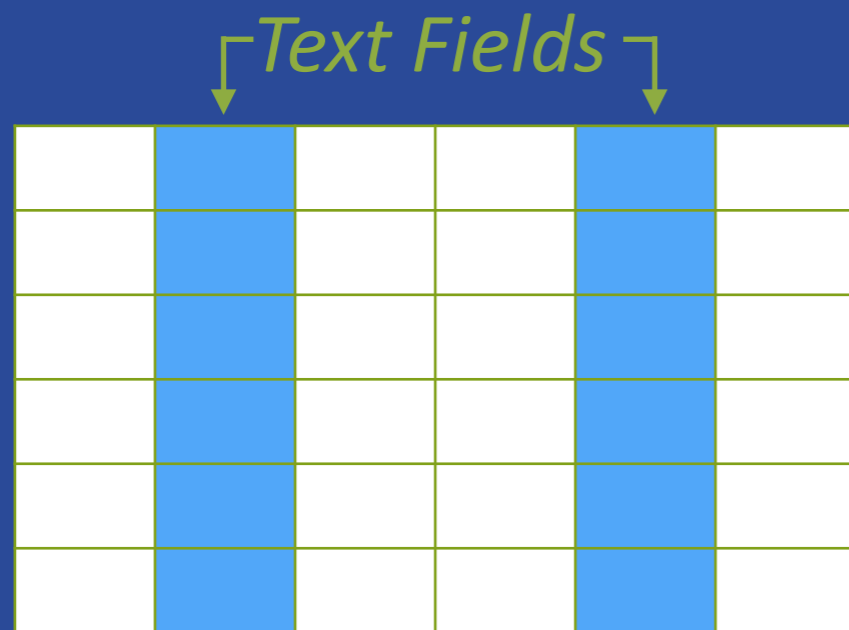# Topic Models
## Discovering **Thematic Meaning** in Text

**Charles Parker**
**VP ML Algorithms, BigML, Inc**

# What is Topic Modeling?

*Text Fields*

- Unsupervised algorithm

- Learns only from **text fields**

- Finds hidden **topics** that **model** the text

**Questions:**

- How is this different from the **Text Analysis** that BigML already offers?

- What does it output and how do we use it

- Unsupervised… model?

# Text Analysis

1. Stem Words -> Tokens

2. Remove tokens that occur too often

3. Remove tokens that do not occur often enough

4. Count occurrences of remaining "interesting" tokens

Be not **afraid** of **great**ness:
some are **born great,** some
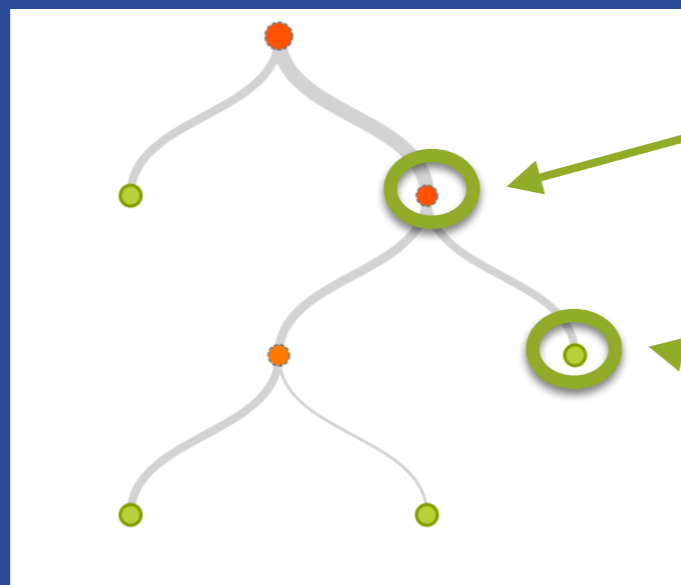**achieve great**ness, and
some have **great**ness
thrust upon 'em.

