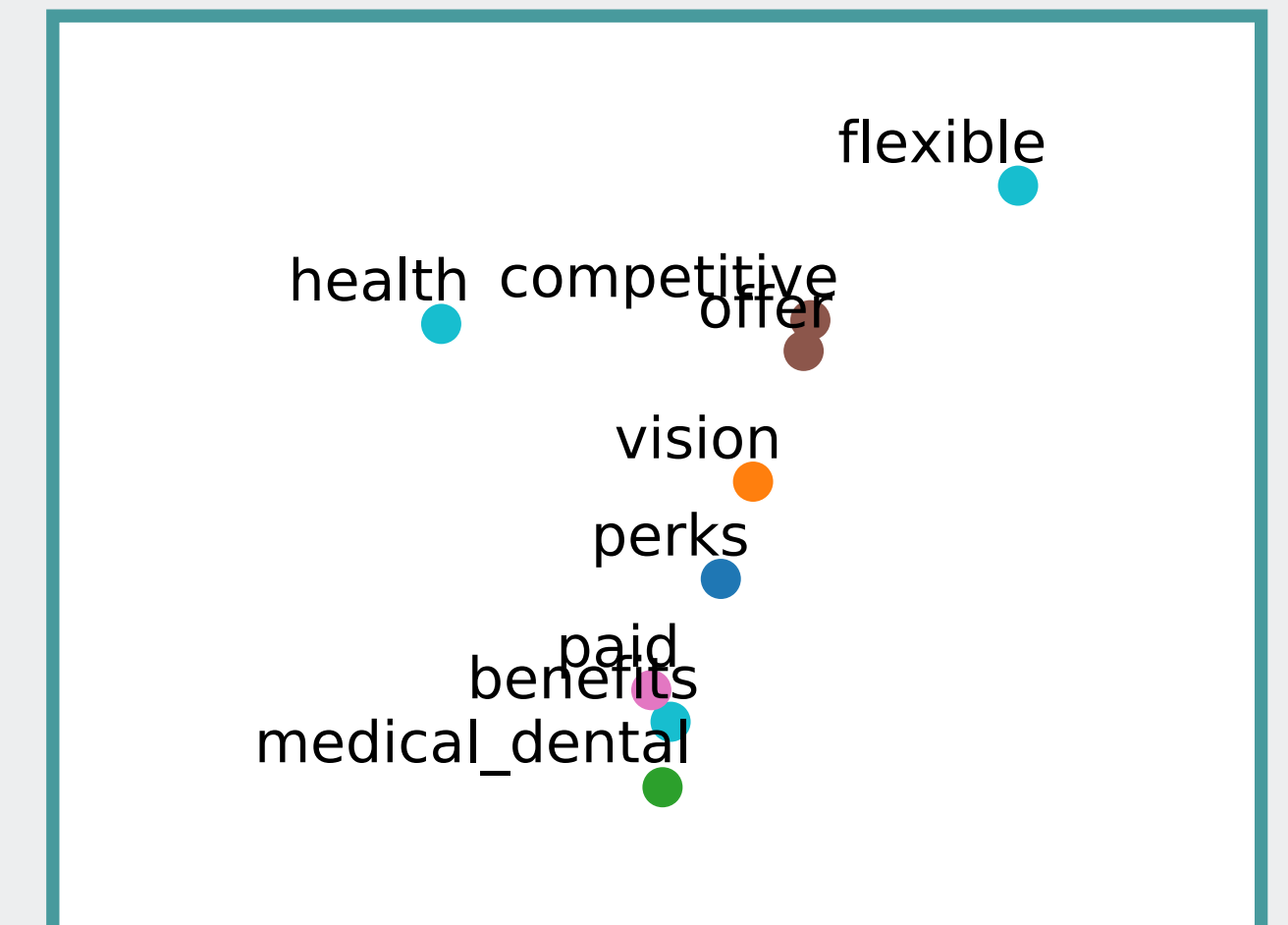
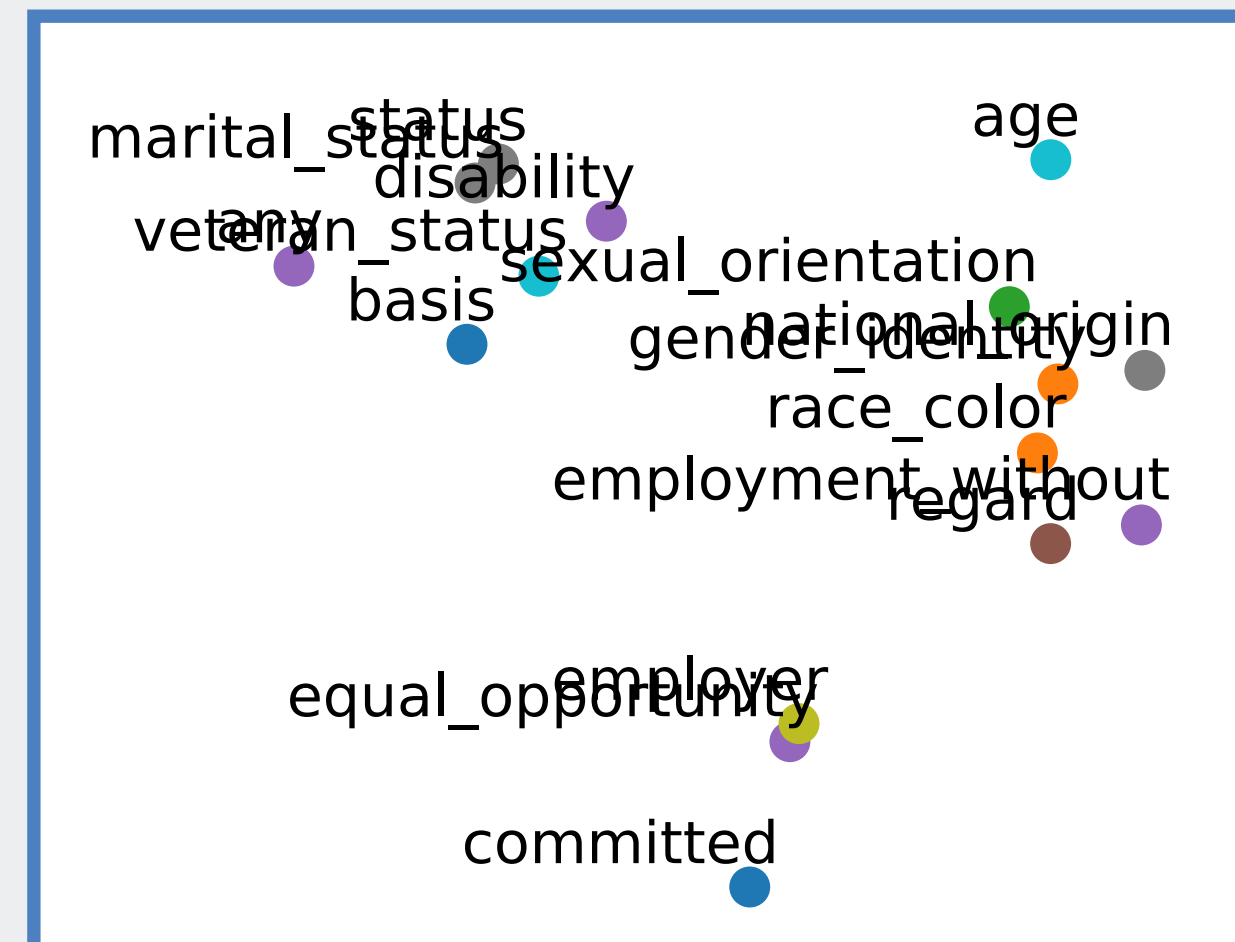
