



**Static models:
A (relevant) message from the
past**



Queen Mary
University of London

The
Alan Turing
Institute



UNIVERSITY OF
CAMBRIDGE

Hello again!



Haim Dubossarsky

Researcher @ Change is Key!

Queen Mary University of London (QMUL)

Language Technology Lab, University of Cambridge

The Alan Turing Institute

h.dubossarsky@qmul.ac.uk

Outline & takehomes

old ≠ not useful

- Brief recap on static models (an introduction to some?)
 - **Static models are not contextualized models**
 - Explicit models: Count-based, PPMI (interpretable dimensions)
 - Predictive models: word2vec and its likes (“opaque” dimensions)
- Doing **semantic change** with static models: features, pros and cons
 - 1-word : all meanings – polysemy
 - Measure change in meaning via cosine-distance
 - Can work well with small corpora
 - Some models provide much more detailed report of change
- Does not cover all models in this overview: e.g., topic models

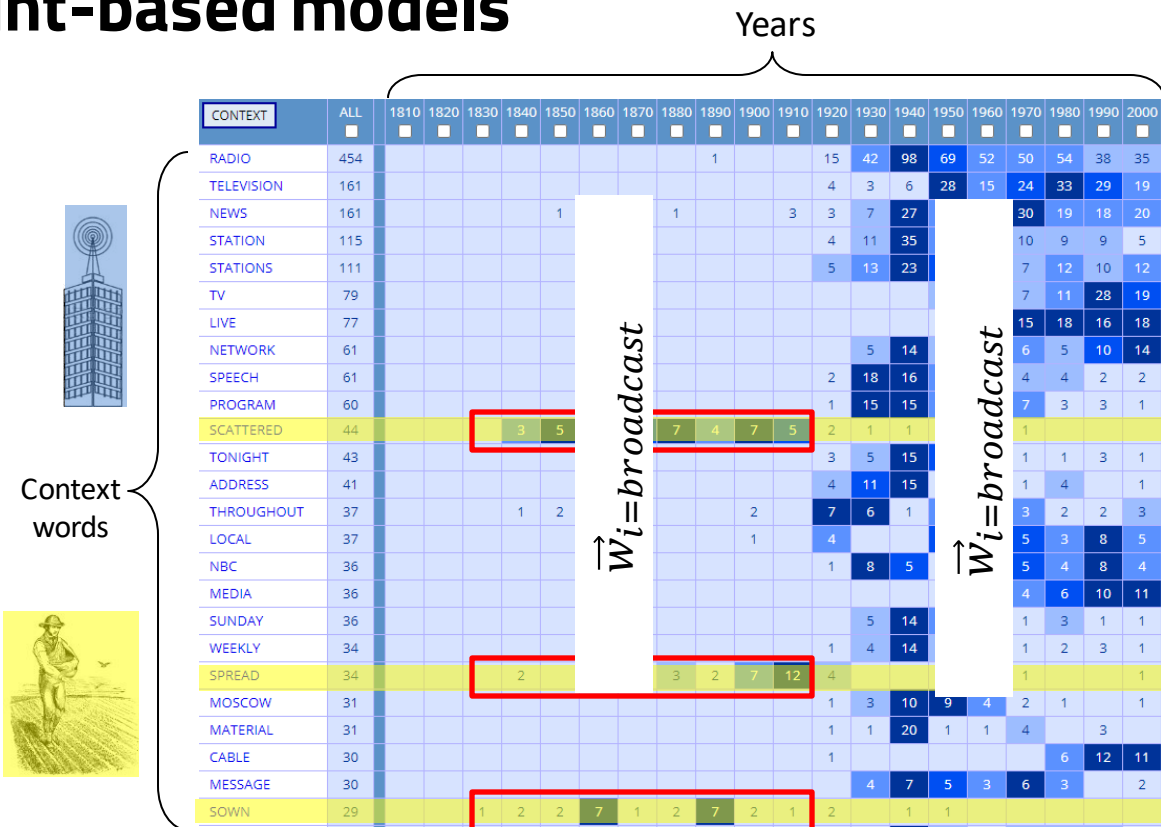
Count-based models

Simple co-occurrence models within a context window

Very sparse

CONTEXT	ALL	1920s	1930s	1940s	1950s	1960s	1970s	1980s	1990s	2000s
COMPUTER	171						1	75	67	28
BIG	170	2	26	8	2	7	31	43	29	22
PIE	148	4	3	22	18	16	29	31	17	8
TREE	94	6	18	16	17	15	9	9	1	3
ADAM	59	8	12	9	7	9	8	3	3	
JOBS	57						1	14	23	19
TREES	51	4	10	18	5	6	3	4	1	
JUICE	45	1	2		3	5	9	7	11	7
MACINTOSH	43							24	13	6
IBM	41							24	17	

Count-based models



Count-based models

- Very rare: Most cases will not be so clear



