

CONCEPTUAL ORGANIZATION IN THE SUPERMARKET

Adam Hornsby (@adamnhornsby), Thomas Evans, Peter Riefer, Rosie Prior & Brad Love

<https://arxiv.org/abs/1810.08577>

INTRODUCTION

WHAT IS A TOMATO?



WHAT IS A TOMATO?



It's a **fruit**, it's red, it's **fleshy** and **juicy**...

WHAT IS A TOMATO?



It's a **fruit**, it's **red**, it's **fleshy** and **juicy**...
...but is this **how customers think**?

WHAT IS A TOMATO?



It's a **fruit**, it's **red**, it's **fleshy** and **juicy**...
...but is this **how customers think**?

The answer helps us to **optimize in-store and online search** for customers

TOMATOES ELUDE US BECAUSE...



TOMATOES ELUDE US BECAUSE...



Objects gain meaning through their **interactions** with other objects (Wittgenstein, 1967, Jones & Love, 2007)

TOMATOES ELUDE US BECAUSE...



Objects gain meaning through their **interactions** with other objects (Wittgenstein, 1967, Jones & Love, 2007)
People **categorise** things in terms of **goals** (tomato → salad), as well as taxonomy (tomato → fruit)

TOMATOES ELUDE US BECAUSE...

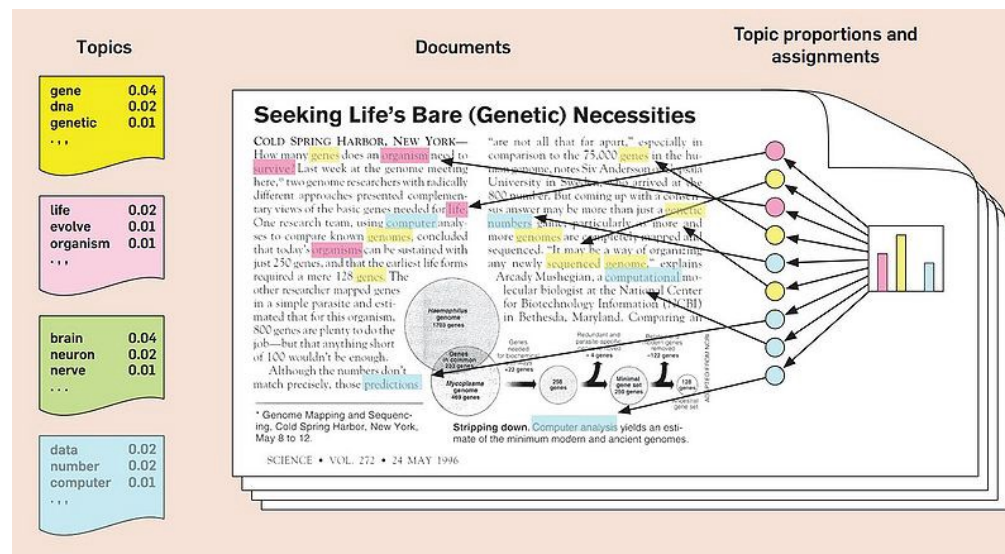


Objects gain meaning through their **interactions** with other objects (Wittgenstein, 1967, Jones & Love, 2007)

