



HOW TO TURN WORDS INTO VECTORS?

Sandeep Soni

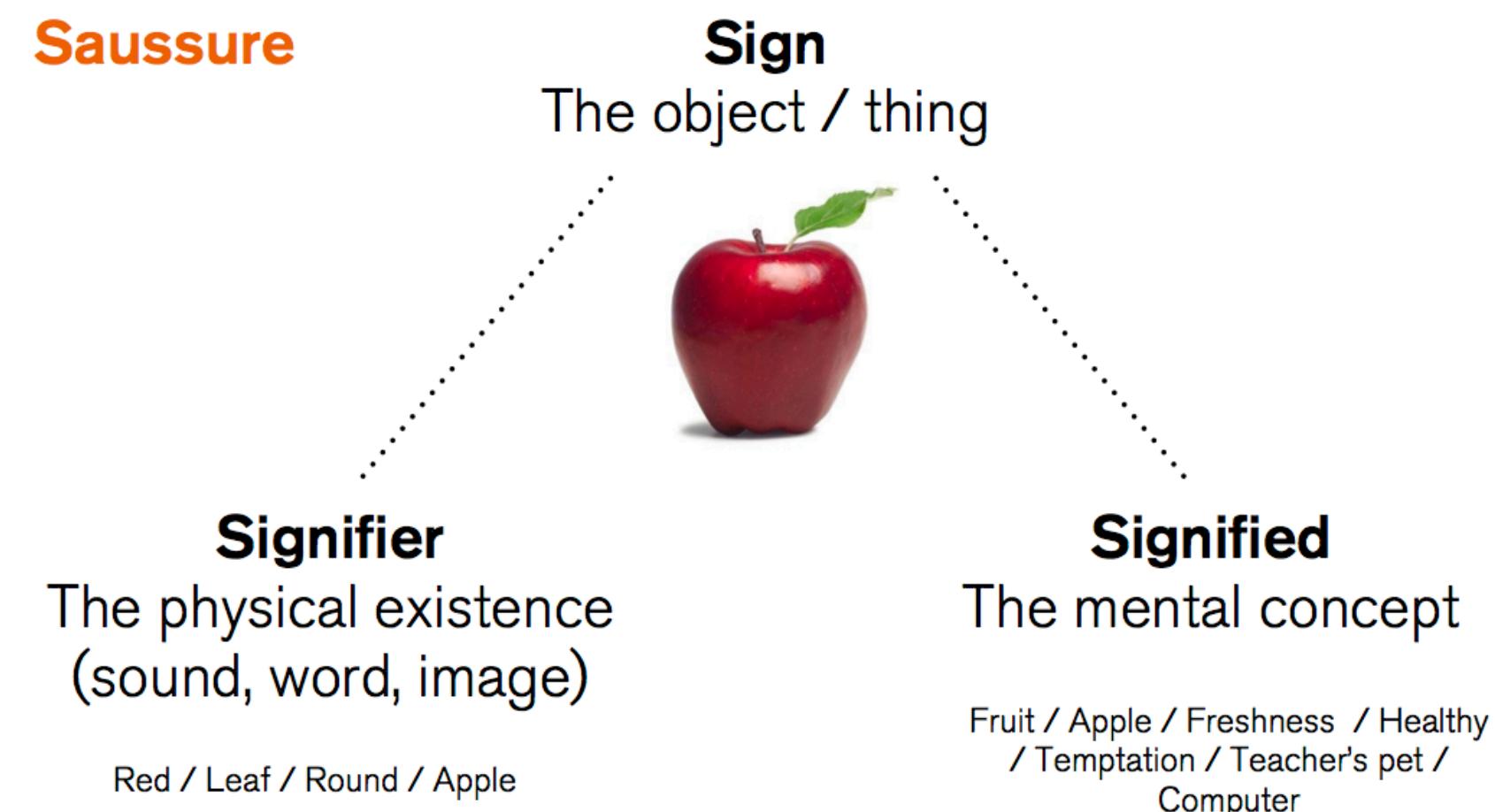
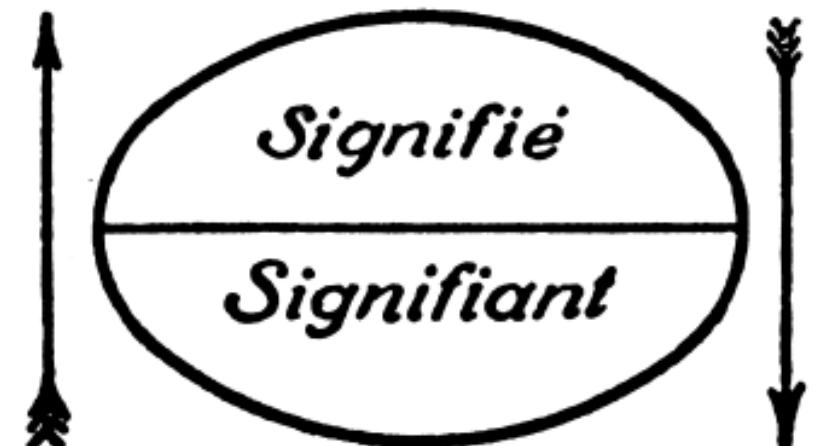
09/05/2024

STORY SO FAR

- Words/Types/Tokens/Stems/Lemmas
- Tokenization
- Frequency based statistics to compare corpora/groups
 - χ^2 test
 - Pointwise mutual information (PMI)

MEANING

- Saussure proposed two characteristics of signs
 - Signified or the “plane of content”
 - Signifier or the “plane of expression”
- Words are a surface representation of concepts



What is the meaning of the word tezgüino?

A bottle of

is on the table

Everybody

likes

Don't have

before you drive

We make

out of corn

Example taken from Eisenstein 2018, which was in turn taken from Lin, 1998

Is tezgüino similar to loud?

A bottle of

is on the table

Everybody

likes

Don't have

before you drive

We make

out of corn

Example taken from Eisenstein, 2018, which was in turn taken by Lin, 1998

Is tezgüino similar to oil?

A bottle of

is on the table

Everybody

likes

Don't have

before you drive

We make

out of corn

Example taken from Eisenstein, 2018, which was in turn taken by Lin, 1998

Is tezgüino similar to tortilla?

A bottle of

is on the table

Everybody

likes

Don't have

before you drive

We make

out of corn

Example taken from Eisenstein, 2018, which was in turn taken by Lin, 1998

Is tezgüino similar to beer?

Contexts

A bottle of

is on the table

Everybody

likes

Don't have

before you drive

We make

out of corn

Example taken from Eisenstein, 2018, which was in turn taken by Lin, 1998

QUESTION FOR THE DAY

“How do we represent word meaning?”