*great: appears 4 times*

Be not afraid of greatness:
some are born great, some achieve
greatness, and some have greatness
thrust upon 'em.

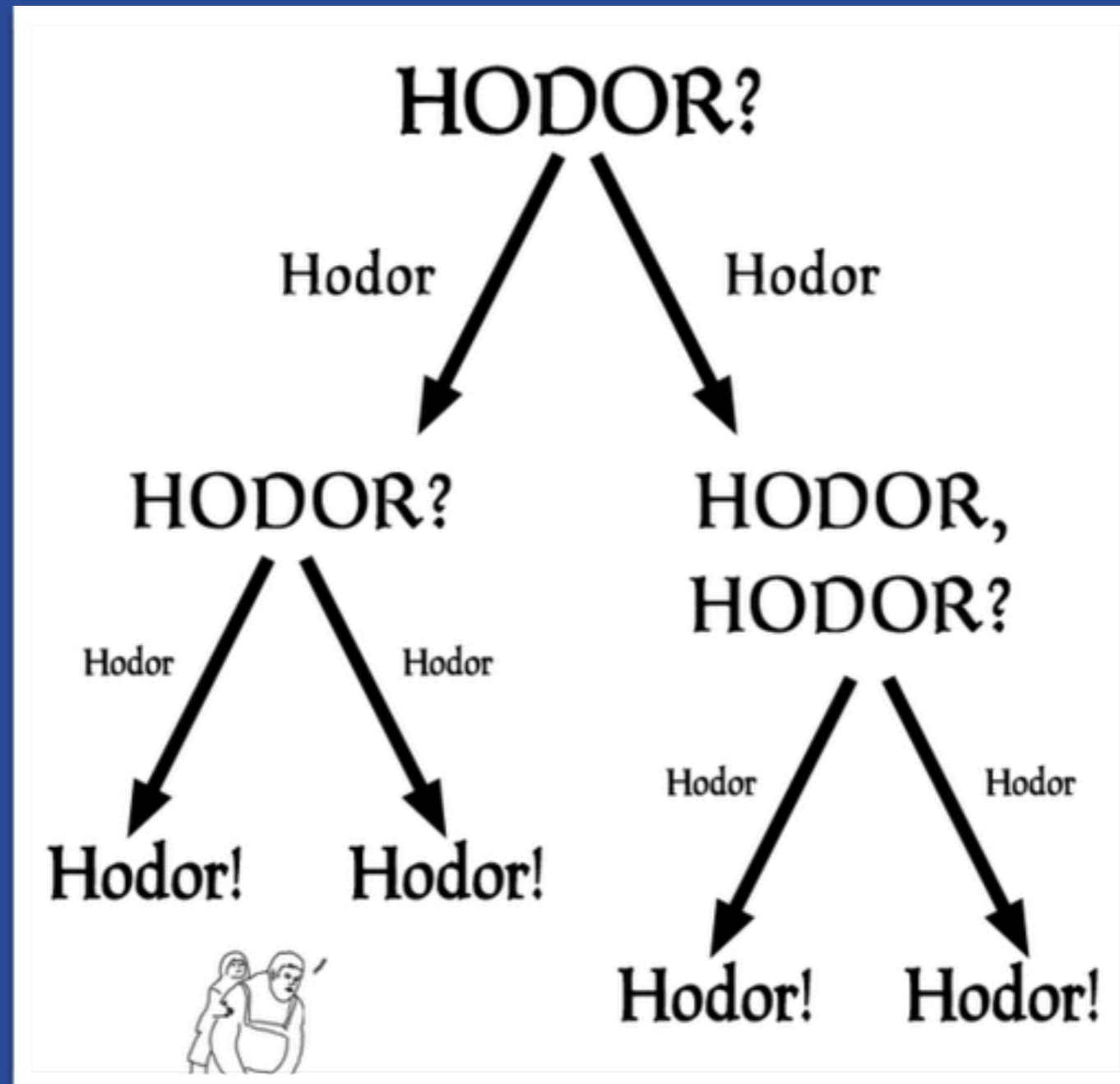| ... | great | afraid | born | achieve | ... | ... |
|---|---|---|---|---|---|---|
| ... | 4 | 1 | 1 | 1 | ... | ... |
| ... | ... | ... | ... | ... | ... | ... |