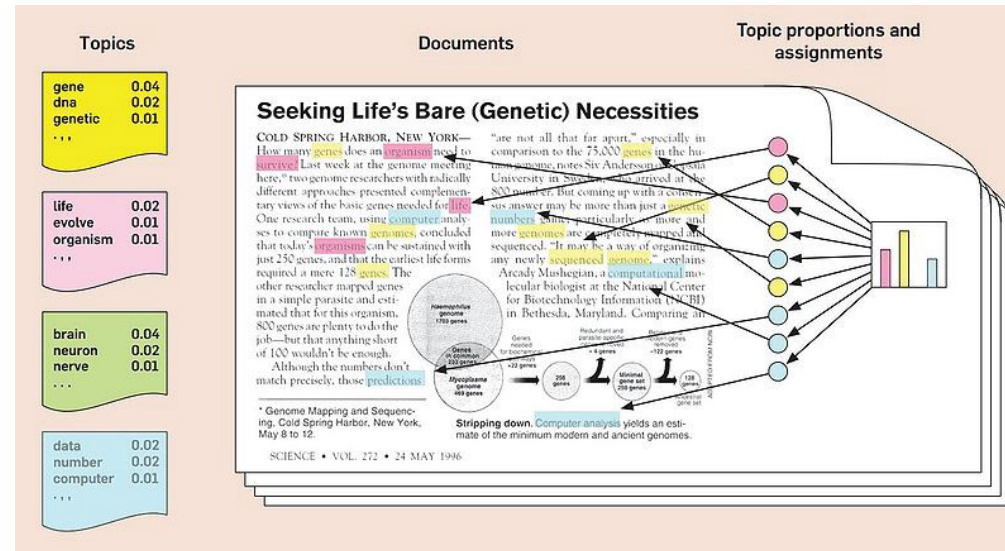
People **categorise** things in terms of **goals** (tomato → salad), as well as taxonomy (tomato → fruit)

Can we **categorise products** in a way that is
more aligned with how customers think?

NLP RESEARCHERS KNOW ABOUT INTERACTIONS

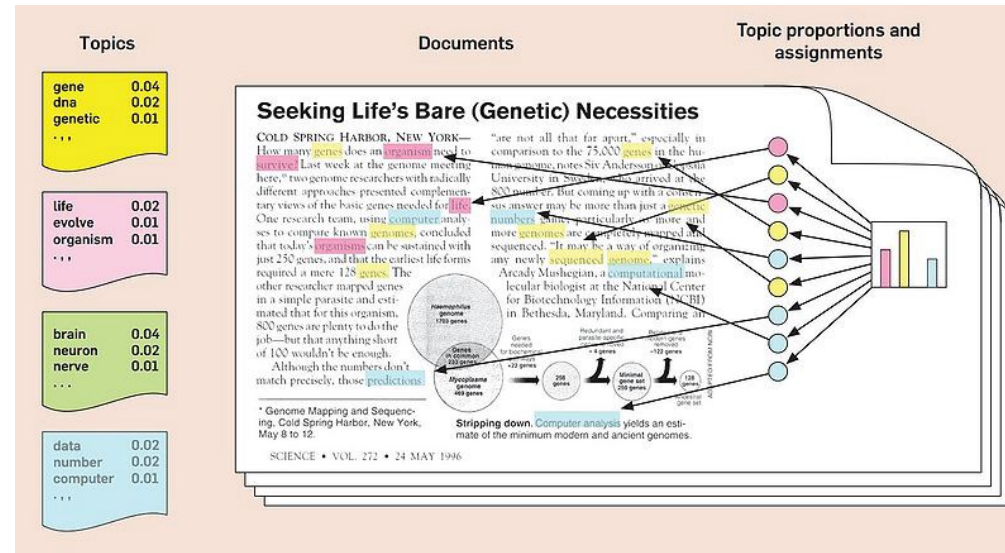


NLP RESEARCHERS KNOW ABOUT INTERACTIONS



"You shall know a word by the company it keeps" (Firth, 1957)
(i.e. Distributional Hypothesis)

