



Lesson4:
Descriptive Modelling of Similarity of Text
Unit3:
Vector space models for similarity

Rene Pickhardt

Introduction to Web Science Part 2
Emerging Web Properties





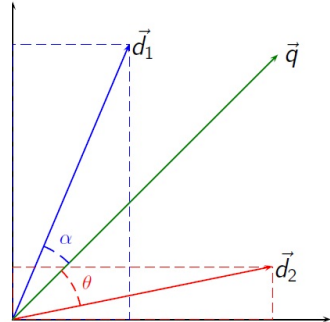
Completing this unit you should ...

- Be familiar with the the vector space model for text documents
- Be aware of term frequency and (inverse) document frequency
- Have reviewed the definitions of base and dimension
- Realize that the angle between two vectors can be seen as a similarity measure

A model based on vector spaces



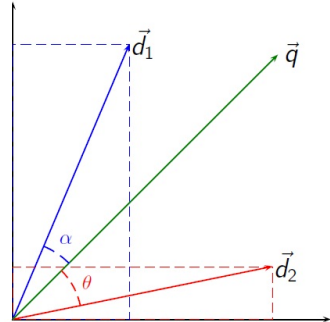
1) Model documents as vectors of words



A model based on vector spaces



1) Model documents as vectors of words



Vector space Model

2) calculating

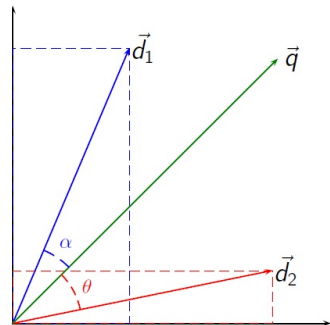
Distance between vectors

A model based on vector spaces



Matthew	Luke
δων δε πολλους των	ελεγεν ουν τοις
φαρισαιων και σαδδουκαιων	εκπορευομενοις οχλοις
ερχομενους επι το βαπτισμα	βαπτισθηται υπ
αυτου ειπεν αυτοις	αυτου
γεννηματα εχιδων τις	γεννηματα εχιδων τις
υπεδειξεν υμιν φυγειν απο	υπεδειξεν υμιν φυγειν απο
της μελλουσης οργης	της μελλουσης οργης
ποιησατε ουν καρπους	ποιησατε ουν καρπους
αξιους της μετανοιας	αξιους της μετανοιας
και μη δοξητε λεγειν εν	και μη αρξησθε λεγειν εν
εαυτοις πατερα εχομεν τον	εαυτοις πατερα εχομεν τον
αβρααμ λεγω γαρ υμιν οτι	αβρααμ λεγω γαρ υμιν οτι
δυναται ο θεος εκ των λιθων	δυναται ο θεος εκ των λιθων
τουτων εγειραι τεκνα του	τουτων εγειραι τεκνα του
αβρααμ ηδη δε και η αβνη	αβρααμ ηδη δε και η αβνη
προς την ριζαν των δενδριων	προς την ριζαν των δενδριων
κειται παν ουν δενδρον μη	κειται παν ουν δενδρον μη
ποιουν καρπον καλον	ποιουν καρπον καλον
εκκοπεται και εις πυρ	εκκοπεται και εις πυρ
βαλλεται	βαλλεται