SEMANTICS

- Three perspectives in semantics
 - Relational
 - Compositional
 - Distributional

SEMANTICS

Relational

John interviews **Mark**. He is a great **tennis player**

Compositional

compose —> composition —> compositional

Distributional

The **paint** is still wet

Paint that wall red

The old **paint** is coming off

“a word is characterized by the company it keeps”

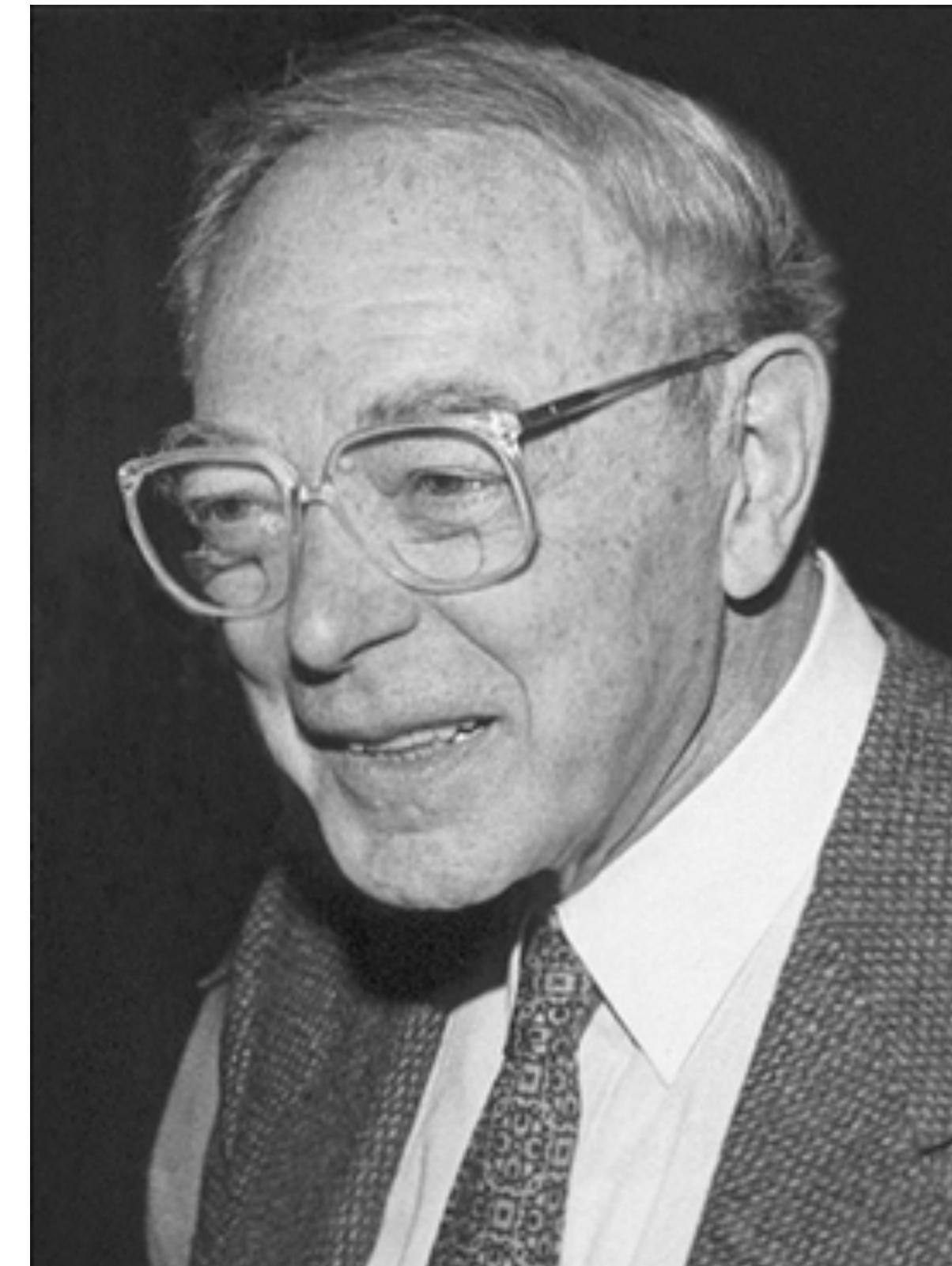
-J.R.Firth

DISTRIBUTIONAL SEMANTICS

- Word meanings can be learned from the contexts in which they appear
- A foundational block in modern NLP



J.R.Firth



Zellig Harris

DISTRIBUTED REPRESENTATIONS

- Distributed representations are vector representations that summarize the distribution of the contexts in which a word appears.
- Words that appear in similar contexts should have similar representations (distributional hypothesis)
- Key questions:
 - What type of contexts?
 - What type of vectors?
 - How to systematically do this?

CORPUS

Imagine that you have a bunch of reviews

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.

In general, such a collection is called a text corpus

VOCABULARY

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.

You can lowercase all reviews and then
tokenize to construct a vocabulary

VOCABULARY

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.

the
temperature
was
cold
and
food
horrible
service
wine
great
beer
ok
...
not
even

You can lowercase all reviews and then
tokenize to construct a vocabulary

VOCABULARY

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.

the	5
temperature	1
was	6
cold	4
and	1
food	1
horrible	1
service	4
wine	6
great	1
beer	4
ok	1
...	...
not	1
even	1

You can lowercase all reviews and then tokenize to construct a vocabulary

TERM-DOCUMENT MATRIX

Vocabulary

Reviews

Reviews

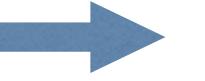
- 1 The temperature was cold and the food was cold ... horrible service.
- 2 Wine was great. Beer was OK. Go here if you are into wine Wine WINE
- 3 I like the restaurant because their service is good ... service is key!
- 4 I mostly order a single wine or a wine and cold beer
- 5 I ordered a beer. The service was pathetic. The beer was not even cold.

TERM-DOCUMENT MATRIX

Vocabulary

Reviews

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.



TERM-DOCUMENT MATRIX

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.



Vocabulary

<i>the</i>	5
<i>temperature</i>	1
<i>was</i>	6
<i>cold</i>	4
<i>and</i>	1
<i>food</i>	1
<i>horrible</i>	1
<i>service</i>	4
<i>wine</i>	6
<i>great</i>	1
<i>beer</i>	4
<i>ok</i>	1
...	...
<i>not</i>	1
<i>even</i>	1

Reviews

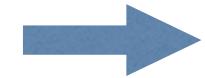
TERM-DOCUMENT MATRIX

Reviews	
1	The temperature was cold and the food was cold ... horrible service.
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer. The service was pathetic. The beer was not even cold.