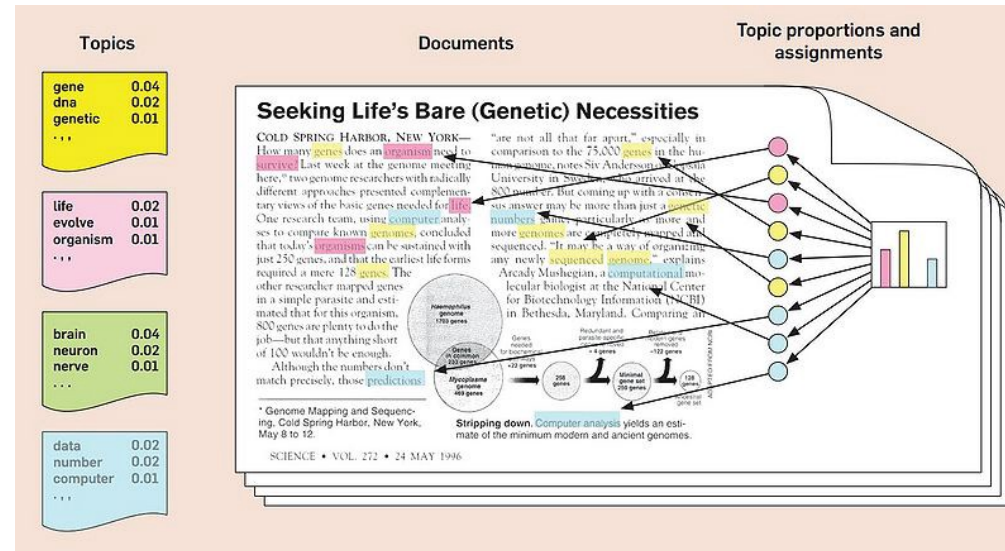
NLP RESEARCHERS KNOW ABOUT INTERACTIONS



"You shall know a word by the company it keeps" (Firth, 1957)
(i.e. Distributional Hypothesis)

Topic models (e.g. LDA) use this premise to learn high-level categories from language data

NLP RESEARCHERS KNOW ABOUT INTERACTIONS



"You shall know a word by the company it keeps" (Firth, 1957)
(i.e. Distributional Hypothesis)

Topic models (e.g. LDA) use this premise to learn high-level categories from language data
So maybe a topic model can learn the mental categories **used by customers?**

IMAGINE A BASKET INSTEAD OF A SENTENCE



Item	Count
Dog	1
Cat	0
Man	1
Bites	1
...	...



Item	Count
Chili	2
Lime	1
Milk	0
Banana	1
...	...

IMAGINE A BASKET INSTEAD OF A SENTENCE

<div><div>NEWS</div><div>MAN BITES DOG</div><div><div></div><div></div></div></div> <table><tr><th>Item</th><th>Count</th></tr><tr><td>Dog</td><td>1</td></tr><tr><td>Cat</td><td>0</td></tr><tr><td>Man</td><td>1</td></tr><tr><td>Bites</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Dog	1	Cat	0	Man	1	Bites	1	<div><div>*****</div><div>RECEIPT</div><div>*****</div><div>APPLE</div><div>BANANA</div><div>BEAN SPROUTS</div><div>CHILI</div><div>CHILI</div><div>LIME</div><div>*****</div></div> <table><tr><th>Item</th><th>Count</th></tr><tr><td>Chili</td><td>2</td></tr><tr><td>Lime</td><td>1</td></tr><tr><td>Milk</td><td>0</td></tr><tr><td>Banana</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Chili	2	Lime	1	Milk	0	Banana	1
Item	Count																								
Dog	1																								
Cat	0																								
Man	1																								
Bites	1																								
...	...																								
Item	Count																								
Chili	2																								
Lime	1																								
Milk	0																								
Banana	1																								
...	...																								

It is straightforward to use existing NLP algorithms on basket data

IMAGINE A BASKET INSTEAD OF A SENTENCE

<div><div>NEWS</div><div>MAN BITES DOG</div><div><div></div><div></div></div></div>	<div><div>*****</div><div>RECEIPT</div><div>*****</div><div>APPLE</div><div>BANANA</div><div>BEAN SPROUTS</div><div>CHILI</div><div>CHILI</div><div>LIME</div><div>*****</div></div>																								
<table><tr><th>Item</th><th>Count</th></tr><tr><td>Dog</td><td>1</td></tr><tr><td>Cat</td><td>0</td></tr><tr><td>Man</td><td>1</td></tr><tr><td>Bites</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Dog	1	Cat	0	Man	1	Bites	1	<table><tr><th>Item</th><th>Count</th></tr><tr><td>Chili</td><td>2</td></tr><tr><td>Lime</td><td>1</td></tr><tr><td>Milk</td><td>0</td></tr><tr><td>Banana</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Chili	2	Lime	1	Milk	0	Banana	1
Item	Count																								
Dog	1																								
Cat	0																								
Man	1																								
Bites	1																								
...	...																								
Item	Count																								
Chili	2																								
Lime	1																								
Milk	0																								
Banana	1																								
...	...																								