1) Model documents as vectors of words



3) interpreting

2) calculating

Vector space Model

Distance between vectors

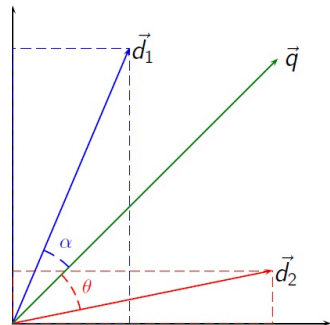
A model based on vector spaces



Be able to make statements about similarity of documents

Matthew	Luke
δὼν δὲ πολλοὺς τῶν	εἶπεν οὖν τοῖς
φαρισαίων καὶ σαδδουκαίων	ἐκπορευομένοις σχολαῖς
ἐρχομένους ἐπὶ τὸ βάπτισμα	βαπτισθῆναι ὑπὲρ
αὐτοῦ εἶπεν αὐτοῖς	αὐτοῦ
γεννητὰ ἐχθρῶν τῆς	γεννητὰ ἐχθρῶν τῆς
ὑπεδείξεν ὑμῖν φυγεῖν ἀπὸ	ὑπεδείξεν ὑμῖν φυγεῖν ἀπὸ
τῆς μελλούσης ὀργῆς	τῆς μελλούσης ὀργῆς
ποιήσατε οὖν καρποὺς	ποιήσατε οὖν καρποὺς
ἀδίκους τῆς μετανάστευσης	ἀδίκους τῆς μετανάστευσης
καὶ μὴ δοξάτε λέγειν ἐν	καὶ μὴ ἀρξήσθε λέγειν ἐν
ἐαυτοῖς πατέρα ἔχομεν τὸν	ἐαυτοῖς πατέρα ἔχομεν τὸν
ἀβραάμ· λέγω γὰρ ὑμῖν οὐ	ἀβραάμ· λέγω γὰρ ὑμῖν οὐ
δυνατὸν ὁ θεὸς ἐκ τῶν λίθων	δυνατὸν ὁ θεὸς ἐκ τῶν λίθων
τοῦτων ἐγενεῖται τέκνα τῶν	τοῦτων ἐγενεῖται τέκνα τῶν
ἀβραάμ· ἤδη δὲ καὶ ἡ ἀσὴν	ἀβραάμ· ἤδη δὲ καὶ ἡ ἀσὴν
πρὸς τὴν ρίζαν τῶν δένδρων	πρὸς τὴν ρίζαν τῶν δένδρων
κεῖται παν οὖν δένδρον μὴ	κεῖται παν οὖν δένδρον μὴ
ποιῶν καρπὸν καλόν	ποιῶν καρπὸν καλόν
ἐκκοπτεται καὶ εἰς τὴν	ἐκκοπτεται καὶ εἰς τὴν
βάλλεται	βάλλεται

1) Model documents as vectors of words



3) interpreting

2) calculating

Vector space Model

Distance between vectors

Pay close attention to notation

w_i  D_j

- Words $W = \{w_1, w_2, \dots, w_n\}$

- Word Vectors $V = \langle \vec{w}_1, \vec{w}_2, \dots, \vec{w}_n \rangle$



- Document $D_j \in W^*$ are a sequence of words
 - So $D_j = w_{i_1} w_{i_2} \dots w_{i_m}$ has a length of m

- Document vector $\vec{d}_j = \sum_{i=1}^n tf(w_i, D_j) \vec{w}_i$



Usually tf-idf is considered instead of tf!

- The document frequency is defined as
$$df(w_i) = |\{D_j | w_i \text{ in } D_j\}|$$
- Inverse document frequency is defined as

- $idf(w_i) = \log \frac{|D|}{df(w_i)}$ resulting in

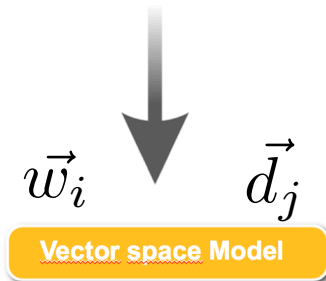
$$tfidf(w_i, D_j) = tf(w_i, D_j) \times \log \frac{|D|}{df(w_i)}$$

- In the videos and slides for simplicity of numbers we will only use the term frequency

Choose a vector space and base

- Let $V = \langle \vec{a}, \vec{b} \rangle$ be the vector space spanned by the words “a” and “b”

w_i  D_j



- $\vec{a} = \begin{pmatrix} 1 \\ 0 \end{pmatrix}$ is a base vector for word “a”
- Similarly for word “b” we have $\vec{b} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$
- Choosing the base was a modelling choice!