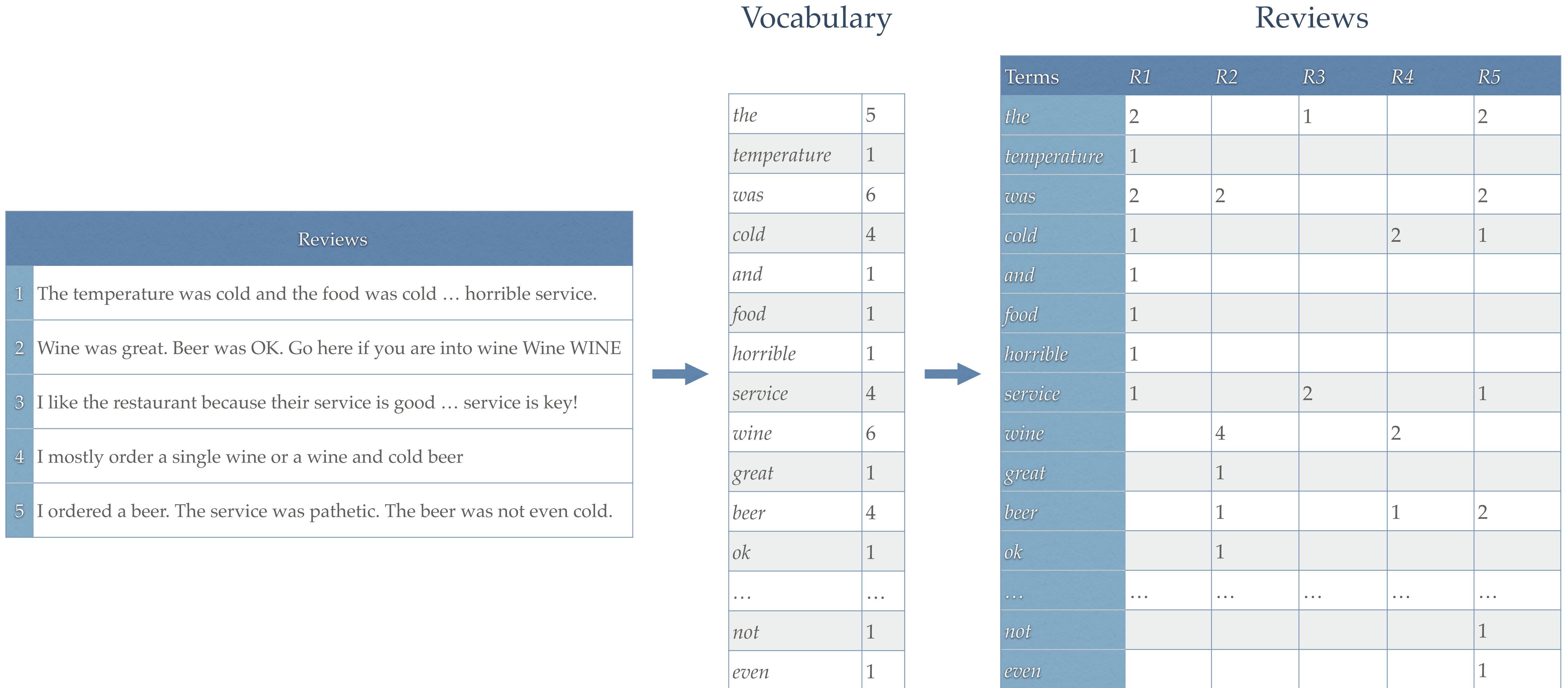
Vocabulary

<i>the</i>	5
<i>temperature</i>	1
<i>was</i>	6
<i>cold</i>	4
<i>and</i>	1
<i>food</i>	1
<i>horrible</i>	1
<i>service</i>	4
<i>wine</i>	6
<i>great</i>	1
<i>beer</i>	4
<i>ok</i>	1
...	...
<i>not</i>	1
<i>even</i>	1

Reviews



TERM-DOCUMENT MATRIX



TERM-DOCUMENT MATRIX

- Every term in the vocabulary is mapped to a row in this matrix
- Each row is a vector encoding a distribution of counts over reviews/documents
- Context = entire document

TERM-DOCUMENT MATRIX

Terms	R1	R2	R3	R4	R5
<i>the</i>	2		1		2
<i>temperature</i>	1				
<i>was</i>	2	2			2
<i>cold</i>	1			2	1
<i>and</i>	1				
<i>food</i>	1				
<i>horrible</i>	1				
<i>service</i>	1		2		1
<i>wine</i>		4		2	
<i>great</i>		1			
<i>beer</i>		1		1	2
<i>ok</i>		1			
...
<i>not</i>					1
<i>even</i>					1

- Every term in the vocabulary is mapped to a row in this matrix
- Each row is a vector encoding a distribution of counts over reviews/documents
- Context = entire document

TERM-DOCUMENT MATRIX

- Blank entries in this matrix correspond to zero, so vectors are sparse!
- The size of each vector is equal to the number of documents

TERM-DOCUMENT MATRIX

Terms	R1	R2	R3	R4	R5
<i>the</i>	2		1		2
<i>temperature</i>	1				
<i>was</i>	2	2			2
<i>cold</i>	1			2	1
<i>and</i>	1				
<i>food</i>	1				
<i>horrible</i>	1				
<i>service</i>	1		2		1
<i>wine</i>		4		2	
<i>great</i>		1			
<i>beer</i>		1		1	2
<i>ok</i>		1			
...
<i>not</i>					1
<i>even</i>					1

- Blank entries in this matrix correspond to zero, so vectors are sparse!
- The size of each vector is equal to the number of documents

CORPUS

Reviews	
1	The temperature was cold and the food was cold ... horrible service .
2	Wine was great. Beer was OK. Go here if you are into wine Wine WINE
3	I like the restaurant because their service is good ... service is key!
4	I mostly order a single wine or a wine and cold beer
5	I ordered a beer . The service was pathetic. The beer was not even cold

Let's focus on these words for now. Is wine more similar to beer or cold?

VECTOR SIMILARITY

- The similarity between two vectors can be calculated by taking the dot product between them.
- If vectors are normalized, then dot products can be interpreted as the cosine of the angle between them.

$$\begin{array}{lll} a \cdot b & \text{dot}(a, b) & \sum_i a_i b_i \\ \langle a, b \rangle & a^T b & \\ b \cdot a & & \langle b, a \rangle \\ b^T a & \|a\| \|b\| \cos \theta & \end{array}$$

These are all the same!

VECTOR SIMILARITY

Reviews

	$R1$	$R2$	$R3$	$R4$	$R5$
\dots	\dots	\dots	\dots	\dots	\dots
<i>wine</i>		4		2	
<i>cold</i>	2			1	1
<i>beer</i>		2		1	2
\dots	\dots	\dots	\dots	\dots	\dots

$\text{dot}(\text{wine}, \text{ cold})$

$$= 0*2 + 4*0 + 0*0 + 2*1 + 0*1$$

$$= 2$$

VECTOR SIMILARITY

Reviews

	$R1$	$R2$	$R3$	$R4$	$R5$
\dots	\dots	\dots	\dots	\dots	\dots
<i>wine</i>		4		2	
<i>cold</i>	2			1	1
<i>beer</i>		2		1	2
\dots	\dots	\dots	\dots	\dots	\dots

$\text{dot}(\text{wine}, \text{ beer})$

$$= 0*0 + 4*2 + 0*0 + 2*1 + 0*2$$

$$= 10$$

COOCCURENCE MATRIX

- Instead of taking the entire document as context, we can define a context window of some size
- We can construct a matrix with number of rows and columns equal to vocabulary size
- Each cell's value is the number of times a word pair corresponding to the row and column appear together.

Dataset

The old wine tastes good

The bottled beer is stale

The red wine is stale

Store the beer for long

Assume context window size of 2

All words that co-appear with wine are:

{the, old, tastes, good, red, is, stale}

Dataset

The old wine tastes good

The bottled beer is stale

The red wine is stale

Store the beer for long

Assume context window size of 2

All words that co-appear with wine are:

{the, bottled, is, stale, store, for, long}

COCCURRENCE MATRIX

Contexts	<i>the</i>	<i>bottled</i>	<i>tastes</i>	<i>good</i>	...	<i>stale</i>	<i>long</i>
...							
<i>wine</i>	2		1	1		1	
<i>beer</i>	2	1				1	
...							

- Every term in the vocabulary is mapped to a row in this matrix
- Each row is a vector encoding a distribution of counts over its cooccurring words
- Context = window around the term

VECTORS

Terms (t)	Contexts (c)						
	<i>the</i>	<i>bottled</i>	<i>tastes</i>	<i>good</i>	...	<i>stale</i>	<i>long</i>
...							
<i>wine</i>	2		1	1		1	
<i>beer</i>	2	1				1	
...							

- The vectors are still sparse!
- The matrix grows quadratically with the size of the vocabulary

WEIGHTING

- Many words cooccur with functional words such as “the”, “a”, etc
- So can you weigh the dimensions based on their informativeness?

TFIDF

- Idea: Upweight the cooccurrence counts for terms that cooccur frequently and exclusively.
- $tfidf(t, c) = tf(t, c) \times \frac{N}{d_t}$
- N is the total number of contexts, $tf(t, c)$ is the number of times t co-appears with c, d_t is the number of times t appears in any context

PMI

- Idea: Upweight the concurrence counts if the pair is more likely to cooccur than chance

- $PMI(t, c) = \log_2 \frac{P(t, c)}{P(t) \cdot P(c)}$

- $P(t, c)$ is the probability of cooccurrence; $P(t)$ and $P(c)$ is probability of independent occurrence

- $PPMI(t, c) = \max(0, \log_2 \frac{P(t, c)}{P(t) \cdot P(c)})$

DENSE REPRESENTATIONS

- Ideally, we want to learn characteristic dimensions of a word, instead of distributions over all words/contexts
- We learn a low-dimensional but compact/dense vectors by formulating the cooccurrence events as prediction tasks

Word	Ball game	Non-square arena	Olympic Game	Indoor
Soccer	1	0	1	0
Javelin Throw	0	1	1	0
Squash	1	0	0	1
Chess	0	1	0	1

SO FAR!