Model

The token "great" occurs more than 3 times

The token "afraid" occurs no more than once

# Topic Model Demo #1

## Text Analysis

**Creates thousands of hidden token counts**

**Token counts are independently uninteresting**

**No semantic importance**

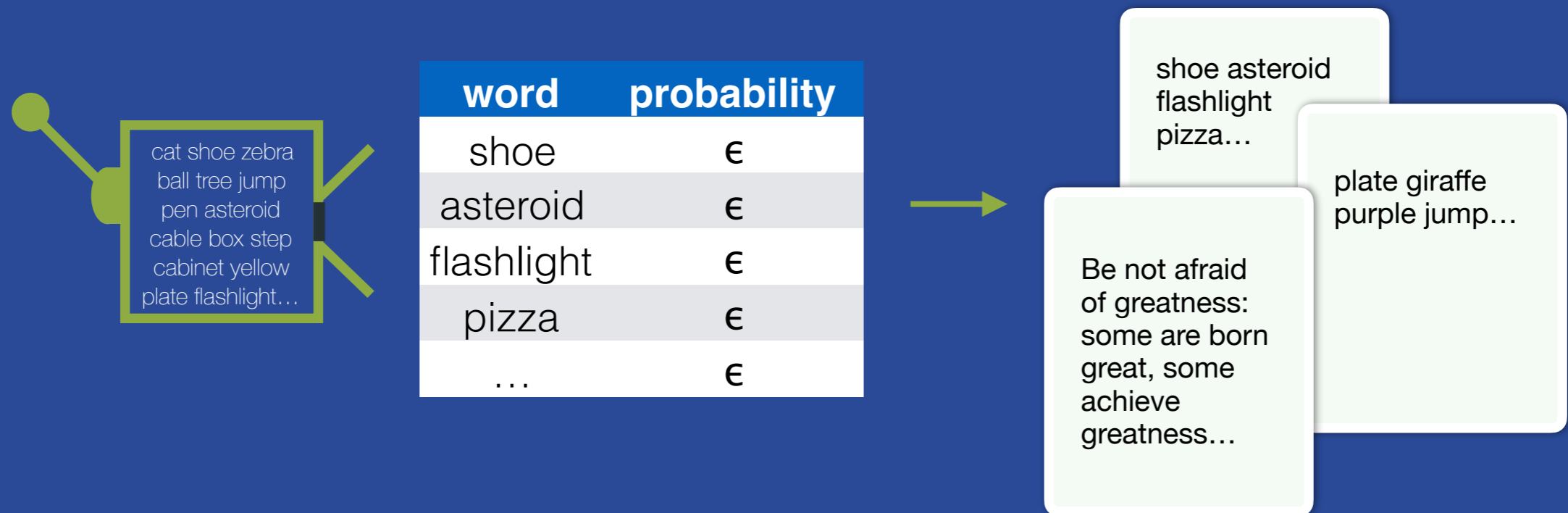**No measure of co-occurrence**

## Topic Model

**Creates tens of topics that model the text**

**Topics are independently interesting**

**Semantic meaning extracted**

**Support for bigrams**

# Generating Documents

cat shoe zebra
ball tree jump
pen asteroid
cable box step
cabinet yellow
plate flashlight…

| word | probability |
|------|-------------|
| shoe | $\epsilon$ |
| asteroid | $\epsilon$ |
| flashlight | $\epsilon$ |
| pizza | $\epsilon$ |
| … | $\epsilon$ |