Calculate the modelled document vectors d_1

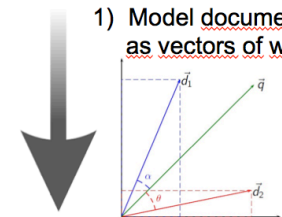
- $D_1 = a a a b b a a b a b$

- $tf(a, D_1) = 6$

- $tf(b, D_1) = 4$



1) Model documents as vectors of words



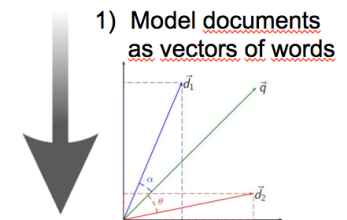
- Let us now create the document vectors

$$\vec{d}_1 = \sum_{i=1}^2 tf(w_i, D_1) \vec{w}_i = tf(a, D_1) \vec{a} + tf(b, D_1) \vec{b}$$

$$= 6 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 4 \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 6 \\ 4 \end{pmatrix}$$

Calculate the modelled document vectors d_2

- $D_2 = a \ b \ b$
 - $tf(a, D_2) = 1$
 - $tf(b, D_2) = 2$



- Let us now create the document vectors

$$\vec{d}_2 = \sum_{i=1}^2 tf(w_i, D_2) \vec{w}_i = tf(a, D_2) \vec{a} + tf(b, D_2) \vec{b}$$

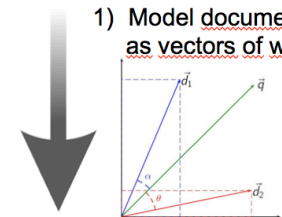
$$= 1 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 2 \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$$

Calculate the modelled document vectors d_3

- $D_3 = a a b$
 - $tf(a, D_3) = 2$
 - $tf(b, D_3) = 1$



1) Model documents as vectors of words



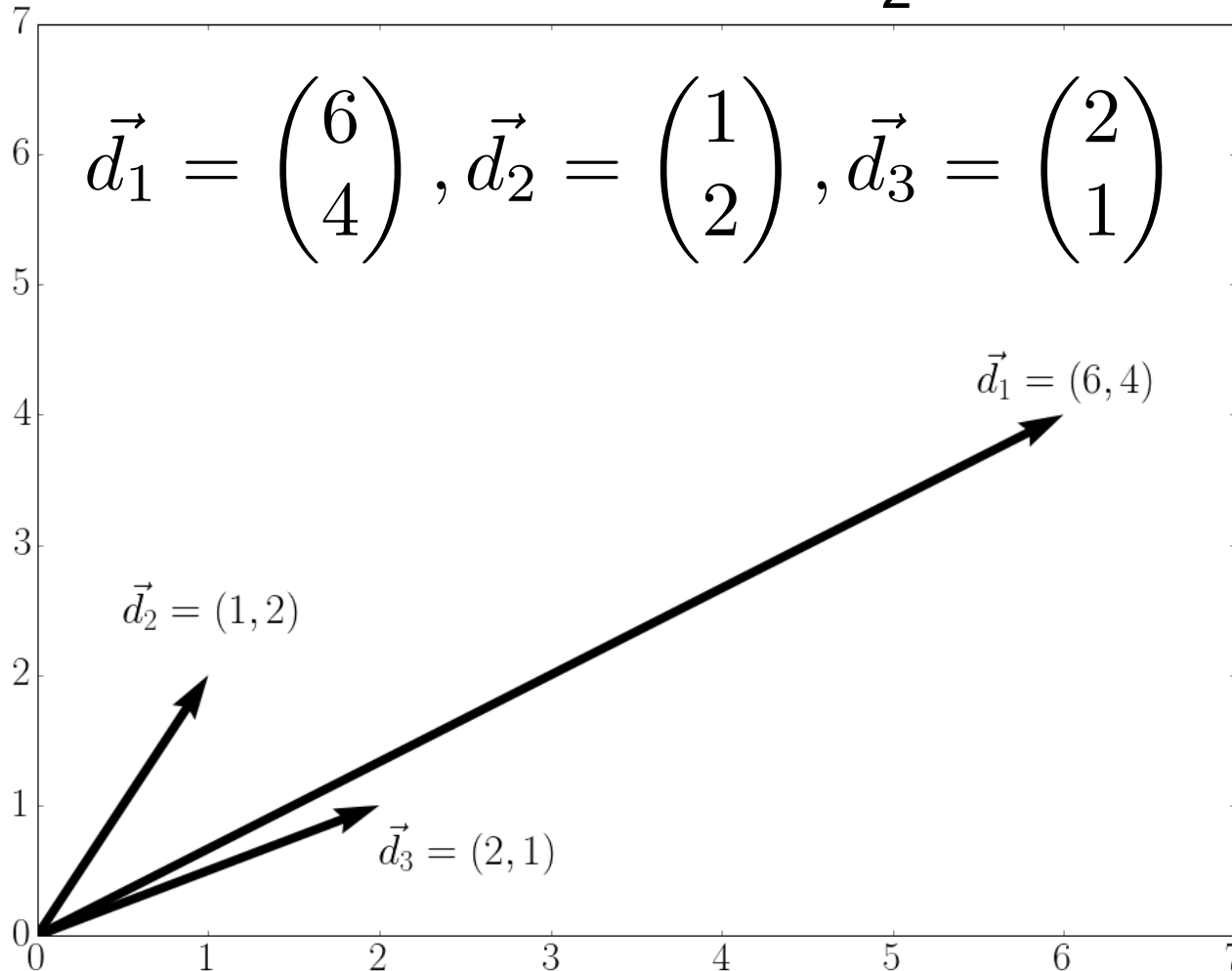
- Let us now create the document vectors