- Vector representations are sparse and long
- They're based on counting

MATRIX FACTORIZATION

- You can factorize the cooccurrence matrix
- The vectors in the factors are low-dimensional and dense

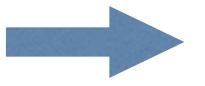
$$\mathbf{M}_{m \times n} = \mathbf{U}_{m \times m} \Sigma_{m \times n} \mathbf{V}^*_{n \times n}$$

Image taken from Wikipedia

WORD2VEC

- Skipgram (Mikolov et. al. 2013) predicts a context word given a target word
- Think of this as a huge multi-class classification task (classes are words) for every word

Cooccurrence Matrix



x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

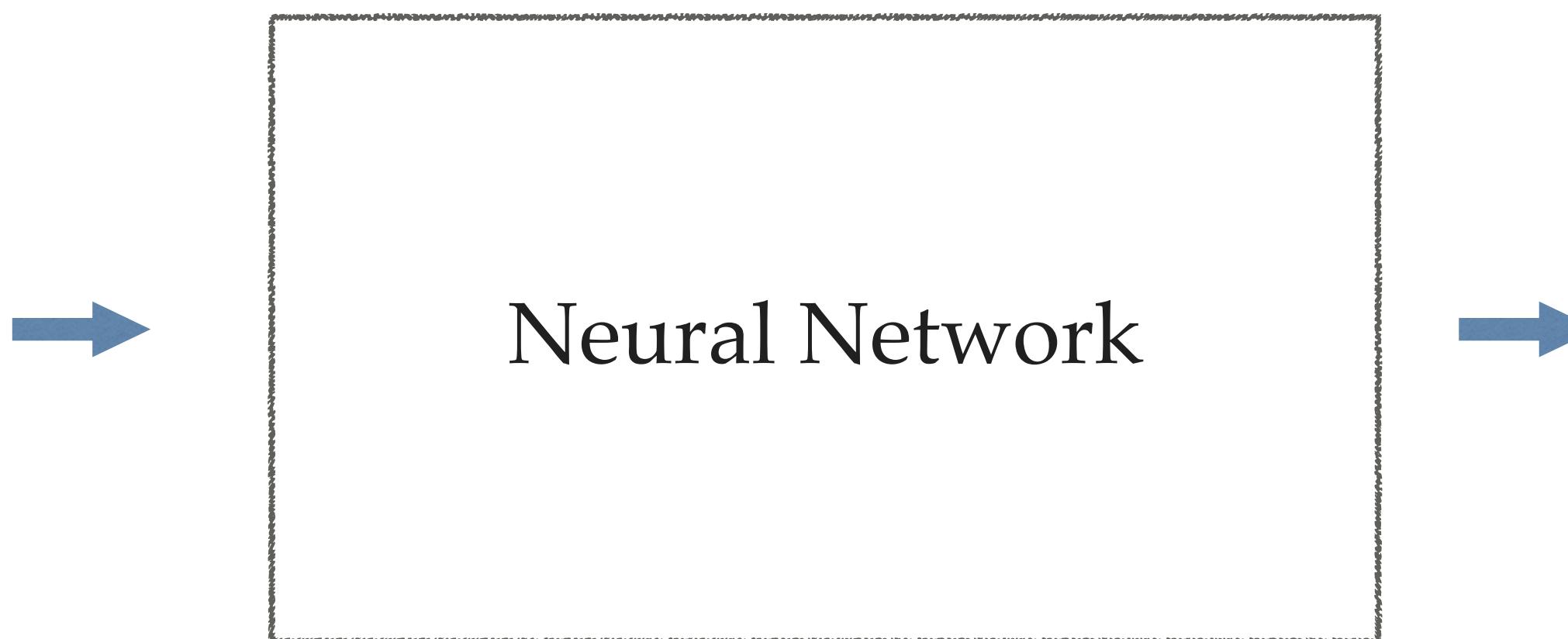
WORD2VEC

- We can estimate the empirical probability as $P(y|x) \propto \exp(\mathbf{y} \cdot \mathbf{x})$
- \mathbf{y} is a vector representation of a word when it appears in the context; \mathbf{x} is a vector representation of a word when it appears in the input

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

WORD2VEC

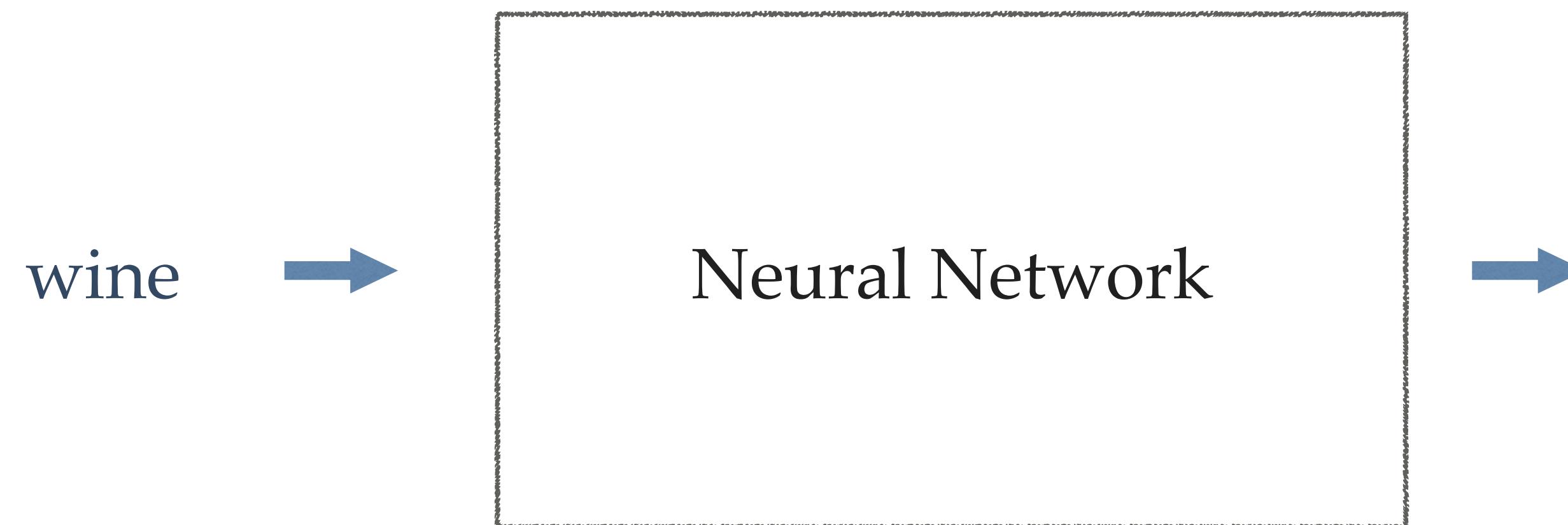
x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...



Take every entry on the left and guess the corresponding entry on the right. Update the network if guess is wrong.