It is straightforward to **use existing NLP algorithms** on basket data
Basket data requires **less preprocessing** and is **unordered**, which suits many NLP algorithms better

IMAGINE A BASKET INSTEAD OF A SENTENCE

<div><div>NEWS</div><div>MAN BITES DOG</div><div><div></div><div></div></div></div>	<div><div>*****</div><div>RECEIPT</div><div>*****</div><div>APPLE</div><div>BANANA</div><div>BEAN SPROUTS</div><div>CHILI</div><div>CHILI</div><div>LIME</div><div>*****</div></div>																								
<table><tr><th>Item</th><th>Count</th></tr><tr><td>Dog</td><td>1</td></tr><tr><td>Cat</td><td>0</td></tr><tr><td>Man</td><td>1</td></tr><tr><td>Bites</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Dog	1	Cat	0	Man	1	Bites	1	<table><tr><th>Item</th><th>Count</th></tr><tr><td>Chili</td><td>2</td></tr><tr><td>Lime</td><td>1</td></tr><tr><td>Milk</td><td>0</td></tr><tr><td>Banana</td><td>1</td></tr><tr><td>...</td><td>...</td></tr></table>	Item	Count	Chili	2	Lime	1	Milk	0	Banana	1
Item	Count																								
Dog	1																								
Cat	0																								
Man	1																								
Bites	1																								
...	...																								
Item	Count																								
Chili	2																								
Lime	1																								
Milk	0																								
Banana	1																								
...	...																								

It is straightforward to **use existing NLP algorithms** on basket data
Basket data requires **less preprocessing** and is **unordered**, which suits many NLP algorithms better

Will a **topic model** recover meaningful **categories** from basket data
directly?

RESULTS

OUR TOPIC MODEL DISCOVERED SEMANTIC GROUPS



OUR TOPIC MODEL DISCOVERED SEMANTIC GROUPS



We tuned an LDA model (through Spark 1.6) on **1.2m** real supermarket transactions

OUR TOPIC MODEL DISCOVERED SEMANTIC GROUPS



We tuned an LDA model (through Spark 1.6) on **1.2m real supermarket transactions**
The final **25 topics** appeared coherent and grouped around **specific** (e.g. *Stir fry*) and **general** (e.g. *Cooking from scratch*), **goal-directed themes**