$$\vec{d}_3 = \sum_{i=1}^2 tf(w_i, D_3) \vec{w}_i = tf(a, D_3) \vec{a} + tf(b, D_3) \vec{b}$$

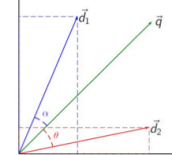
$$= 2 \begin{pmatrix} 1 \\ 0 \end{pmatrix} + 1 \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 2 \\ 1 \end{pmatrix}$$

Digression: Drawing the document vectors

- $D_1 = a a a b b a a b a b$, $D_2 = a b b$, $D_3 = a a b$

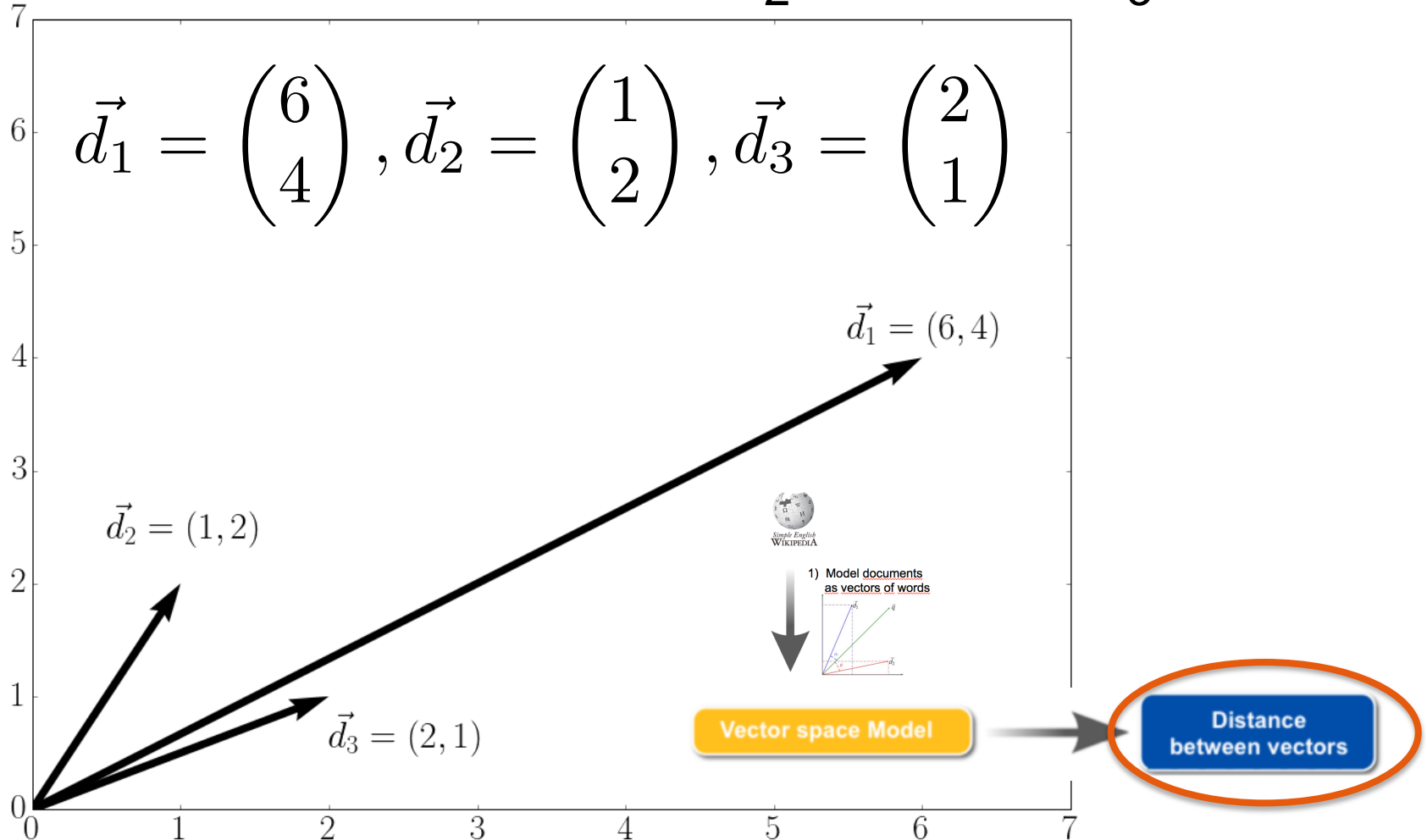


1) Model documents as vectors of words



Which vectors are closest too each other?

- $D_1 = a a a b b a a b a b$, $D_2 = a b b$, $D_3 = a a b$





Two ways of calculating the distance between two vectors d_i and d_j ?

- **Euclidean distance**
 - “Take a rule and measure”
 - Take difference: $\vec{d} = \vec{d}_i - \vec{d}_j$
 - Calculate length of difference $||\vec{d}||$

- **Cosine distance**
 - “Take the angle between vectors”



Calculate Euclidean distance between two document vectors d_i and d_j