WORD2VEC

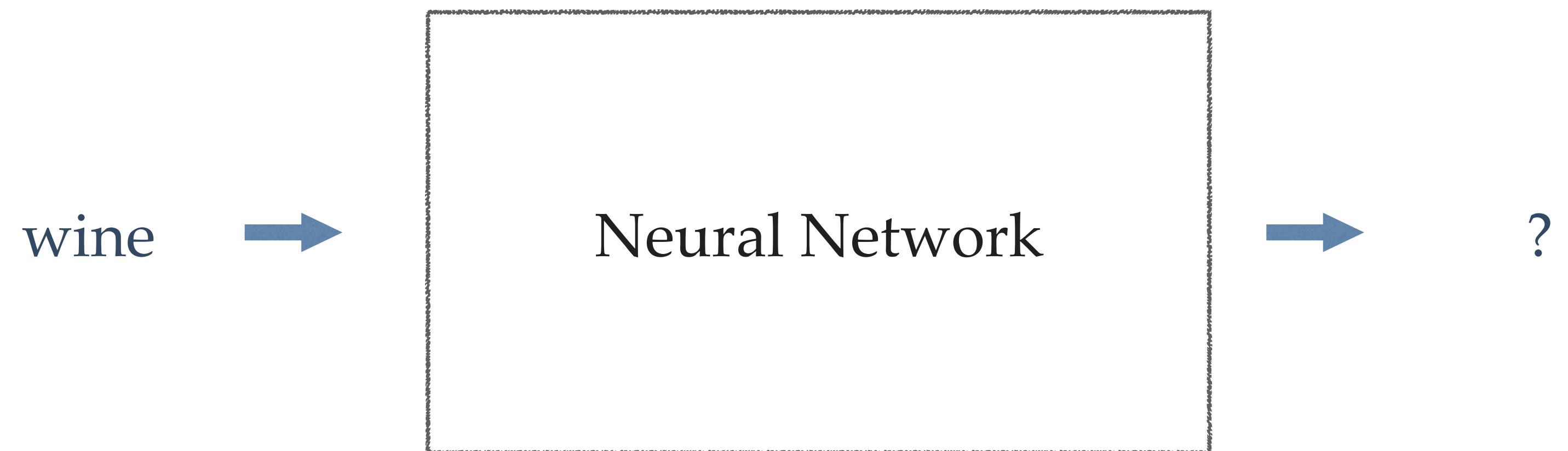
x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...



Take every entry on the left and guess the corresponding entry on the right. Update the network if guess is wrong.

WORD2VEC

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...



Take every entry on the left and guess the corresponding entry on the right. Update the network if guess is wrong.