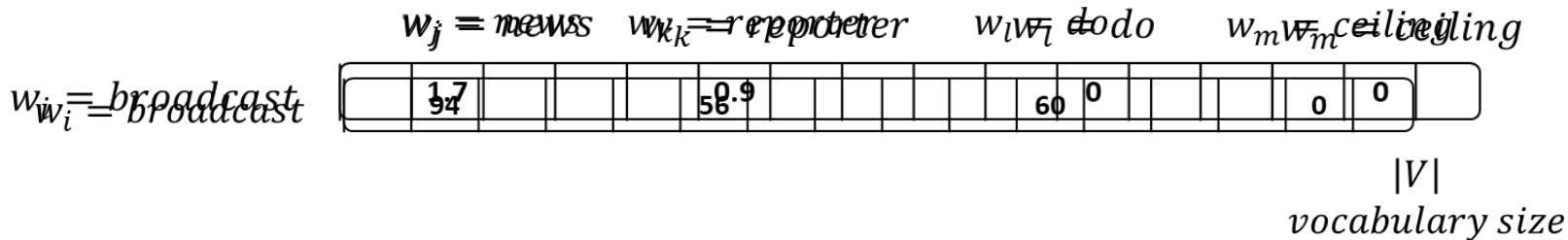
- Highly dependent and reflects the meaning of the corpus/domain
 - True for **all static & contextualized** models
 - More apparent here as static models are not “pre-trained”

- Problem:** Highly skewed for frequent collocates
 - Prepositions, function words (stopwords)
 - Solution: ????



Positive Pointwise Mutual Information (PPMI)

- Co-occurrence models within a context window **with a twist**
 - Twist:** Mutual information measures the strength of association between the target word and its co-occurring words



- Learn associativity by informativity**

Positive Pointwise Mutual Information (PPMI)

- Only “strong” co-occurring words are retained, hence “positive” PMI

$$\text{PPMI}(w, c) = \max\left(\log_2 \frac{P(w, c)}{P(w)P(c)}, 0\right)$$

		Count(w,context)					
		computer	data	pinch	result	sugar	
apricot apple digital informatio			p(w,context)				
			computer	data	pinch	result	sugar
			PPMI(w,context)				
			computer	data	pinch	result	sugar
			apricot	-	-	2.25	-
		apple	-	-	2.25	-	2.25
		digital	1.66	0.00	-	0.00	-
		information	0.00	0.57	-	0.47	-
		p(context)					

Advantages of explicit models (count-based & PPMI)

$w_j = \text{news}$ $w_k = \text{reporter}$ $w_l = \text{do}$ $w_m = \text{ceiling}$
 $w_i = \text{broadcast}$

	1.7				0.9					0			0	
--	-----	--	--	--	-----	--	--	--	--	---	--	--	---	--

- Enables a finer analysis of change (association level)
 - Used in research: Stefanowitsch & Gries Collostructions Analysis (2003)