- Take the difference $\vec{d} = \vec{d}_i - \vec{d}_j$
- Calculate the length of the difference

$$\|\vec{d}\|^2 = \sum_{k=1}^n (d_k)^2 = \sum_{k=1}^n ((\vec{d}_i)_k - (\vec{d}_j)_k)^2$$



Calculate Euclidean distance between two document vectors d_i and d_j

- Take the difference $\vec{d} = \vec{d}_i - \vec{d}_j$
- Calculate the length of the difference

$$\|\vec{d}\|^2 = \sum_{k=1}^n (d_k)^2 = \sum_{k=1}^n ((\vec{d}_i)_k - (\vec{d}_j)_k)^2$$

The k-th component of vector \vec{d}_i
 $(\vec{d}_i)_k = tf(w_k, D_i)$ by definition



Calculate Euclidean distance between two document vectors d_i and d_j

- Take the difference $\vec{d} = \vec{d}_i - \vec{d}_j$
- Calculate the length of the difference

$$\begin{aligned}\|\vec{d}\|^2 &= \sum_{k=1}^n (d_k)^2 = \sum_{k=1}^n ((\vec{d}_i)_k - (\vec{d}_j)_k)^2 \\ &= \sum_{k=1}^n (tf(w_k, D_i) - tf(w_k, D_j))^2\end{aligned}$$