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

Vocab

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

wine

Vocab

x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

Vocab

X

0

wine

1

0

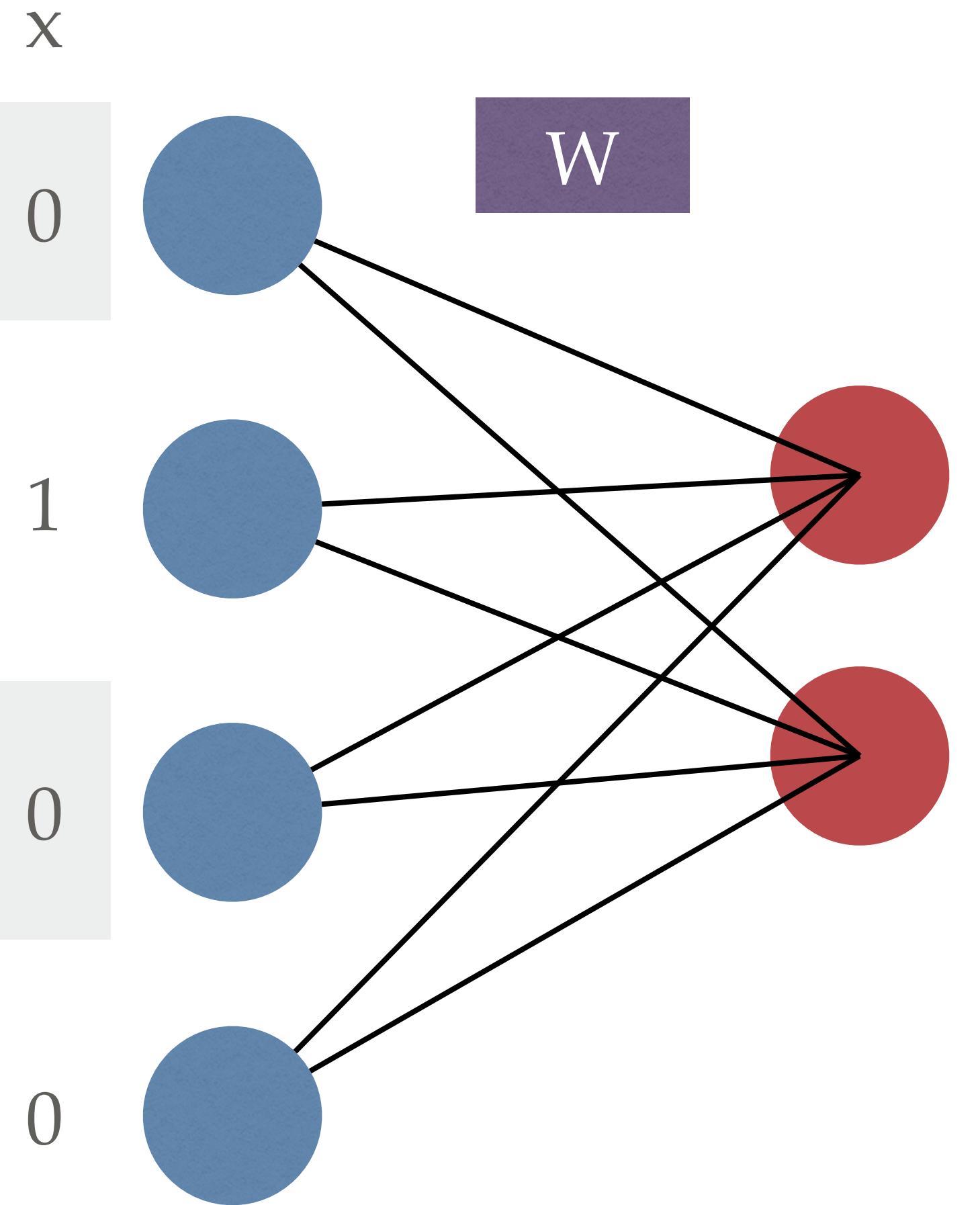
0

x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

<i>service</i>
<i>wine</i>
<i>cold</i>
<i>great</i>

Vocab

wine

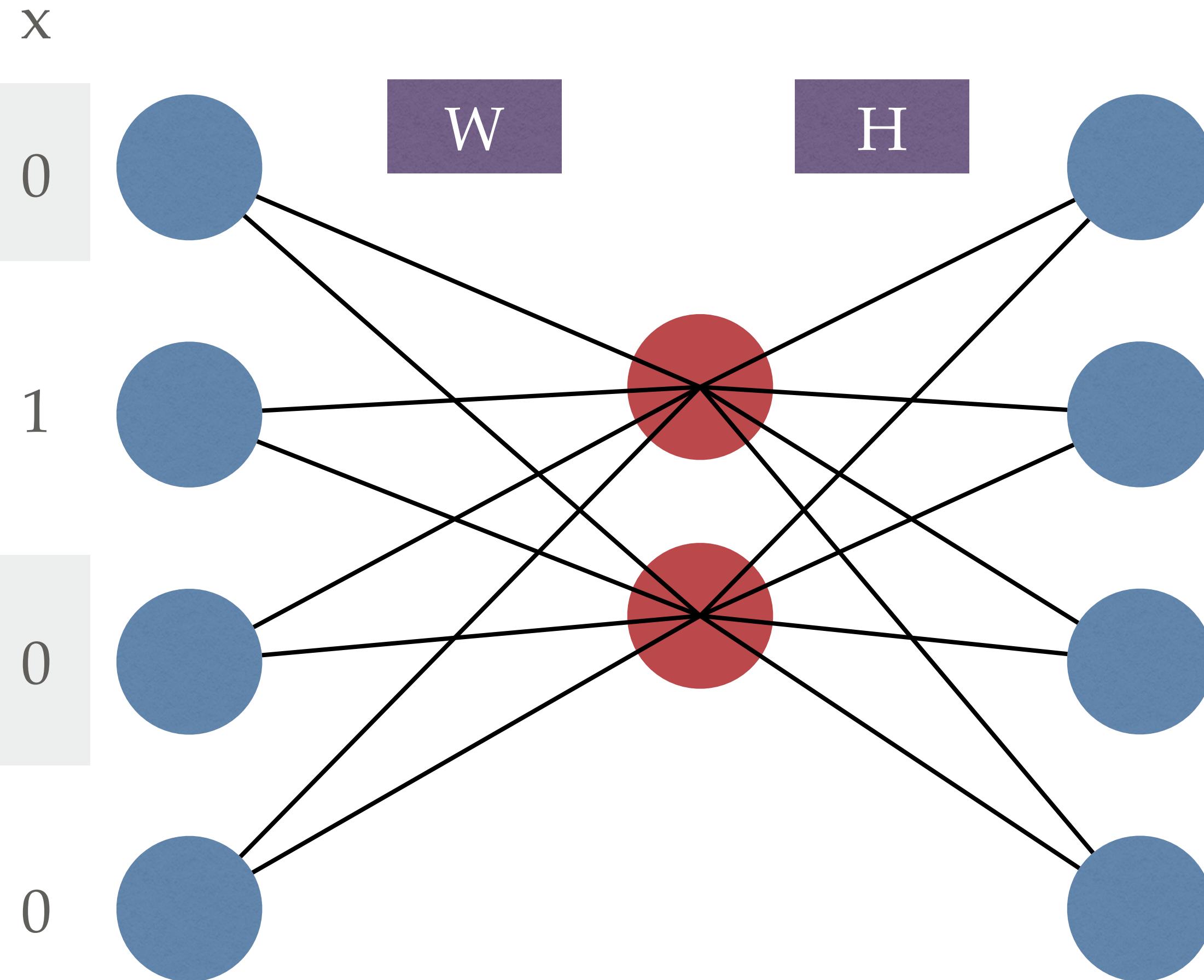


x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