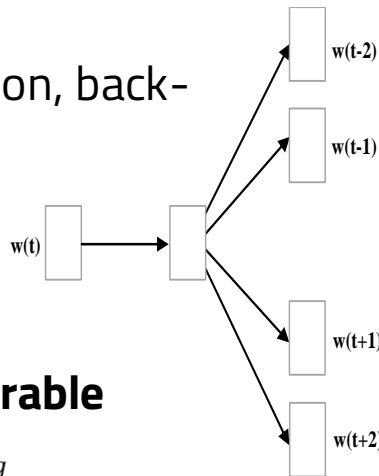
PPMI for <i>prime</i> ¹⁹⁹⁰		
ministe	PPMI for <i>prime</i> ²⁰¹⁰	
suspec	minister	11.26
cut	numbers	9.51
numbe	cut	10.1

PPMI for <i>heart</i> ^{medical}		
attack	PPMI for <i>heart</i> ^{standard}	
chest	attack	13.4
pacemaker	emotion	4.9
	central	4.5
	warmth	3.2

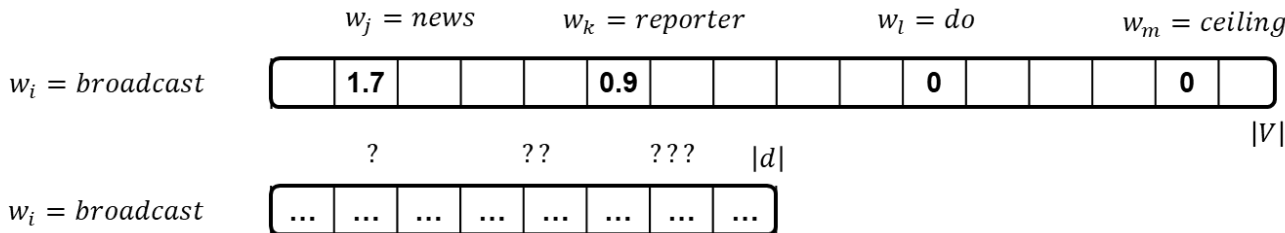
Predictive models (word2vec)

- Word2vec (Mikolov et al. 2013) is a Neural Network model
 - Shallow network: 1 layer
 - Uses known NN machinery: MLM, objective function, back-propagation, SGD, etc.

$$J_{\theta} = \frac{1}{T} \sum_{t=1}^T \sum_{-n \leq j \leq n, j \neq 0} \log p(w_{t+j} | w_t)$$



- Vectors are now opaque & **vector spaces are incomparable**



Predictive models (word2vec)