shoe asteroid flashlight pizza…

plate giraffe purple jump…
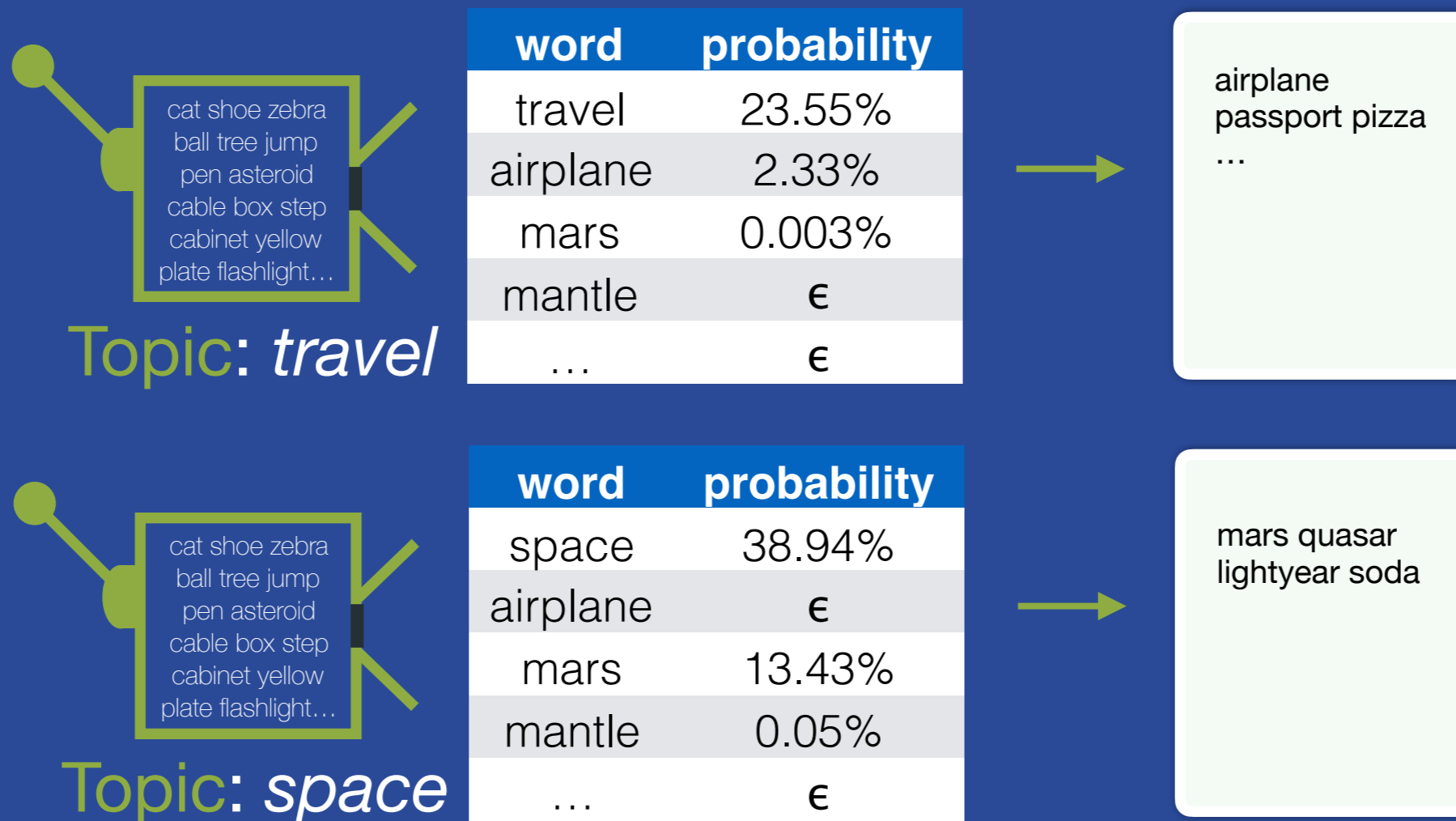
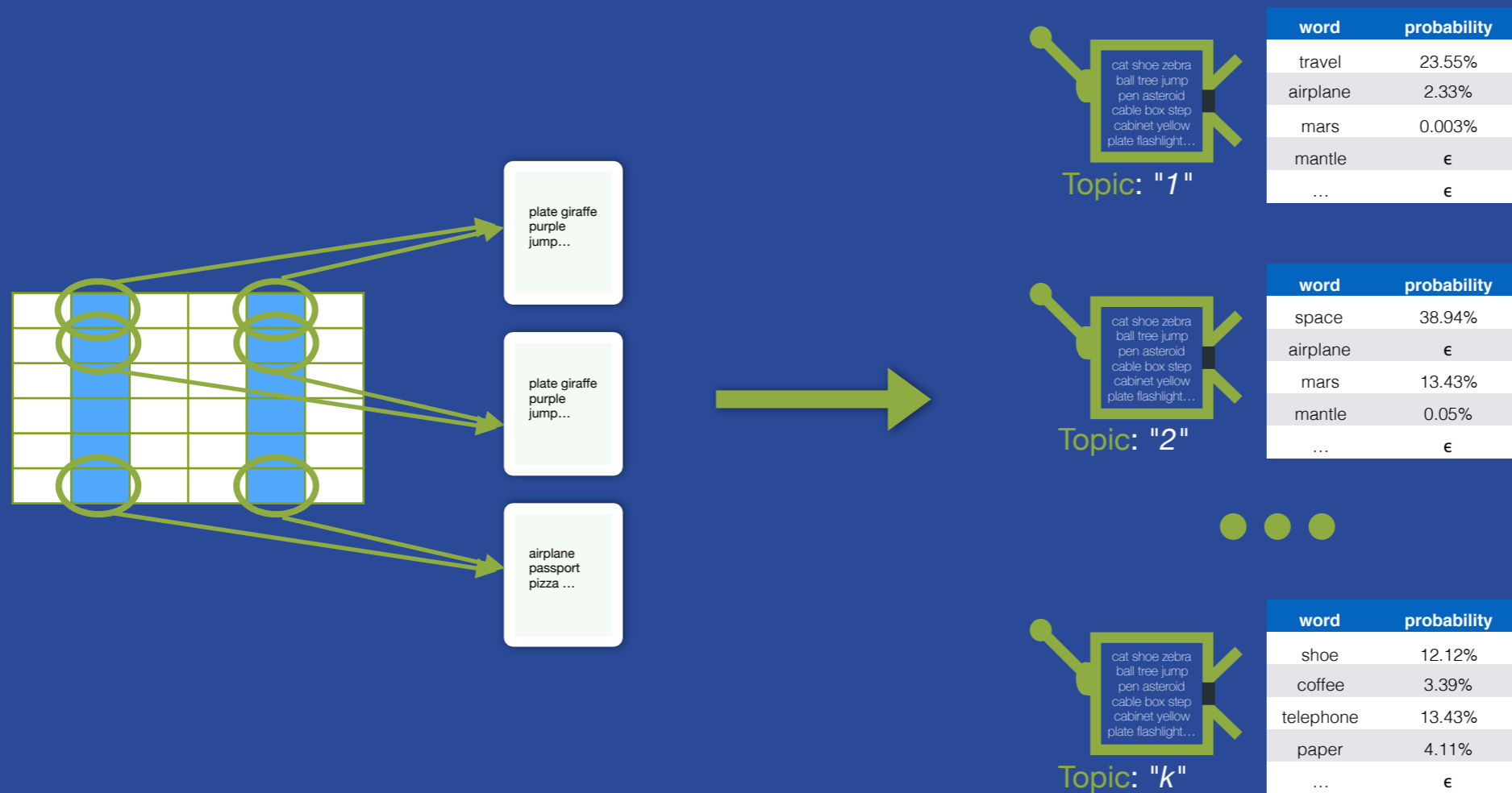Be not afraid of greatness: some are born great, some achieve greatness…

- "Machine" that generates a random word with equal probability with each pull.