service	
wine	
cold	
great	

Vocab

wine

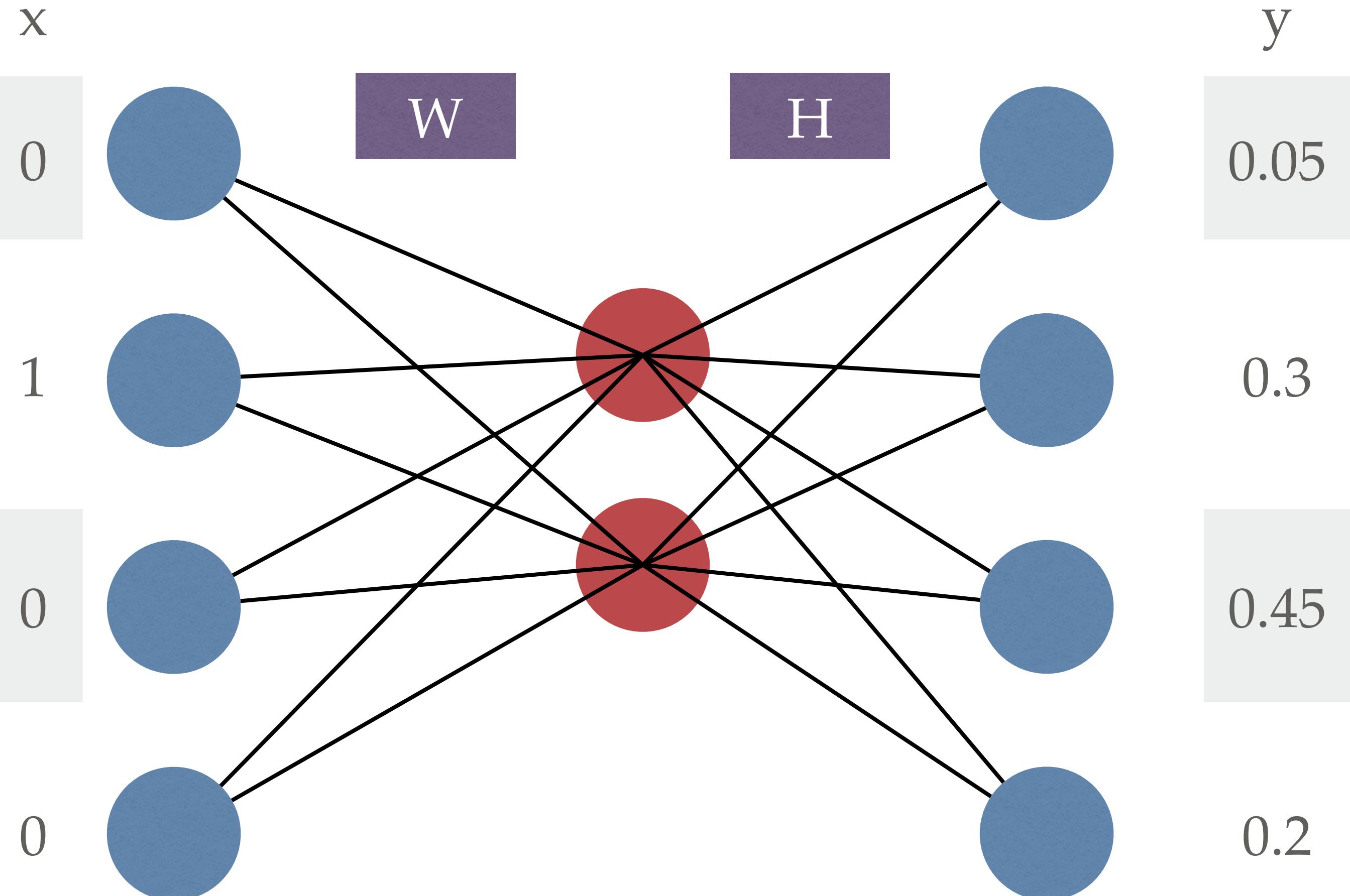


x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

service
wine
cold
great

Vocab

wine

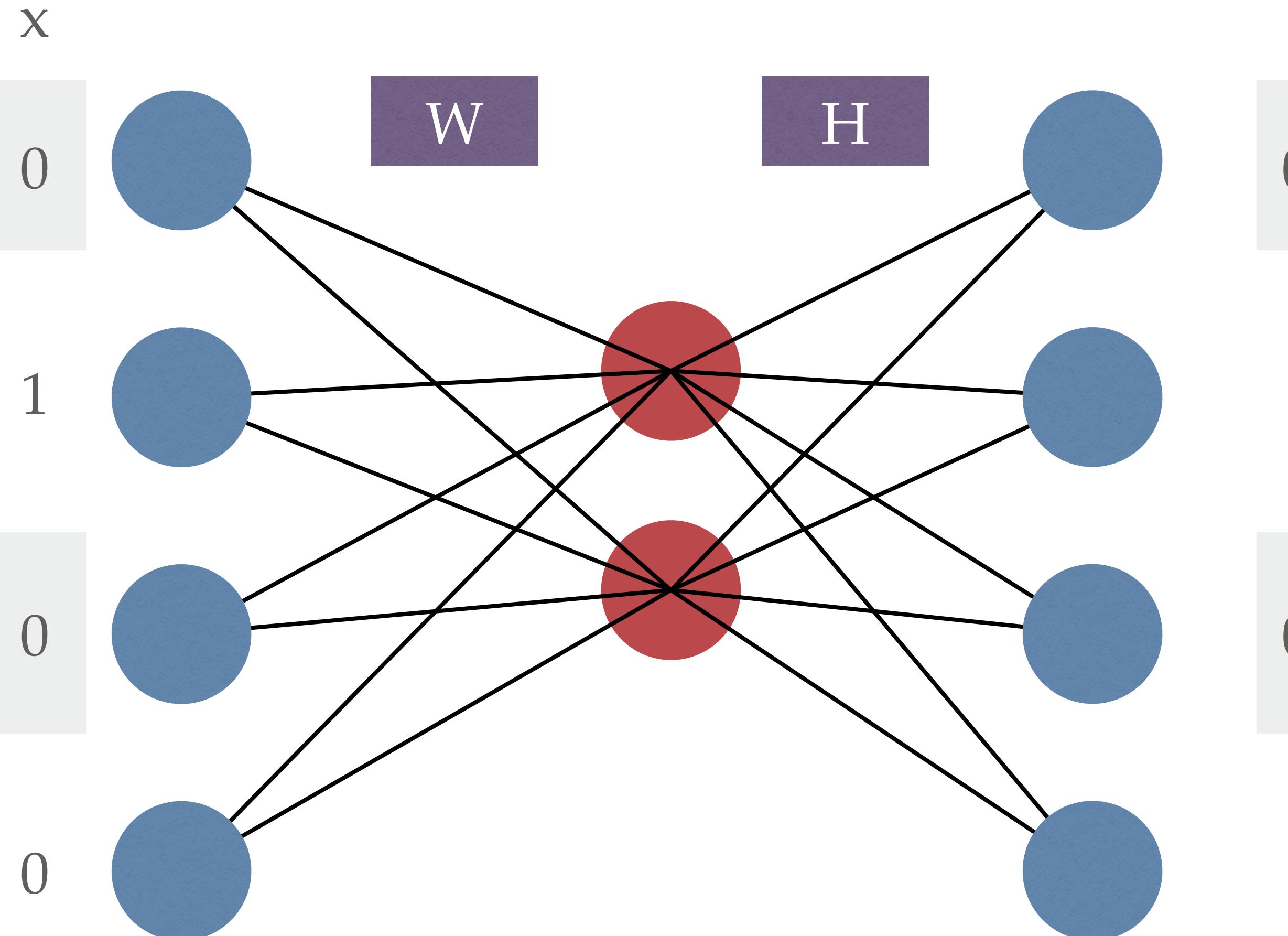


x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

<i>service</i>	
<i>wine</i>	
<i>cold</i>	
<i>great</i>	

Vocab

wine



x	y
...	...
<i>wine</i>	<i>cold</i>
<i>wine</i>	<i>sweet</i>
<i>wine</i>	<i>spirit</i>
<i>wine</i>	<i>drink</i>
...	...