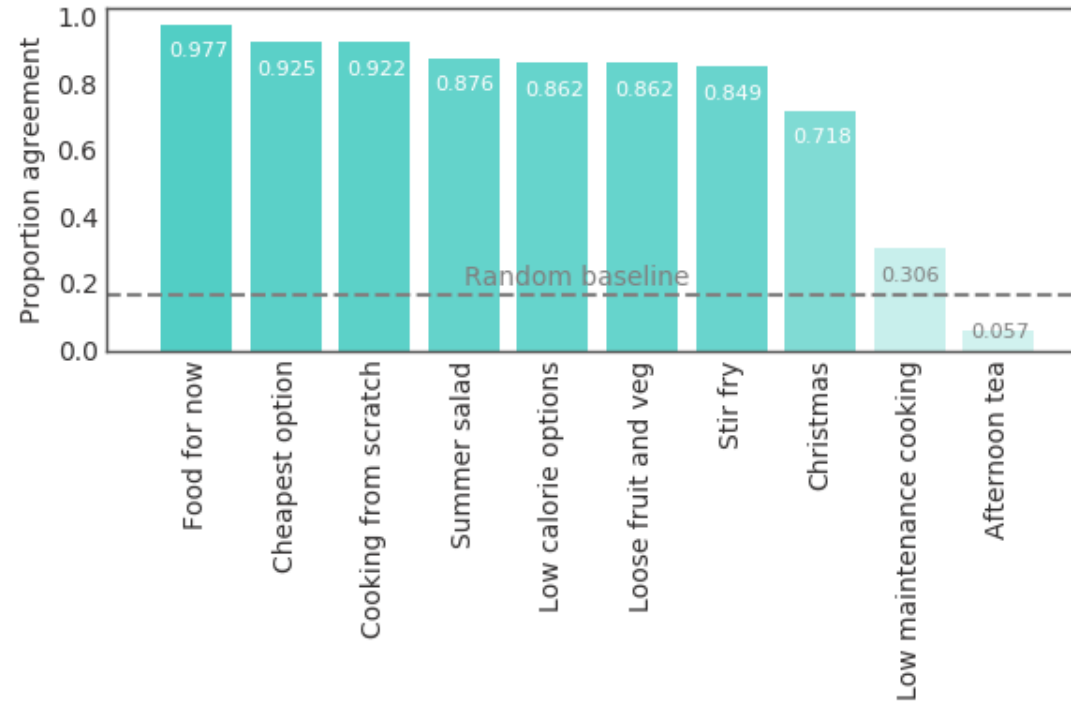
OUR TOPIC MODEL DISCOVERED SEMANTIC GROUPS



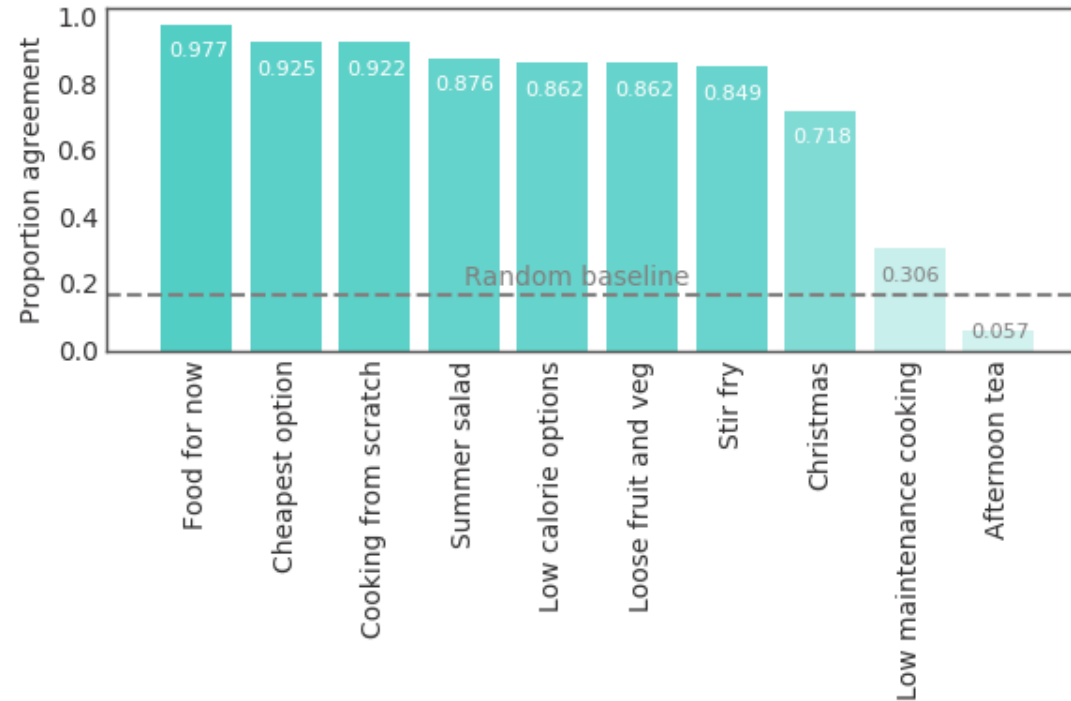
We tuned an LDA model (through Spark 1.6) on **1.2m real supermarket transactions**
The final **25 topics** appeared coherent and grouped around **specific** (e.g. *Stir fry*) and **general** (e.g. *Cooking from scratch*), **goal-directed themes**

So did they make sense to consumers?