- For every word w_k we compare how often it appears in document D_i and document D_j



Calculate Euclidean distance between two document vectors d_i and d_j

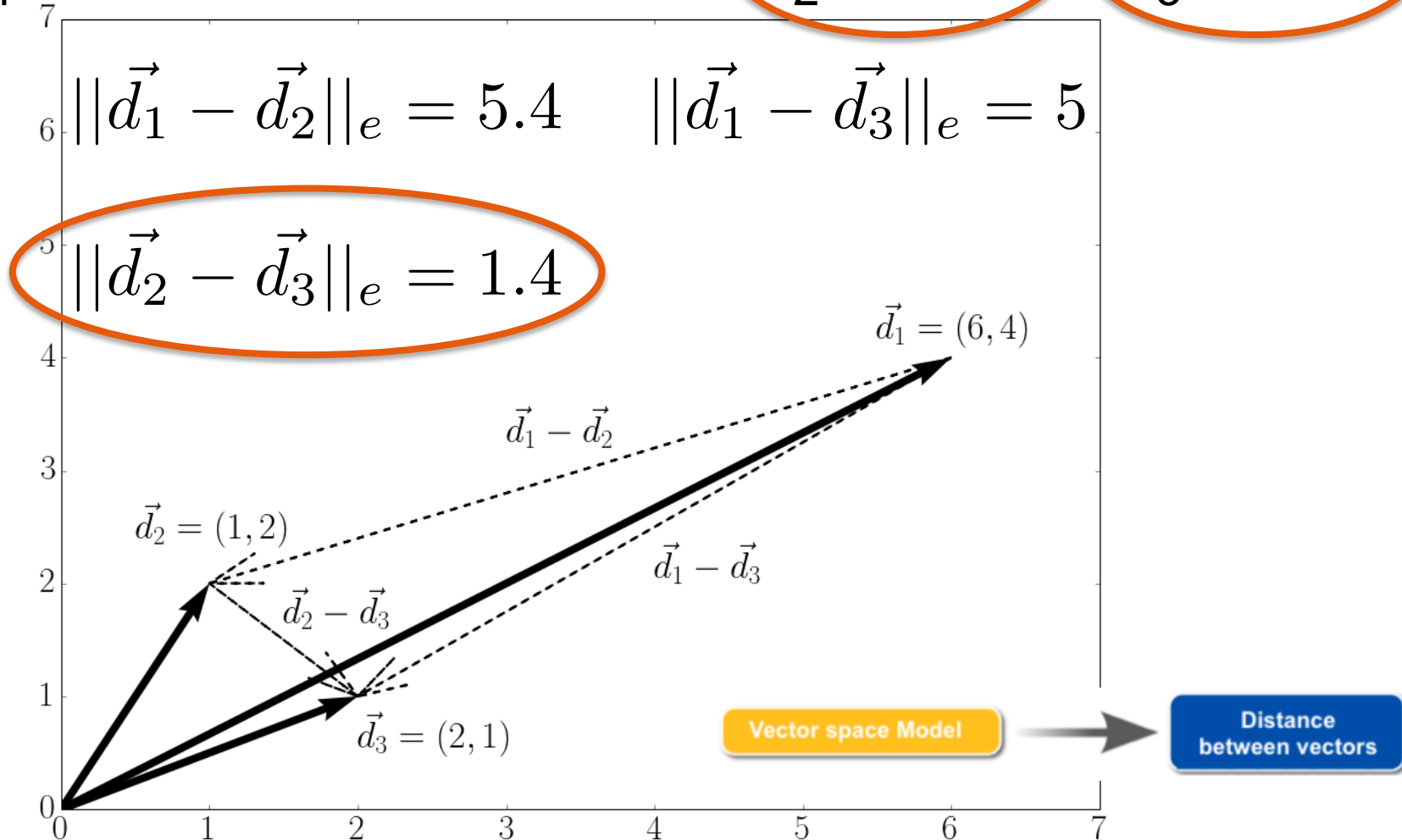
- Take the difference $\vec{d} = \vec{d}_i - \vec{d}_j$
- Calculate the length of the difference

$$\begin{aligned}\|\vec{d}\|^2 &= \sum_{k=1}^n (d_k)^2 = \sum_{k=1}^n ((\vec{d}_i)_k - (\vec{d}_j)_k)^2 \\ &= \sum_{k=1}^n (tf(w_k, D_i) - tf(w_k, D_j))^2\end{aligned}$$