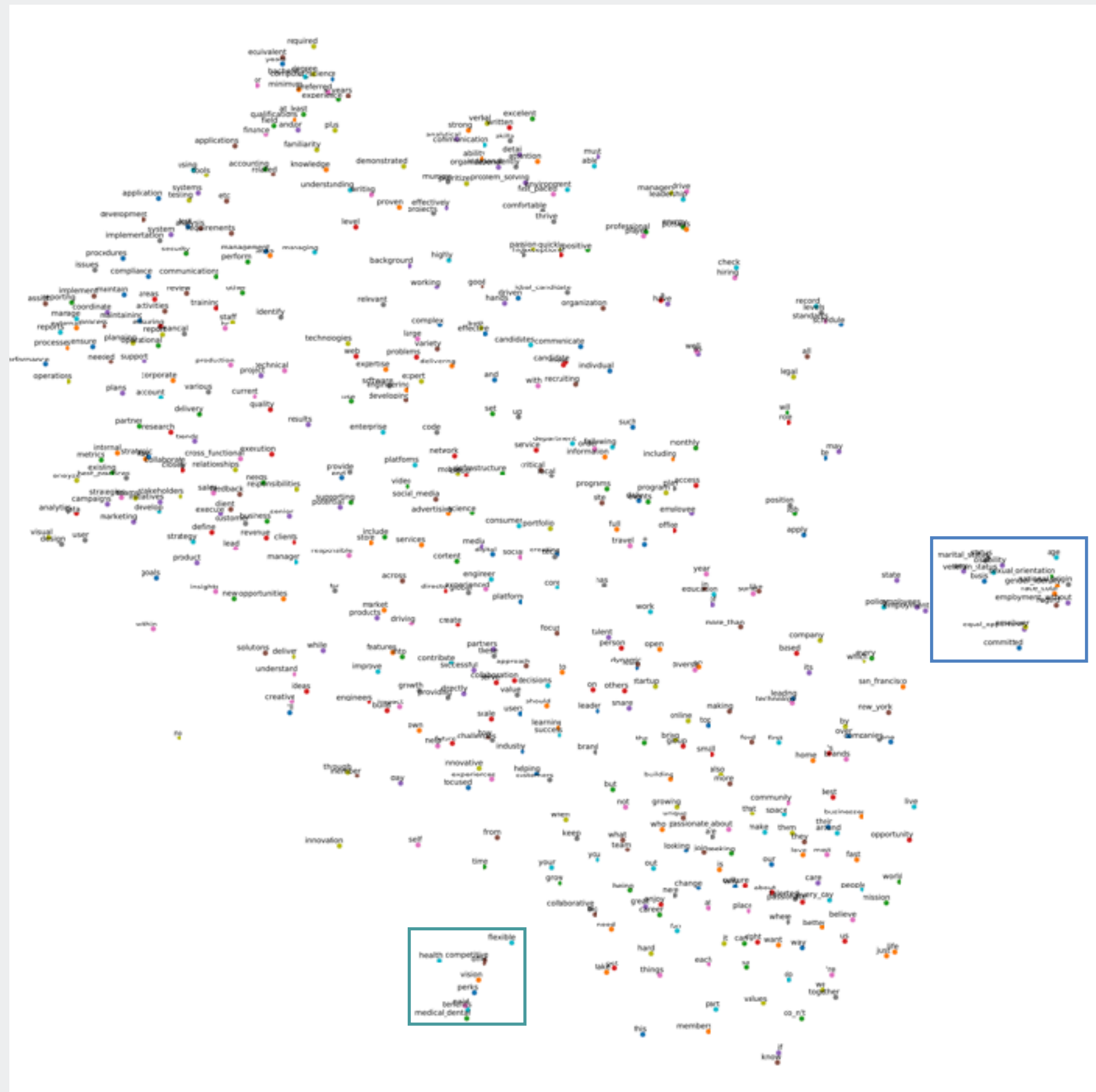
Tokenization, POS tagging, Sentence
Segmentation, Dependency Parsing