Even worse models
are sometimes better

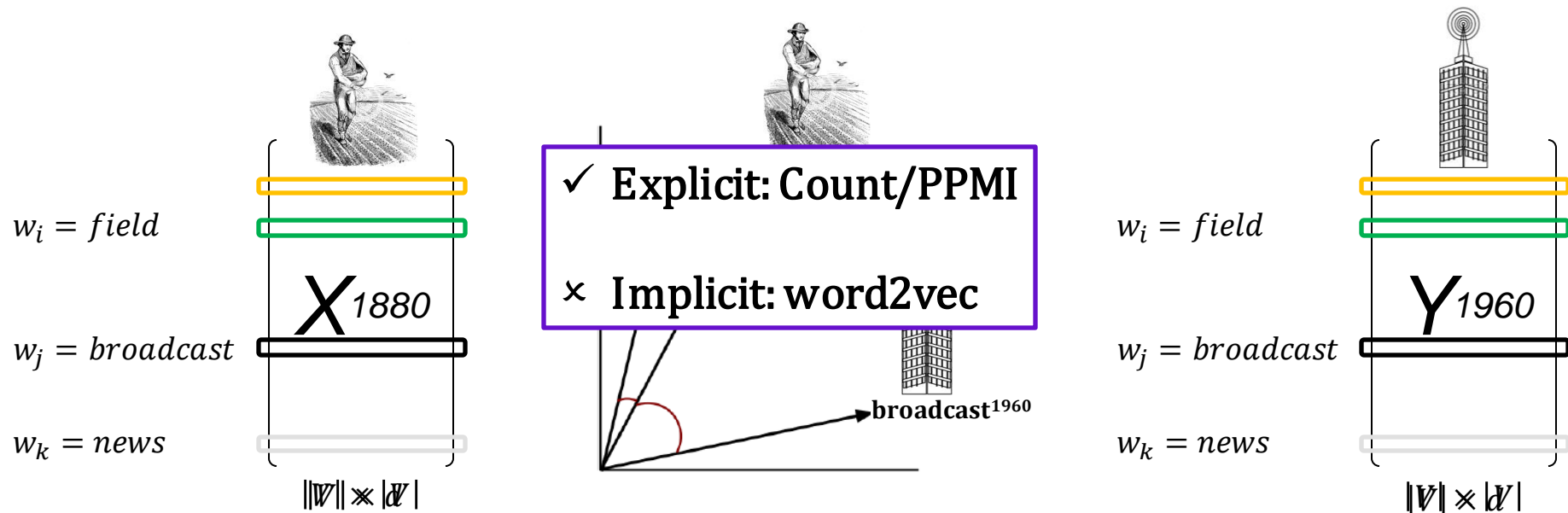
- Why word2vec is more popular than PPMI?
 - Easier and more efficient implementation
 - PR: Nice demonstration of abilities (analogy solving etc.)
 - Simply because of sheer numbers of users
- Is word2vec better than PPMI? Sometimes, but often not.

Method	WordSim Similarity	WordSim Relatedness	Bruni et al. MEN	Radinsky et al. M. Turk	Luong et al. Rare Words	Hill et al. SimLex	Google Add / Mul	MSR Add / Mul
PPMI	.755	.697	.745	.686	.462	.393	.553 / .679	.306 / .535
SVD	.793	.691	.778	.666	.514	.432	.554 / .591	.408 / .468
SGNS	.793	.685	.774	.693	.470	.438	.676 / .688	.618 / .645

From Levy et. al. 2015

- Word2vec-like models are mathematically equivalent to PPMI (Levy et. al., 2014, 2015)
- Not the right question: **Is word2vec better for Semantic change?**

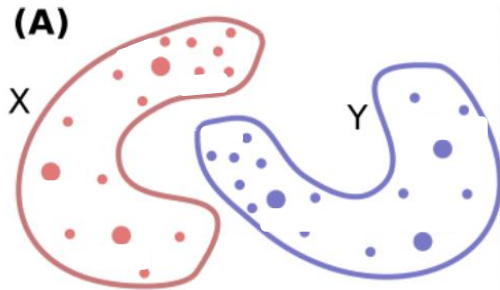
Reminder: measuring change computationally



$$Change = \frac{\vec{w}^{t1} \cdot \vec{w}^{t2}}{\|\vec{w}^{t1}\| \cdot \|\vec{w}^{t2}\|}$$

Lexical semantic change with w2v-like models

- Word2vec models are initiated with random parameters. Hence, **if we don't do something about it**, their vectors lie in difference spaces, and are incomparable.



From Conneau et al. 2018

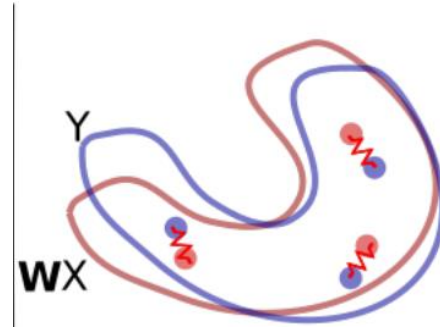
- Solutions:
 - Aligning the vector spaces prior to comparison
 - Avoiding the need for alignment