- For every word w_k we compare how often it appears in document D_i and document D_j

Euclidean distances for our example

- $D_1 = a a a b b a a b a b$, $D_2 = a b b$, $D_3 = a a b$





Calculate Cosine distance between two document vectors \vec{d}_i and \vec{d}_j

- Calculate the scalar product $s = \langle \vec{d}_i, \vec{d}_j \rangle$

$$\begin{aligned}\langle \vec{d}_i, \vec{d}_j \rangle &= \sum_{k=1}^n (\vec{d}_i)_k (\vec{d}_j)_k \\ &= \sum_{k=1}^n tf(w_k, D_i) tf(w_k, D_j)\end{aligned}$$

Remember: $(\vec{d}_i)_k = tf(w_k, D_i)$
is zero most of the time.



Calculate Cosine distance between two document vectors d_i and d_j

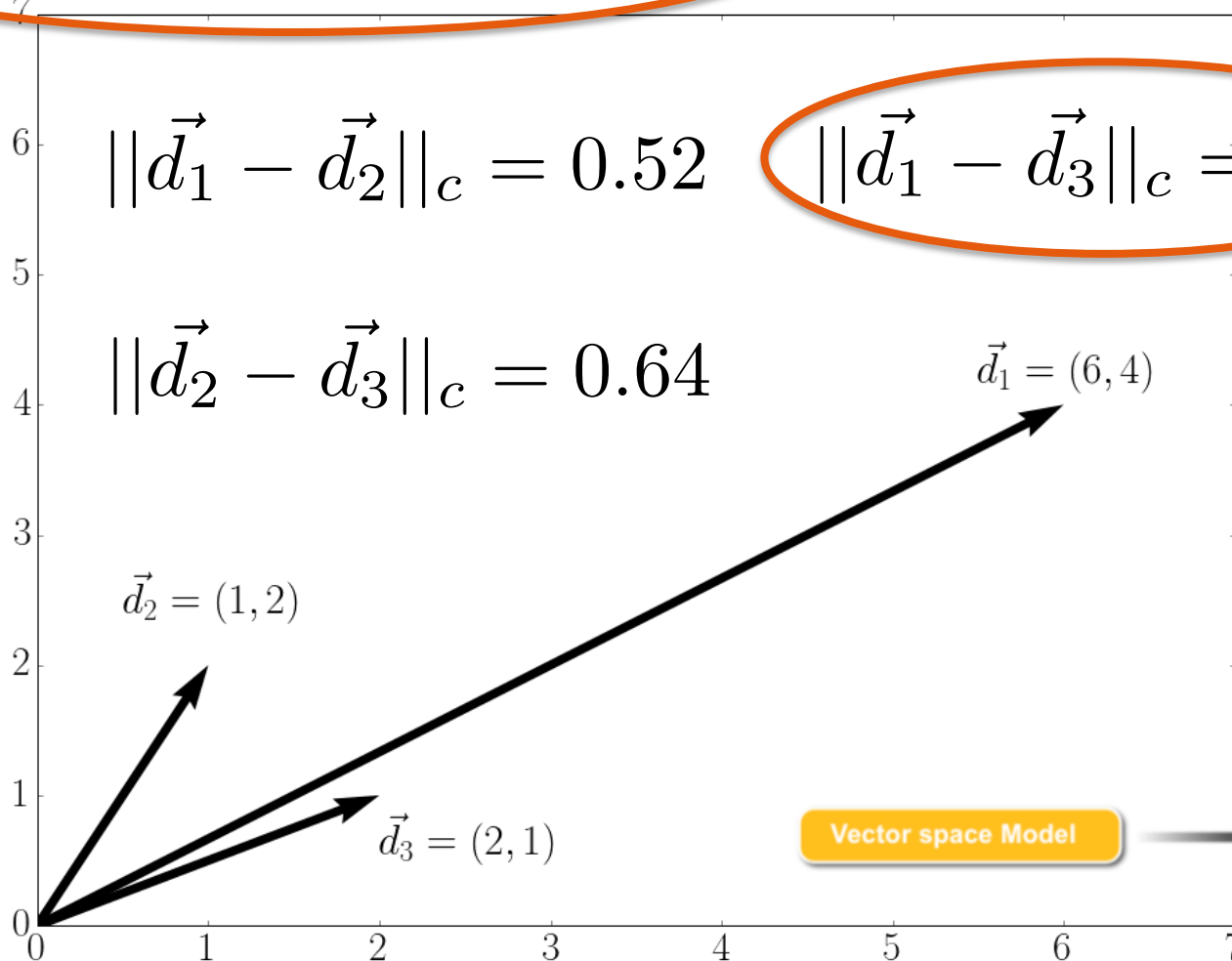
- Calculate the scalar product $s = \langle \vec{d}_i, \vec{d}_j \rangle$
- Divide it by the product of lengths of both vectors (length as in Euclidean distance)