CONSUMERS AGREED WITH LDA'S TOPICS

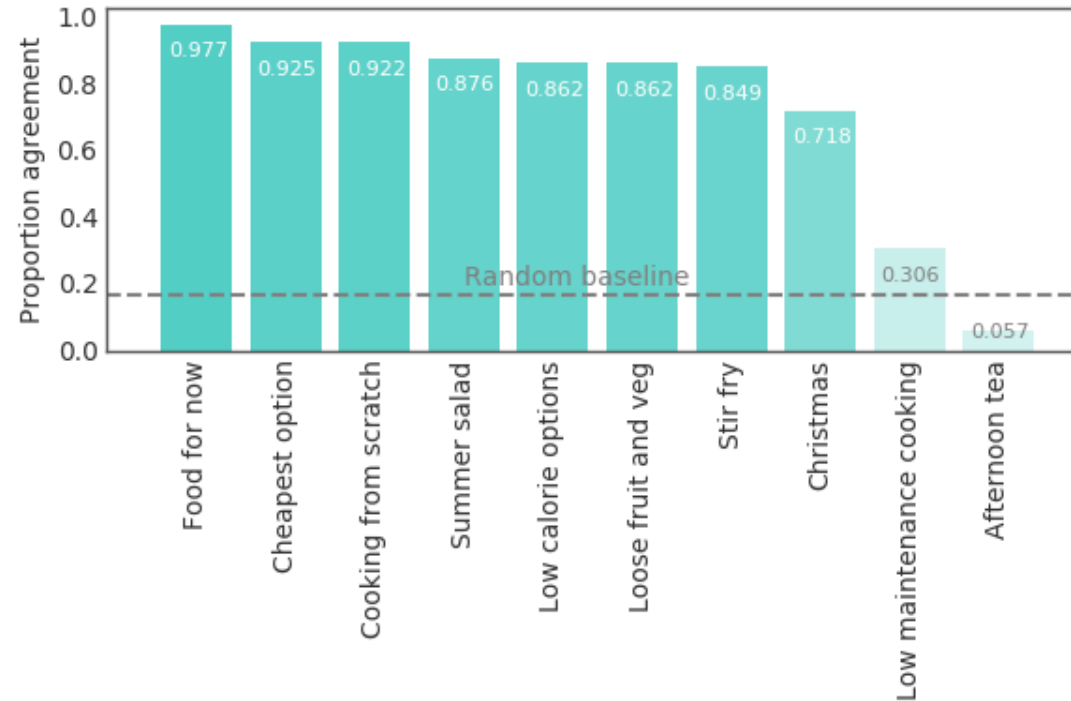


CONSUMERS AGREED WITH LDA'S TOPICS



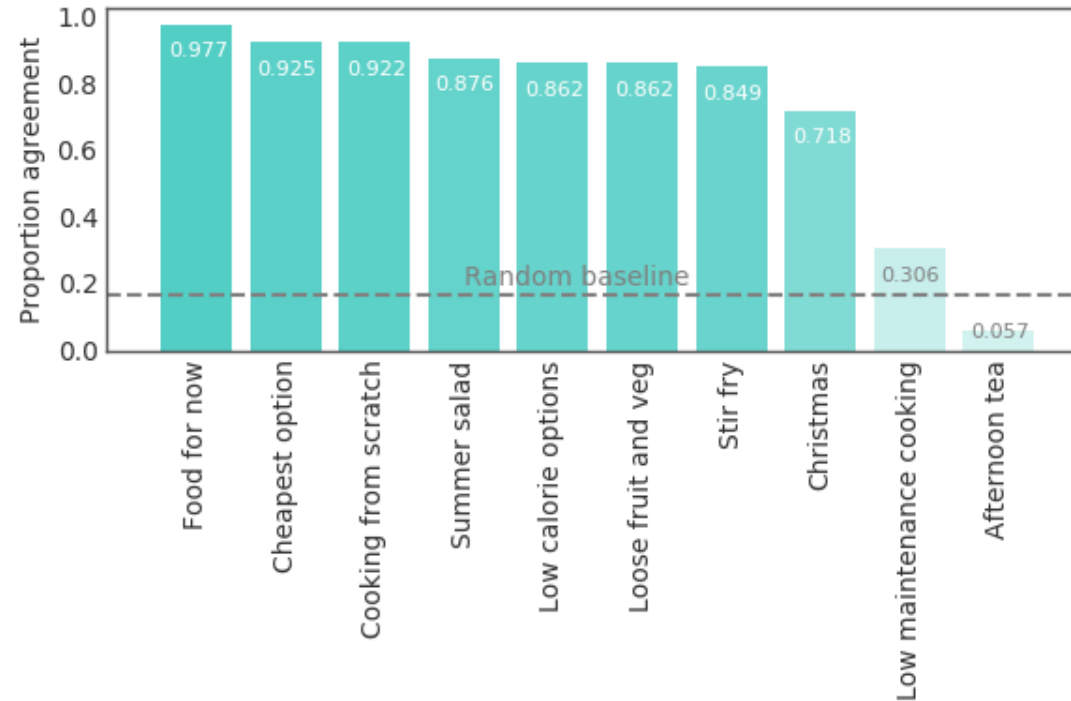
An experiment with **1000 real consumers** showed that they **could identify "intruder" products** accurately