Aligning vector spaces

We need to find $\varphi(X) \rightarrow Y$

$$W^* = \operatorname{argmin}_W \|WX - Y\|_2$$

$$X, Y \in \mathbb{R}^d$$



Under orthogonal constraint ($W^T W = I$) the solutions is:

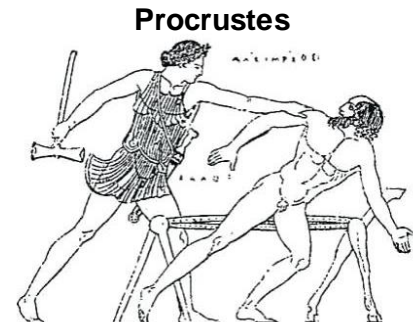
$$U \Sigma V^T = \operatorname{SVD}(YX^T)$$

$$W = UV^T$$

Assumptions

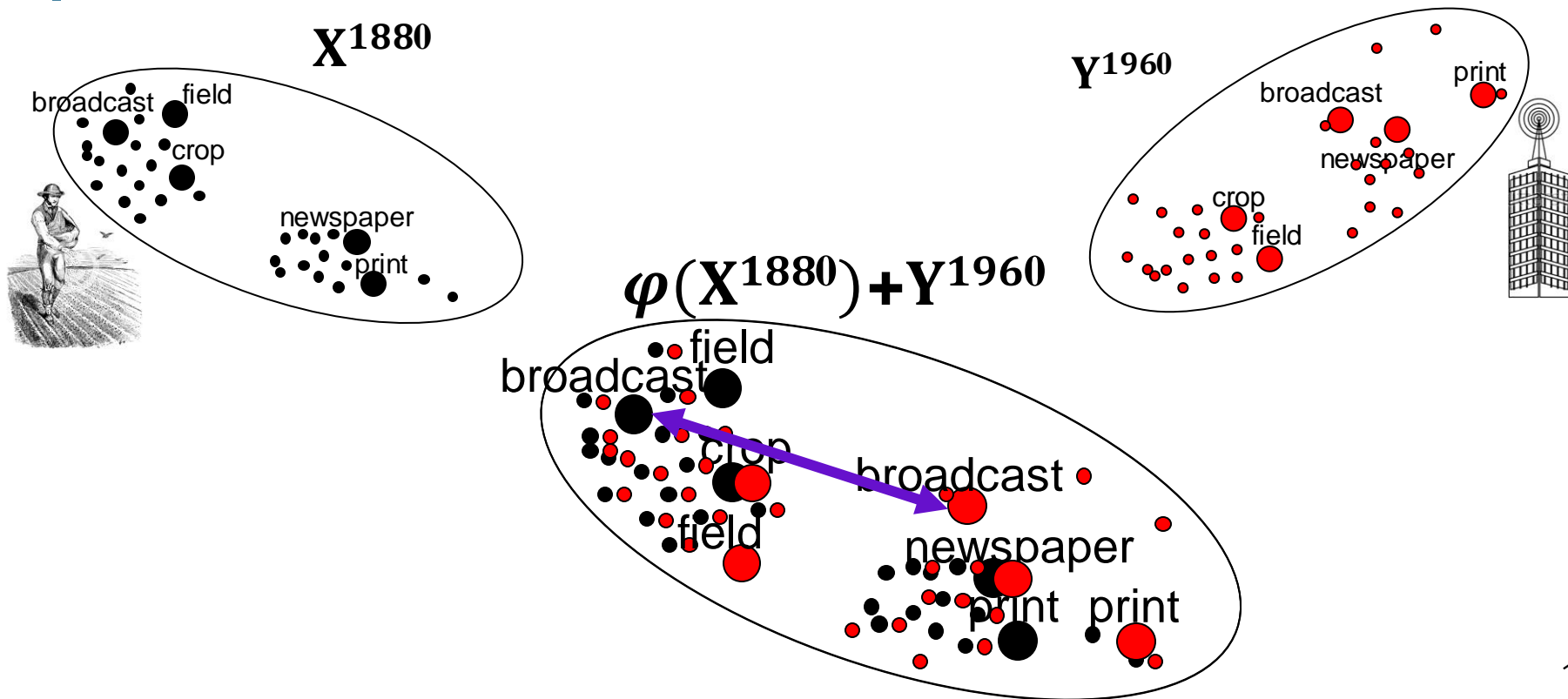
We have 1:1 mapping (dictionaries)