Language Modeling

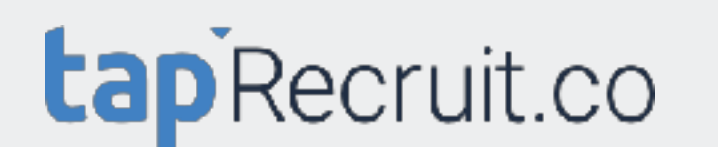
Different word embedding models
(GLoVE, word2vec, fastText)

Career language embedding model

Identified equal opportunity and perks language



Identified 'soft' skills and language around experience

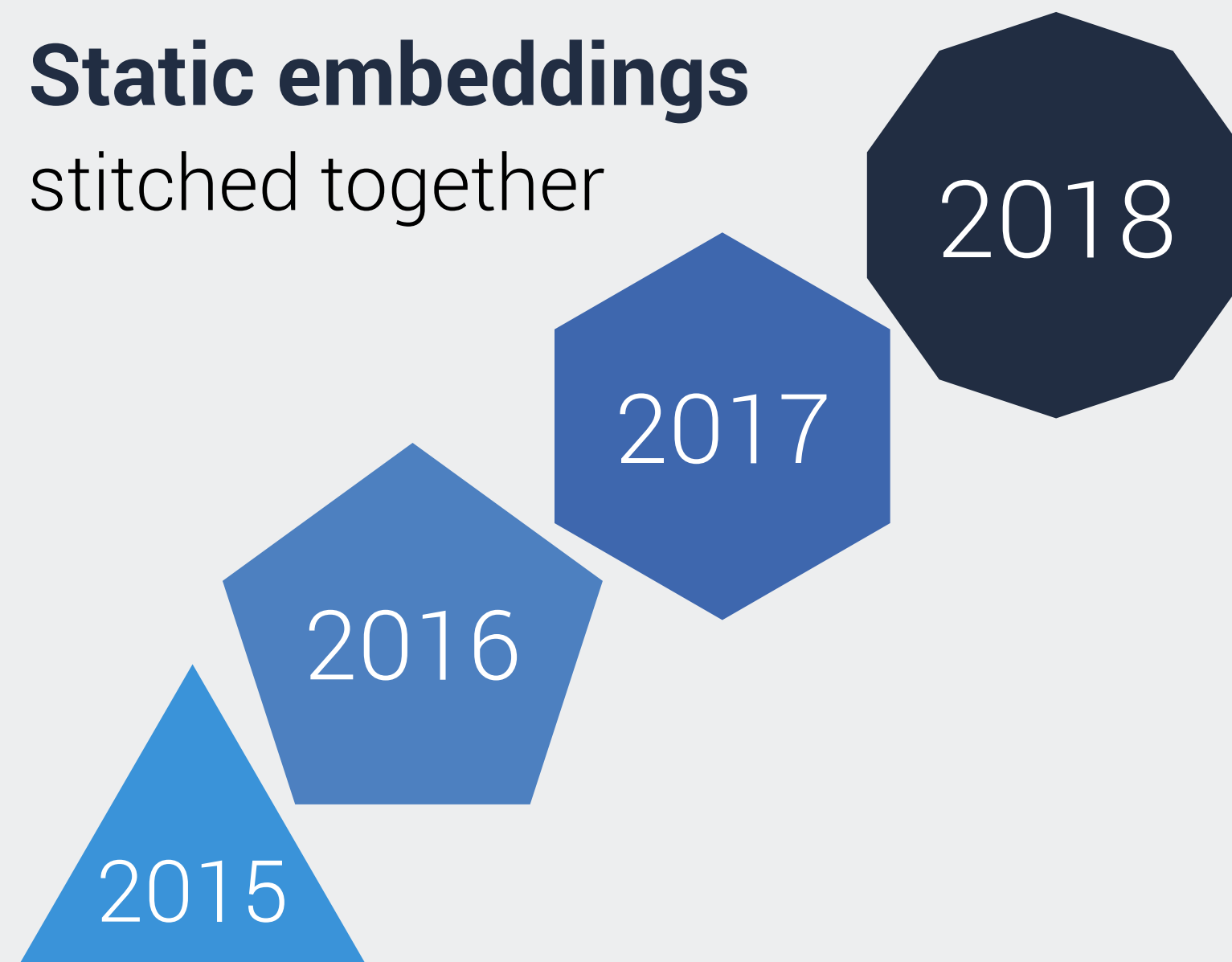


I've got 300 dimensions...
but time ain't one

Two approaches to connect embeddings

Static embeddings

stitched together

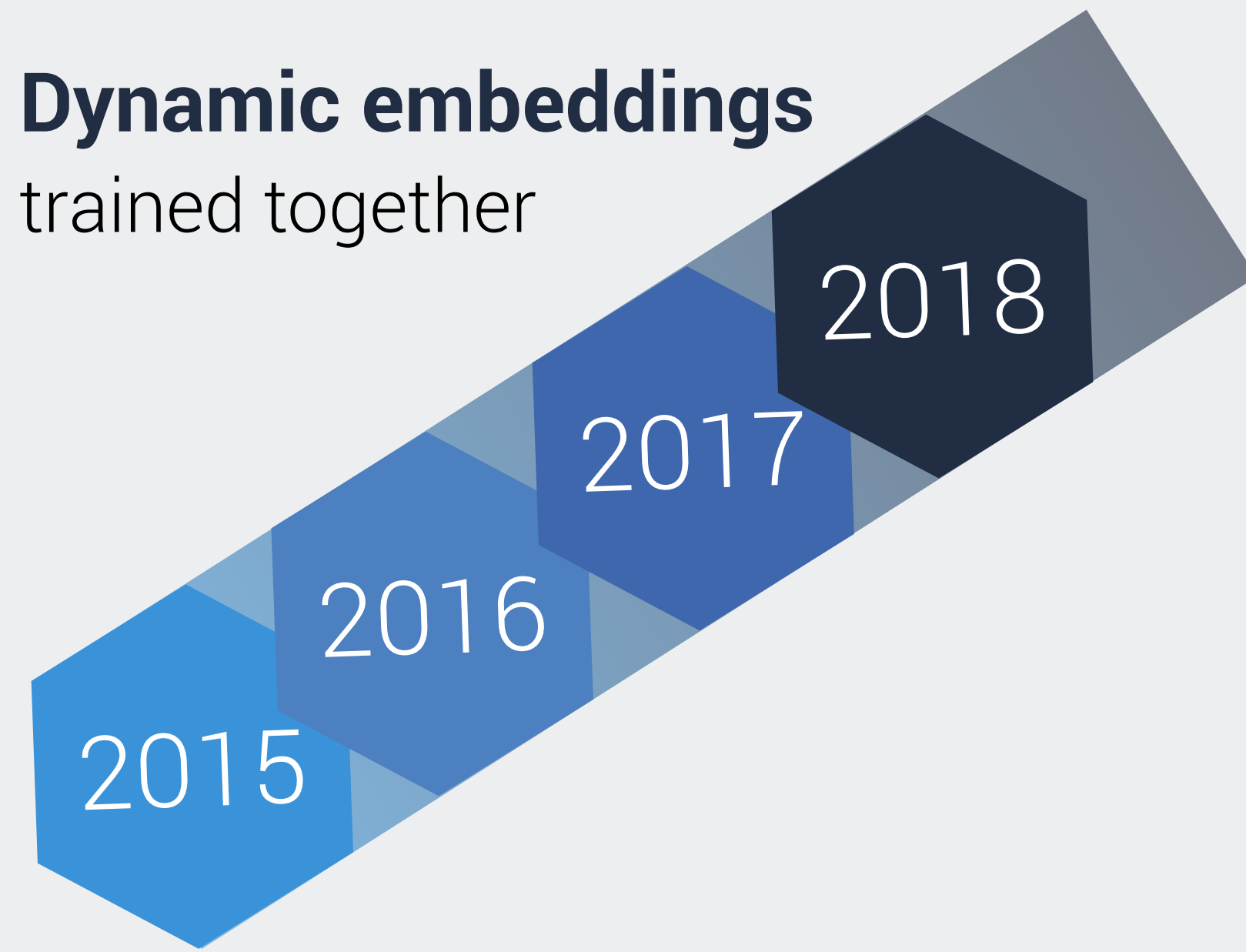


Data hungry
Requires alignment

Kim, Chiu, Kaneki, Hedge and Petrov, [arXiv: 1405:3515](#).
Kulkarni, Al-Rfou, Perozzi and Skiena, [arXiv: 1411:3315](#).

Dynamic embeddings

trained together



Data efficient
Does not require alignment

Balmer and Mandt, [arXiv: 1702:08359](#)
Yao, Sun, Ding, Rao and Xiong, [arXiv: 1703:00607](#)
Rudolph and Blei, [arXiv: 1703:08052](#)

Dynamic Bernoulli embeddings

Outputs facilitate quick analysis of trends

Absolute drift

Identifies top words whose usage changes over time course

| words with largest drift (Senate) | | | |
|-----------------------------------|------|-----------------|------|
| IRAQ | 3.09 | coin | 2.39 |
| tax cuts | 2.84 | social security | 2.38 |
| health care | 2.62 | FINE | 2.38 |
| energy | 2.55 | signal | 2.38 |
| medicare | 2.55 | program | 2.36 |
| DISCIPLINE | 2.44 | moves | 2.35 |
| text | 2.41 | credit | 2.34 |
| VALUES | 2.40 | UNEMPLOYMENT | 2.34 |

Embedding neighborhoods

Extract semantic changes by nearest neighbors of drifting words

| UNEMPLOYMENT | | |
|--------------|--------------|--------------|
| 1858 | 1940 | 2000 |
| unemployment | unemployment | unemployment |
| unemployed | unemployed | jobless |
| depression | depression | rate |
| acute | alleviating | depression |
| deplorable | destitution | forecasts |
| alleviating | acute | crate |
| destitution | reemployment | upward |
| urban | deplorable | lag |
| employment | employment | economists |
| distressing | distress | predict |