CONSUMERS AGREED WITH LDA'S TOPICS



An experiment with **1000 real consumers** showed that they **could identify "intruder" products** accurately
Suggests that most topics were **similar to mental categories** held by consumers

CONSUMERS AGREED WITH LDA'S TOPICS



An experiment with **1000 real consumers** showed that they **could identify "intruder" products** accurately
Suggests that most topics were **similar to mental categories** held by consumers

Some were difficult (e.g. *Afternoon tea*), perhaps
due to **individual differences**

dunnhumby

TOPICS ALSO PREDICTED INDIVIDUAL DIFFERENCES



TOPICS ALSO PREDICTED INDIVIDUAL DIFFERENCES



The topics customers purchased **predicted** their self-reported age, location and gender

TOPICS ALSO PREDICTED INDIVIDUAL DIFFERENCES



The topics customers purchased **predicted their self-reported age, location and gender**
Suggests that people's **mental representations of products** may differ between individuals and groups

TOPICS ALSO PREDICTED INDIVIDUAL DIFFERENCES



The topics customers purchased **predicted their self-reported age, location and gender**
Suggests that people's **mental representations of products** may differ between individuals and groups

This can help us to **personalise search algorithms**

CONCLUSIONS

CONCLUSIONS



CONCLUSIONS



Consumers think about products in terms of their **roles with other products**
(e.g. "goes well in a salad")

CONCLUSIONS



Consumers think about products in terms of their **roles with other products**
(e.g. "goes well in a salad")

Differences in **experience** may lead to **different mental categories**

CONCLUSIONS



Consumers think about products in terms of their **roles with other products**
(e.g. "goes well in a salad")

Differences in **experience** may lead to **different mental categories**

Topic models apply to **other data sources** (e.g. baskets)

CONCLUSIONS



Consumers think about products in terms of their **roles with other products**
(e.g. "goes well in a salad")

Differences in **experience** may lead to **different mental categories**

Topic models apply to **other data sources** (e.g. baskets)

We are using this insight to **optimize customer's routes through the store** (both online and offline) (e.g. dual siting)

CONCLUSIONS



Consumers think about products in terms of their **roles with other products**
(e.g. "goes well in a salad")

Differences in **experience may lead to different mental categories**

Topic models apply to **other data sources** (e.g. baskets)

We are using this insight to **optimize customer's routes through the store** (both online and offline) (e.g. dual siting)

An **alternative source of data** for evaluating NLP models
(data & code available soon)

dunnhumby

THANK YOU

arxiv.org/abs/1810.08577

Adam Hornsby (@adamnhornsby), Thomas Evans, Peter Riefer, Rosie Prior & Brad Love



National Institutes of Health

The Leverhulme Trust

wellcometrust



dunnhumby