Vector spaces are comparable (isometric)



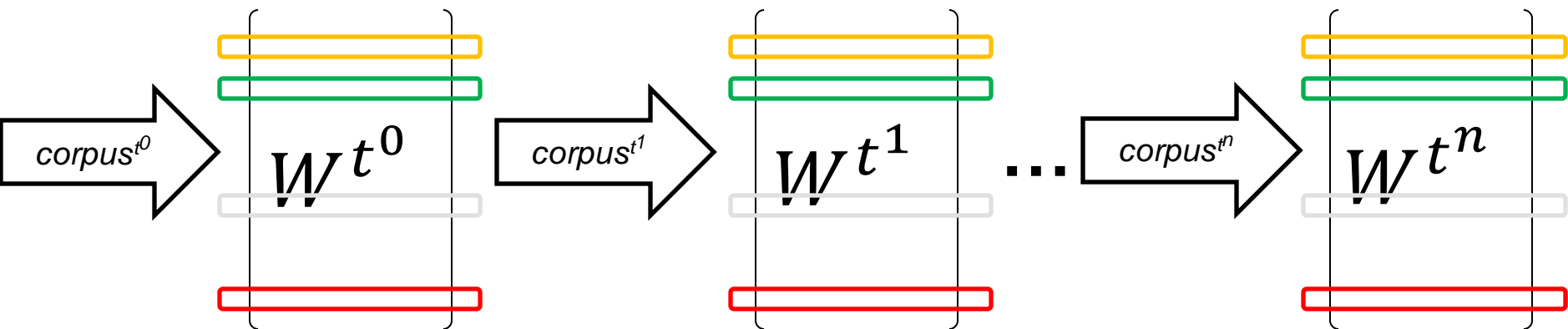
Alignment is not noise free.
What is the nature of the noise?

Aligning vector spaces



Avoiding alignment I

- Incremental training (Kim et al., 2014)
 - For every time step, model is initiated with the parameters of the trained model from the previous step.
 - Causes drift (noise) for the entire vector space

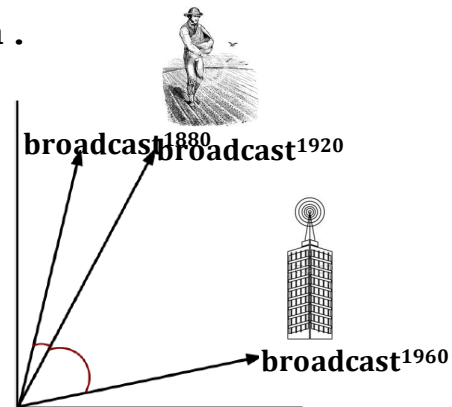


Avoiding alignment II

- Temporal referencing (Dubossarsky et al. 2019)
 - Words are tagged according to the time of corpus
 - Observed the least amount of noise

Example

Silken cauliflowers sown broadcast¹⁸⁷⁰ over the land.
The dramatic broadcast¹⁹⁷⁰ stunned the nation.



Take homes

old \neq not useful

- Brief recap on static models (an introduction to some?)
 - **Static models are not contextualized models**
 - Explicit models: Count-based, PPMI (interpretable dimensions)
 - Predictive models: word2vec and its likes (“opaque” dimensions)
- Doing **semantic change** with static models: features, pros and cons
 - 1-word : all meanings – polysemy
 - Measure change in meaning via cosine-distance
 - Can work well with small corpora
 - Some models provide much more detailed report of change