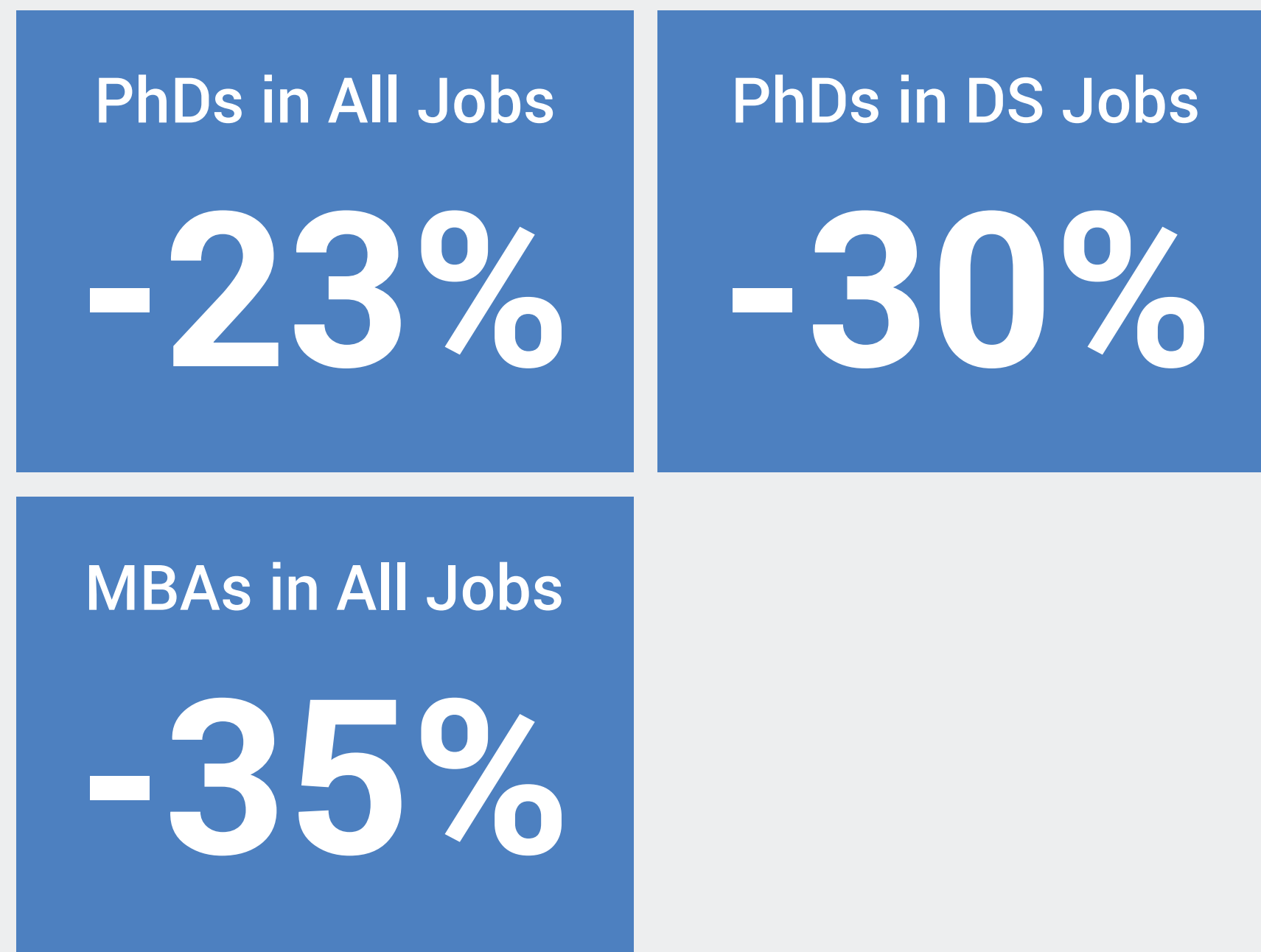
Experiments with dynamic embeddings

| | Small Corpus | Large Corpus |
|-----------------------------|------------------|------------------|
| Job Types | All | All |
| Time Slices | 3 (2016-2018) | 3 (2016-2018) |
| Number of Documents | 50 k | 500 k |
| Vocabulary Size | 10 k | 10 k |
| Data Preprocessing | Basic | Basic |
| Embedding Dimensions | 100 d | 100 d |