$$\cos(\theta) = \frac{\langle \vec{d}_i, \vec{d}_j \rangle}{\|\vec{d}_i\| * \|\vec{d}_j\|}$$

$$\Rightarrow \theta = \cos^{-1} \left(\frac{\langle \vec{d}_i, \vec{d}_j \rangle}{\|\vec{d}_i\| * \|\vec{d}_j\|} \right)$$

Cosine distances for our example

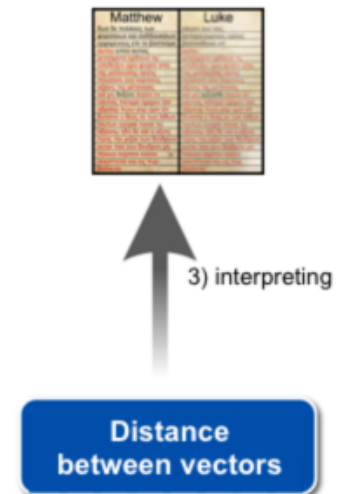
- $D_1 = a a a b b a a b a b$, $D_2 = a b b$, $D_3 = a a b$



Comparing cosine and Euclidean distance

- $D_1 = a a a b b a a b a b$, $D_2 = a b b$, $D_3 = a a b$

	Euklid	Cosine
$ \vec{d}_1 - \vec{d}_2 $	5.4	0.52
$ \vec{d}_2 - \vec{d}_3 $	1.4	0.64
$ \vec{d}_1 - \vec{d}_3 $	5	0.12



- Different choices of metric can yield very different results
- Choice of metric is part of the model!
- **Usually cosine distance is considered**



Thank you for your attention!



Contact:

Rene Pickhardt
Institute for Web Science and Technologies
Universität Koblenz-Landau
rpickhardt@uni-koblenz.de

WeST 
People and Knowledge Networks



Copyright:

- This Slide deck is licensed under creative commons 3.0. share alike attribution license. It was created by Rene Pickhardt. You can use share and modify this slide deck as long as you attribute the author and keep the same license.
- https://commons.wikimedia.org/wiki/File:Synoptic_word-for-word.png By Alecmconroy (Own work) [GFDL (<http://www.gnu.org/copyleft/fdl.html>) or CC BY 3.0 (<http://creativecommons.org/licenses/by/3.0/>)], via Wikimedia Commons
- <https://commons.wikimedia.org/wiki/File:Inner-product-angle.png> CC-BY-SA by CSTAR & Oleg Alexandrov