service

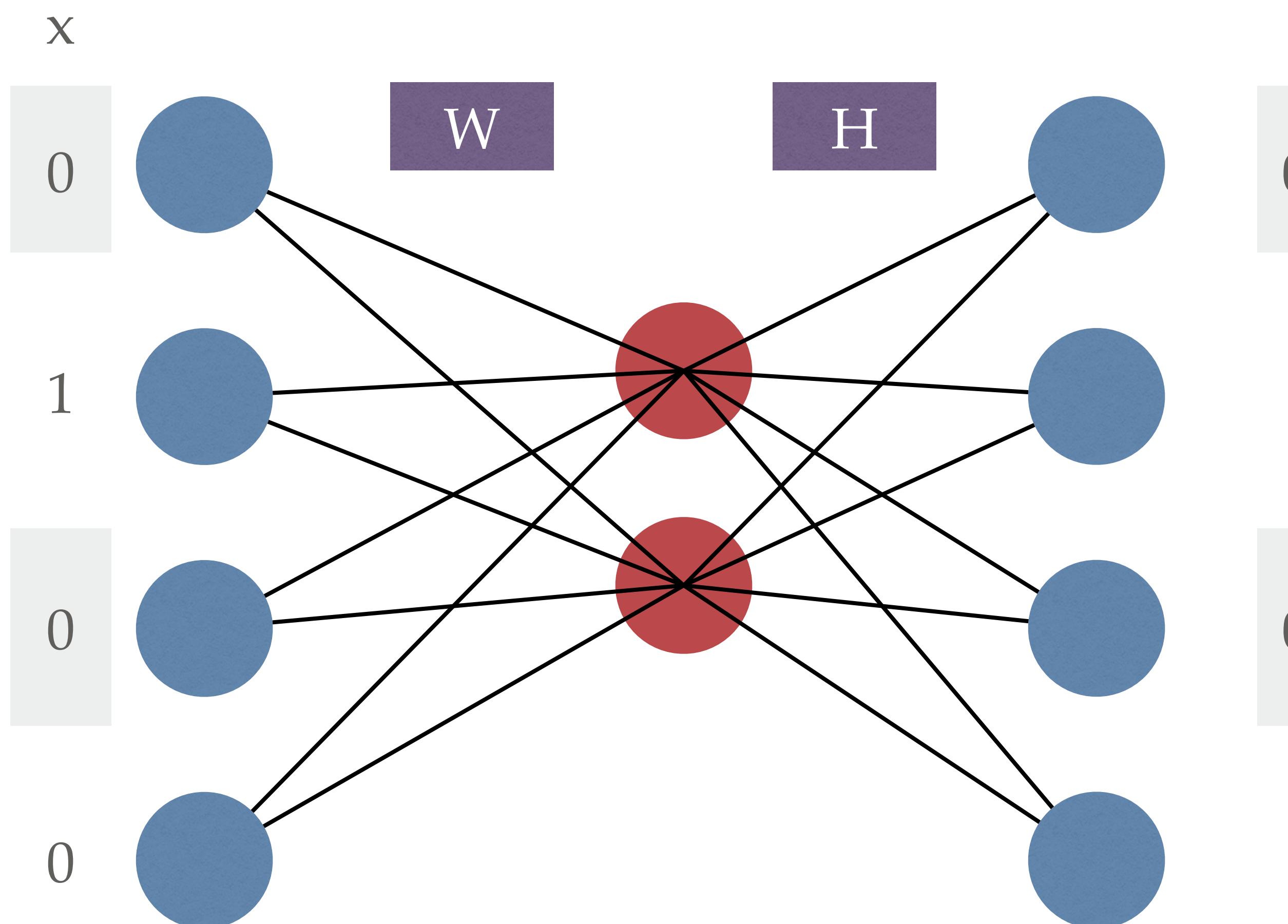
wine

cold

great

Vocab

wine



Word embeddings as columns

W

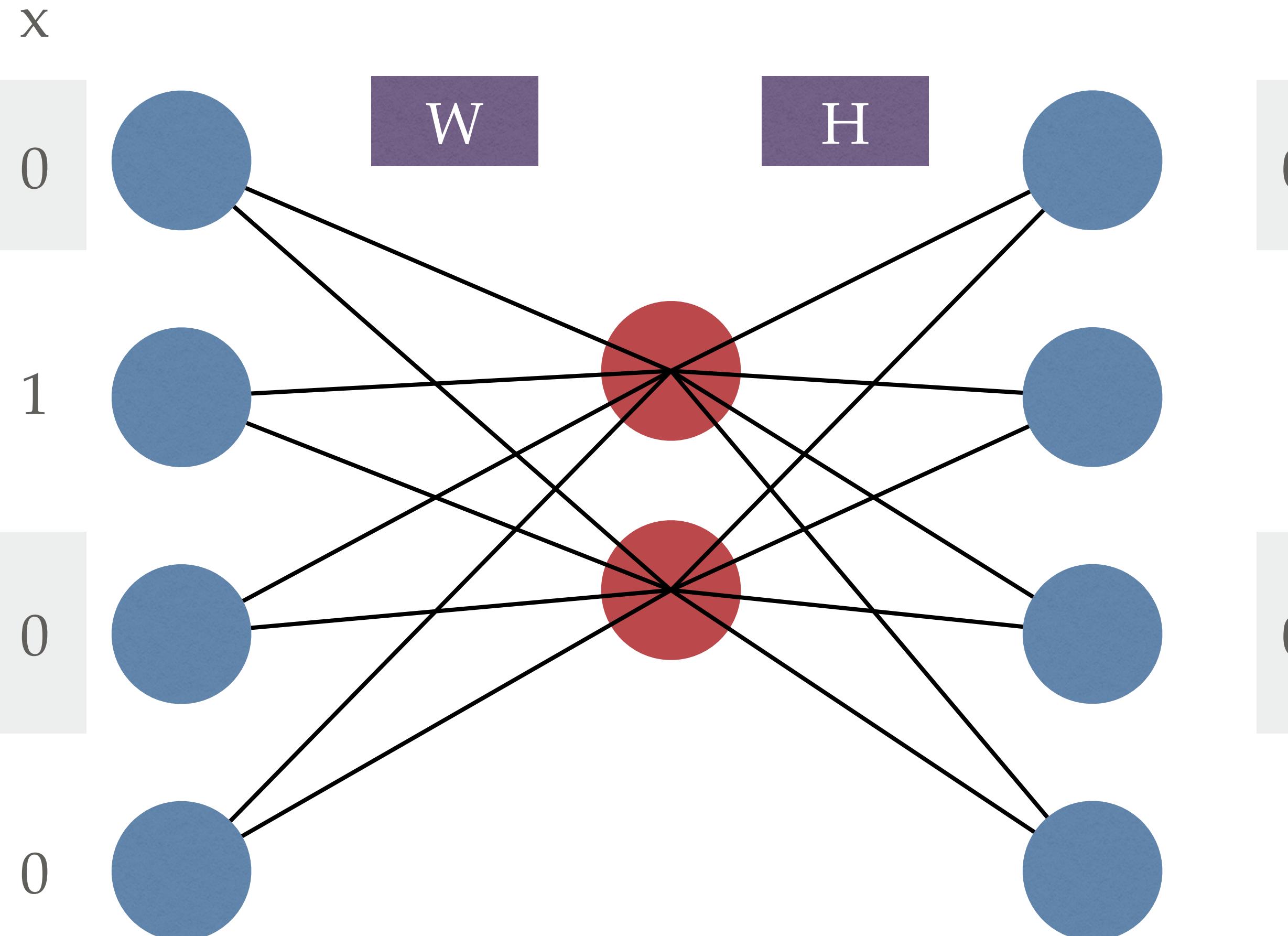
-0.3	1.2	0.5	-0.6
0.2	0.9	0.1	-0.4

x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

	x
service	0
wine	1
cold	0
great	0

Vocab

wine



x	y
...	...
wine	cold
wine	sweet
wine	spirit
wine	drink
...	...

Word embeddings as columns

W

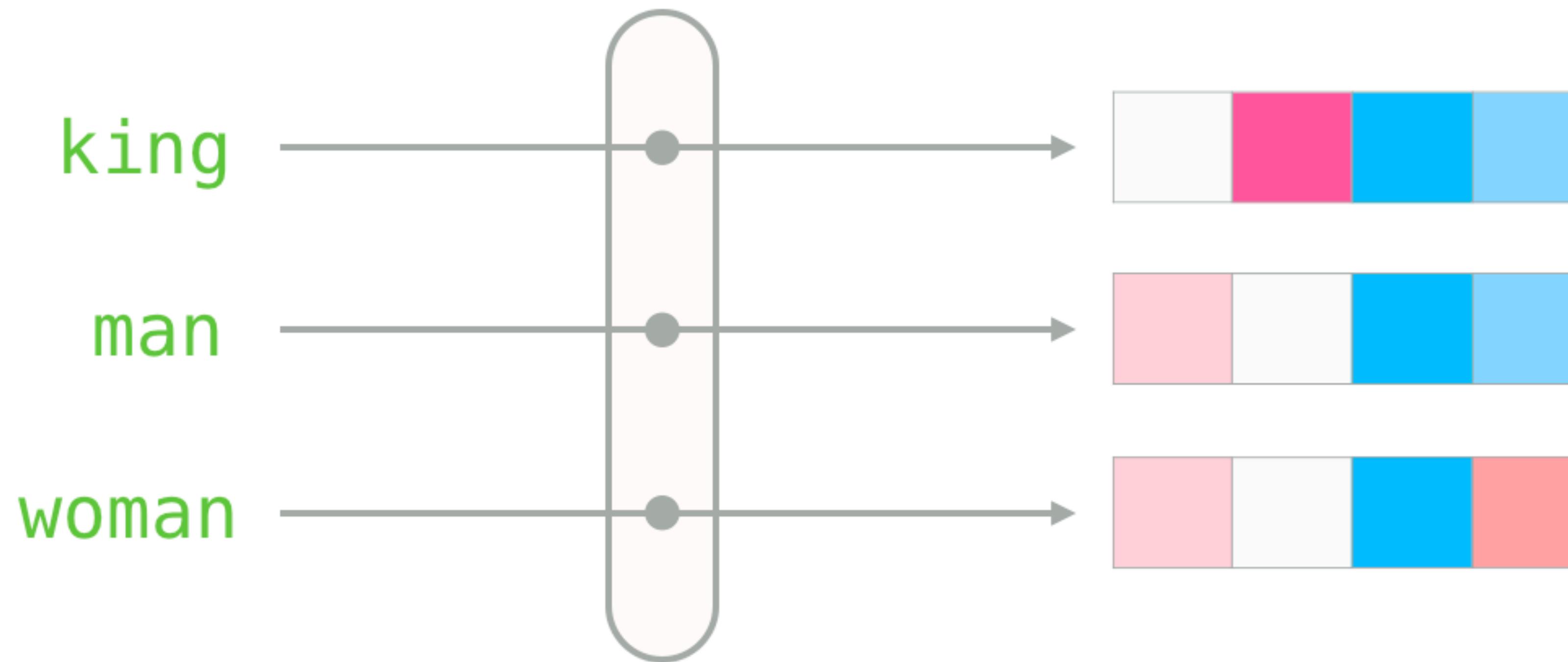
-0.3	1.2	0.5	-0.6
0.2	0.9	0.1	-0.4

Context embeddings as rows

H

0.1	-0.4
0.4	-0.5
0.3	-0.1
0.2	0.1

Word2vec



<http://jalammar.github.io/illustrated-word2vec/>

Trait 2

King

Queen

Trait 1

Man

Woman



GLOVE

- Pennington et. al. 2014 used a similar idea but directly model some statistic of the cooccurrence
- Think of this as a regression task

x	y	statistic
wine	cold	3
wine	sweet	5
wine	spirit	1
wine	drink	2

EVALUATION

EVALUATION

- How do we know if the embeddings we learned are any good?

EVALUATION

EVALUATION

- Intrinsic evaluation
 - Word relatedness: the similarity between vector representations should correlate with human judgments of relatedness of pairs of words
 - Analogical reasoning: King:Queen::Man:?
- Extrinsic evaluation
 - Plug in the embeddings as features in some downstream task

IN CLASS

- Word2Vec demo