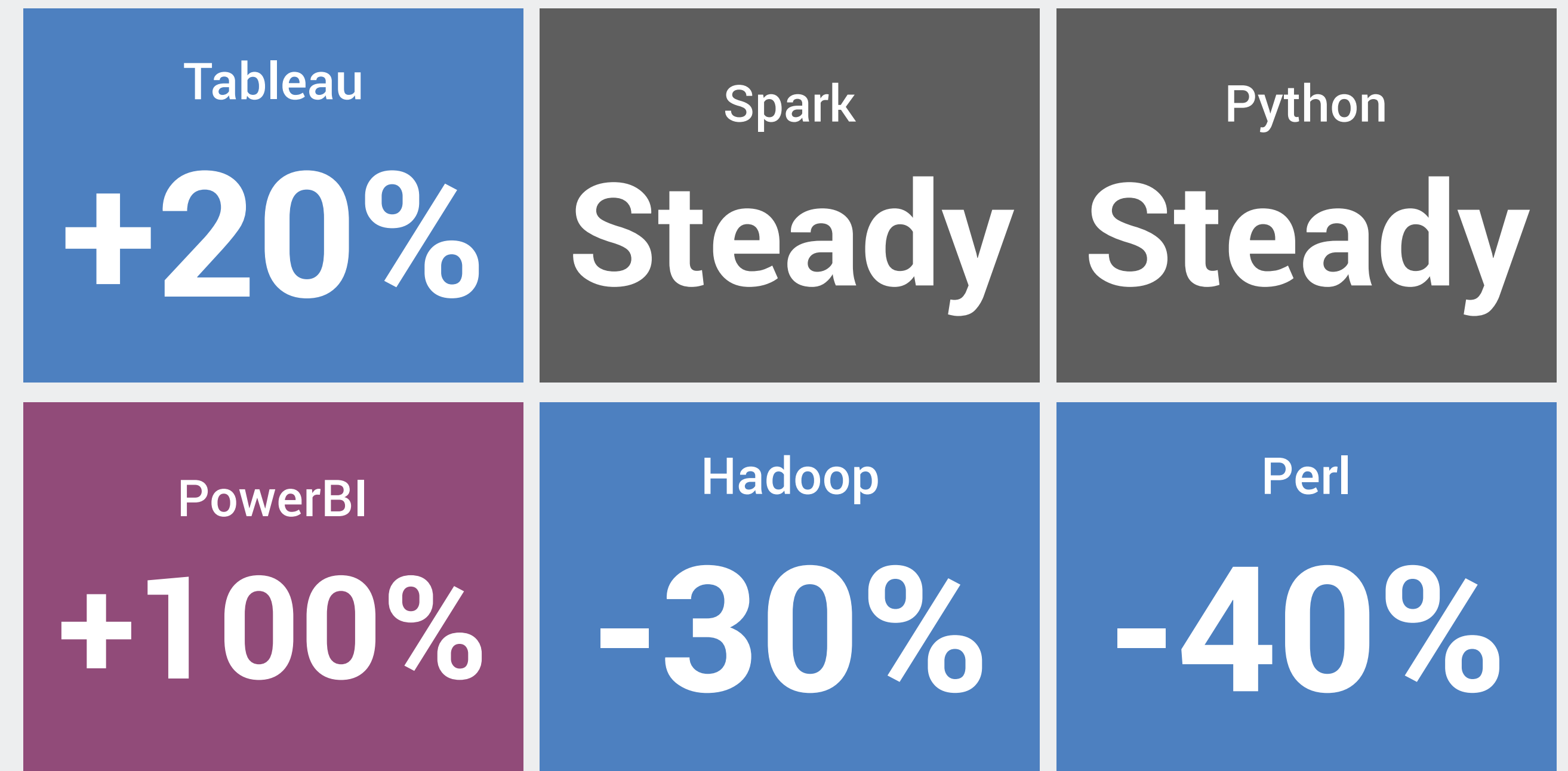
Dynamic embeddings

Small corpus identified gains and losses

Demand for PhDs and MBAs is Falling



Data Science skills showing significant shifts

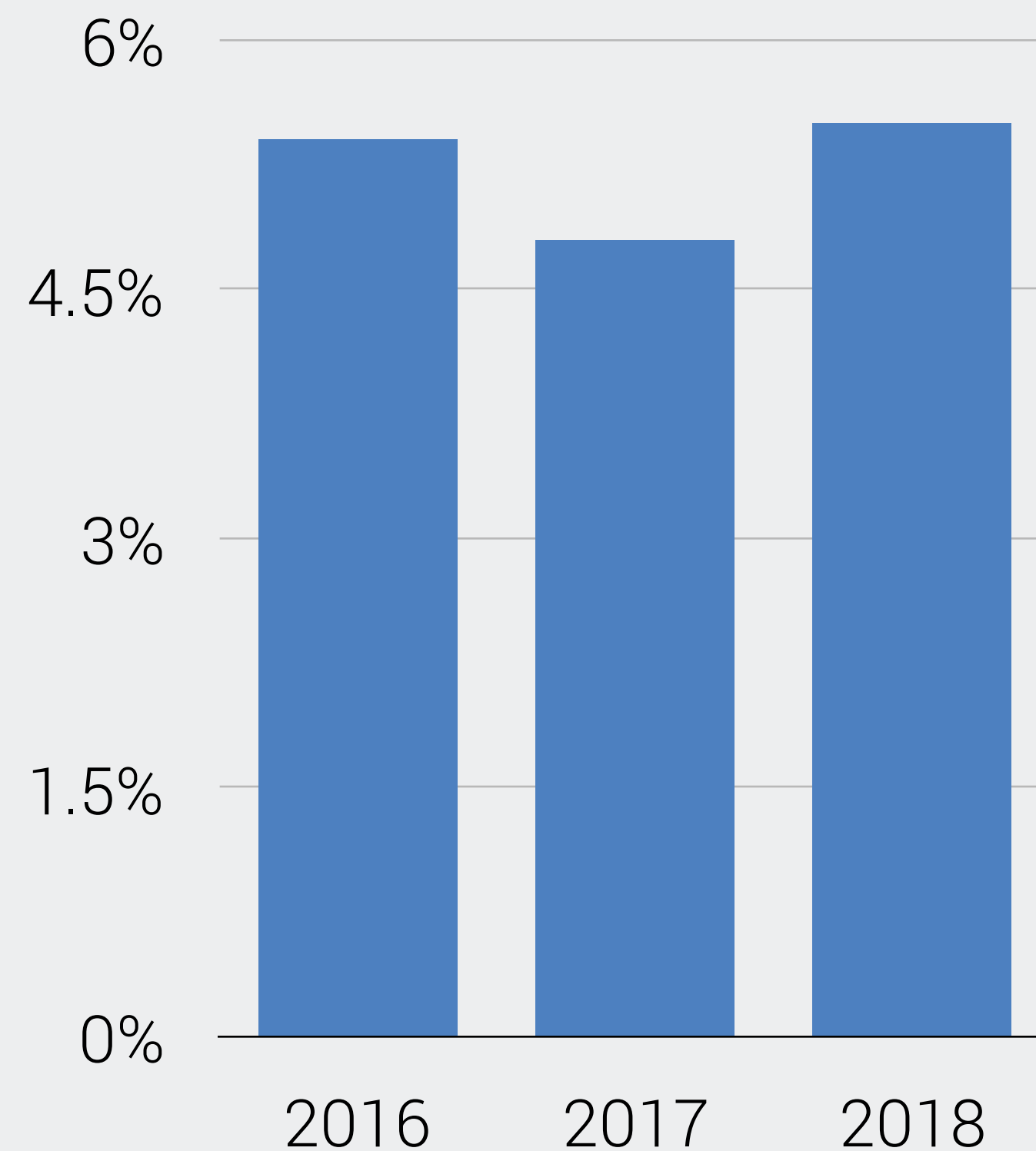


Blue boxes indicate phrases identified from top drifting words analysis.
Grey boxes indicate 'control' skills.

Dynamic embeddings

Large corpus identified role-type dependent shifts in requirements

No change to SQL demand



SQL requirement increases in specific functions



Blue boxes indicate phrases identified from top drifting words analysis.
Grey boxes indicate 'control' skills.

How have data science skills changed over time?

- Flavors of static word embeddings: The Corpus Issue
- Considerations for developing custom embedding models
- Flavors of dynamic models: Dynamic Bernoulli embeddings

Thank you Women in Analytics!

Maryam Jahanshahi Ph.D.

Research Scientist

 @mjahanshahi

 maryam-j

tapRecruit.co

<http://bit.ly/wia-2019>