- Pull random number of times to generate a document.

- All documents can be generated, but most are nonsense.

# Topic Model

**Intuition:**

- Written documents have *meaning* - one way to describe meaning is to assign a topic.

- For our random machine, the topic can be thought of as increasing the probability of certain words.



cat shoe zebra ball tree jump pen asteroid cable box step cabinet yellow plate flashlight…

**Topic:** *travel*

| word | probability |
| --- | --- |
| travel | 23.55% |
| airplane | 2.33% |
| mars | 0.003% |
| mantle | ε |
| … | ε |

airplane
passport pizza
…

cat shoe zebra ball tree jump pen asteroid cable box step cabinet yellow plate flashlight…

**Topic:** *space*

| word | probability |
| --- | --- |
| space | 38.94% |
| airplane | ε |
| mars | 13.43% |
| mantle | 0.05% |
| … | ε |

mars quasar
lightyear soda

# Topic Model



- Each text field in a row is concatenated into a document

- The documents are analyzed to generate "k" related topics

- Each topic is represented by a distribution of term probabilities

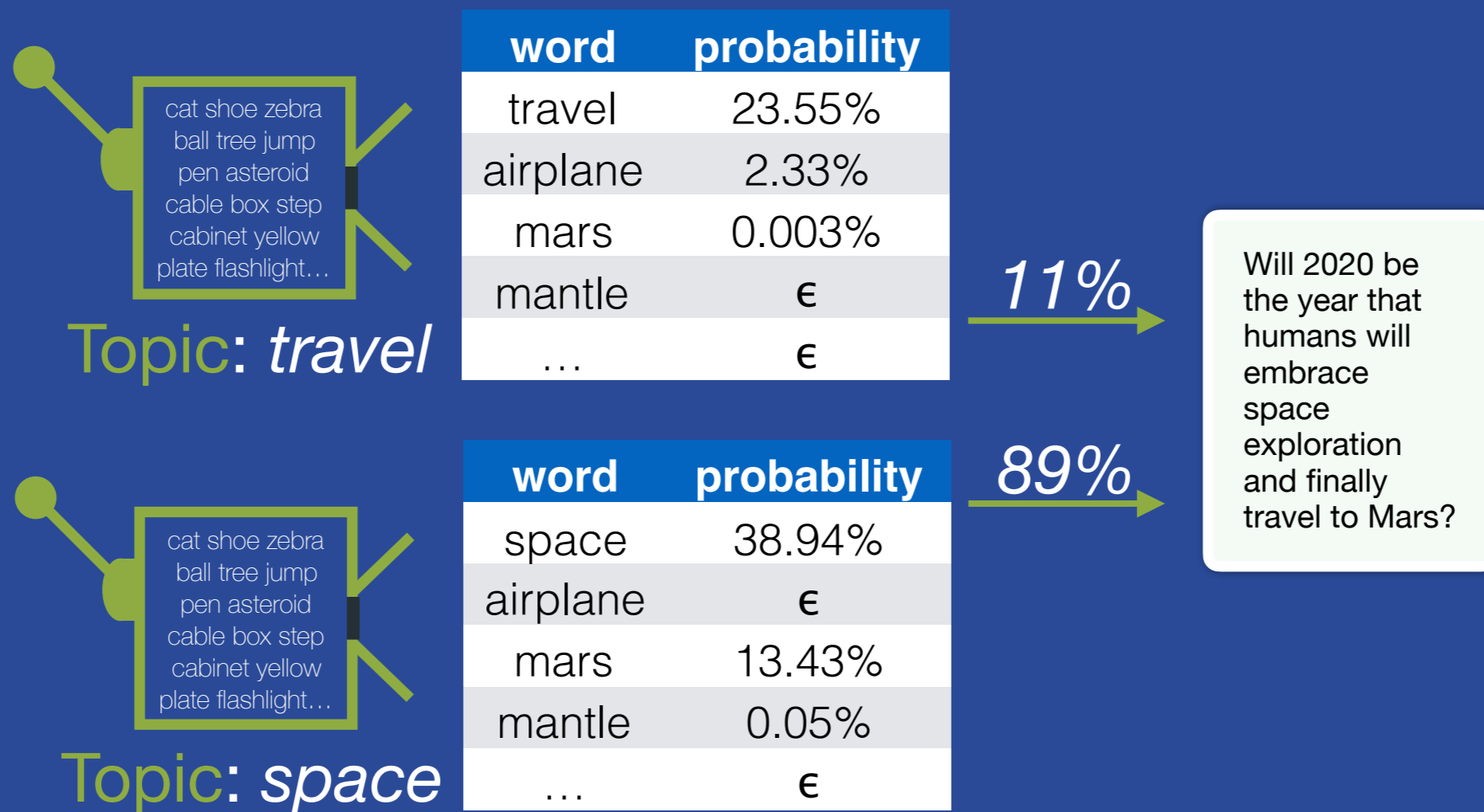# Topic Model Demo #2

# Topic Distribution

**Intuition:**

- Any given document is likely a mixture of the modeled topics…

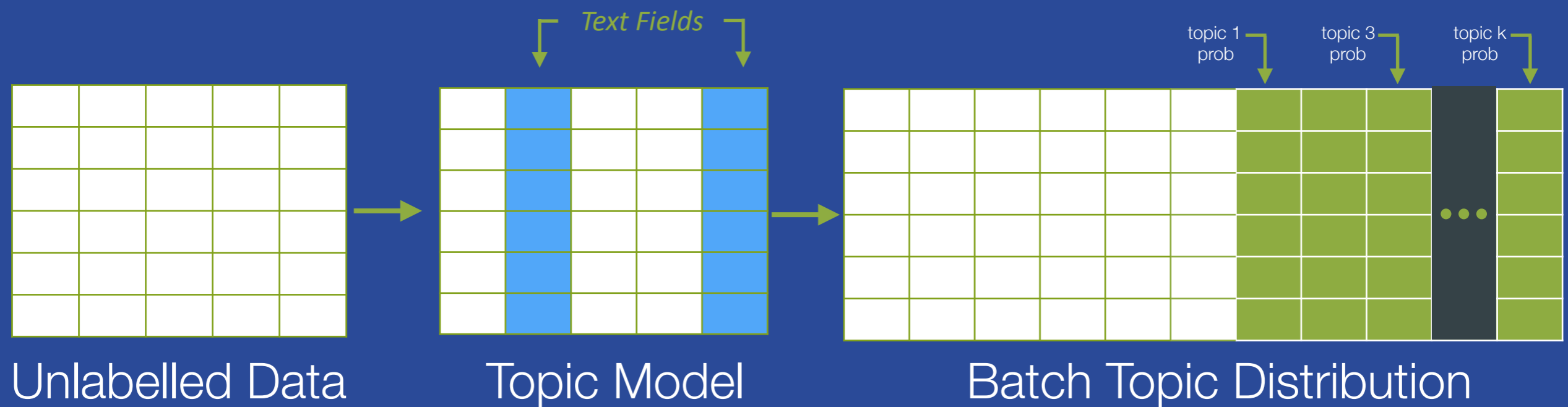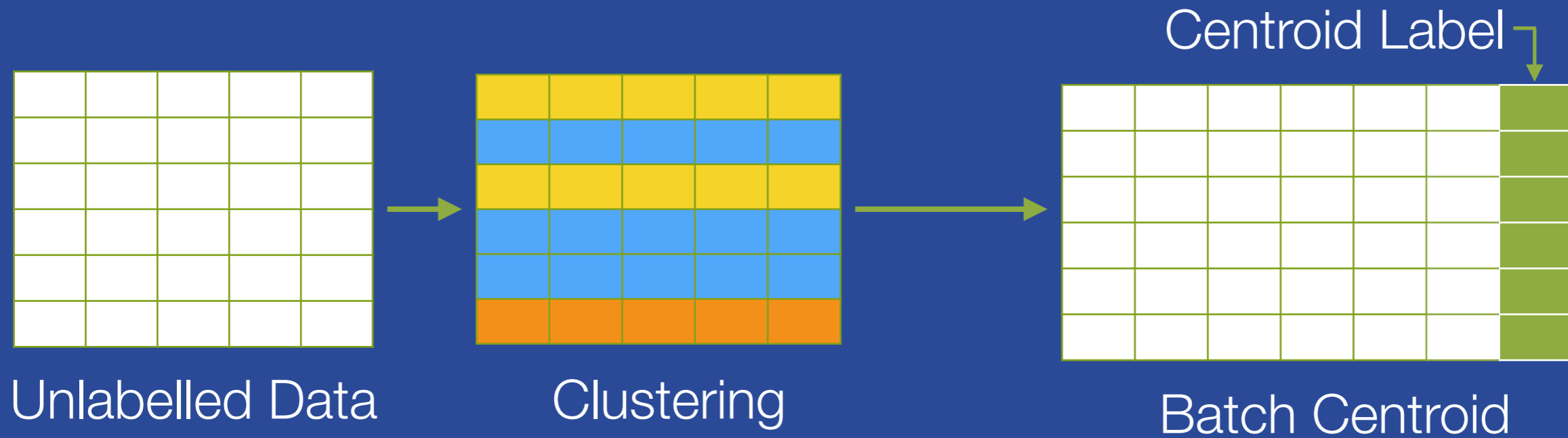- This can be represented as a distribution of topic probabilities

cat shoe zebra ball tree jump pen asteroid cable box step cabinet yellow plate flashlight…

**Topic:** *travel*

| word | probability |
|------|-------------|
| travel | 23.55% |
| airplane | 2.33% |
| mars | 0.003% |
| mantle | $\epsilon$ |
| … | $\epsilon$ |

*11%* →

cat shoe zebra ball tree jump pen asteroid cable box step cabinet yellow plate flashlight…

**Topic:** *space*

| word | probability |
|------|-------------|
| space | 38.94% |
| airplane | $\epsilon$ |
| mars | 13.43% |
| mantle | 0.05% |
| … | $\epsilon$ |

*89%* →

Will 2020 be the year that humans will embrace space exploration and finally travel to Mars?

# Topic Model Demo #3

# Clustering?



Unlabelled Data → Clustering → Batch Centroid (Centroid Label)

Unlabelled Data → Topic Model (Text Fields) → Batch Topic Distribution (topic 1 prob, topic 3 prob, topic k prob)

# Topic Model Demo #4

# Topic Model Use Cases

- As a preprocessor for other techniques

  - Building better models

- Bootstrapping categories for classification

- Recommendation

- Discovery in large, heterogeneous text datasets

# Topic Model Tips

- ## Setting *k*

  - Much like k-means, the best value is data specific

  - Too few will agglomerate unrelated topics, too many will partition highly related topics

  - I tend to find the latter more annoying than the former

- ## Tuning the Model

  - Remove common, useless terms

  - Set term limit higher, use bigrams

- Create a Source and a Dataset from the StumbleUpon tsv

- Configure a Topic Model (not a 1-click) using:

  - Maximum n-grams=2

  - Exclude non-dictionary words

  - Exclude non-language characters

  - Removing HTML tags

  - Exclude numeric digits

- What is the primary topic for the phrase boilerplate = "No soup for you!"