

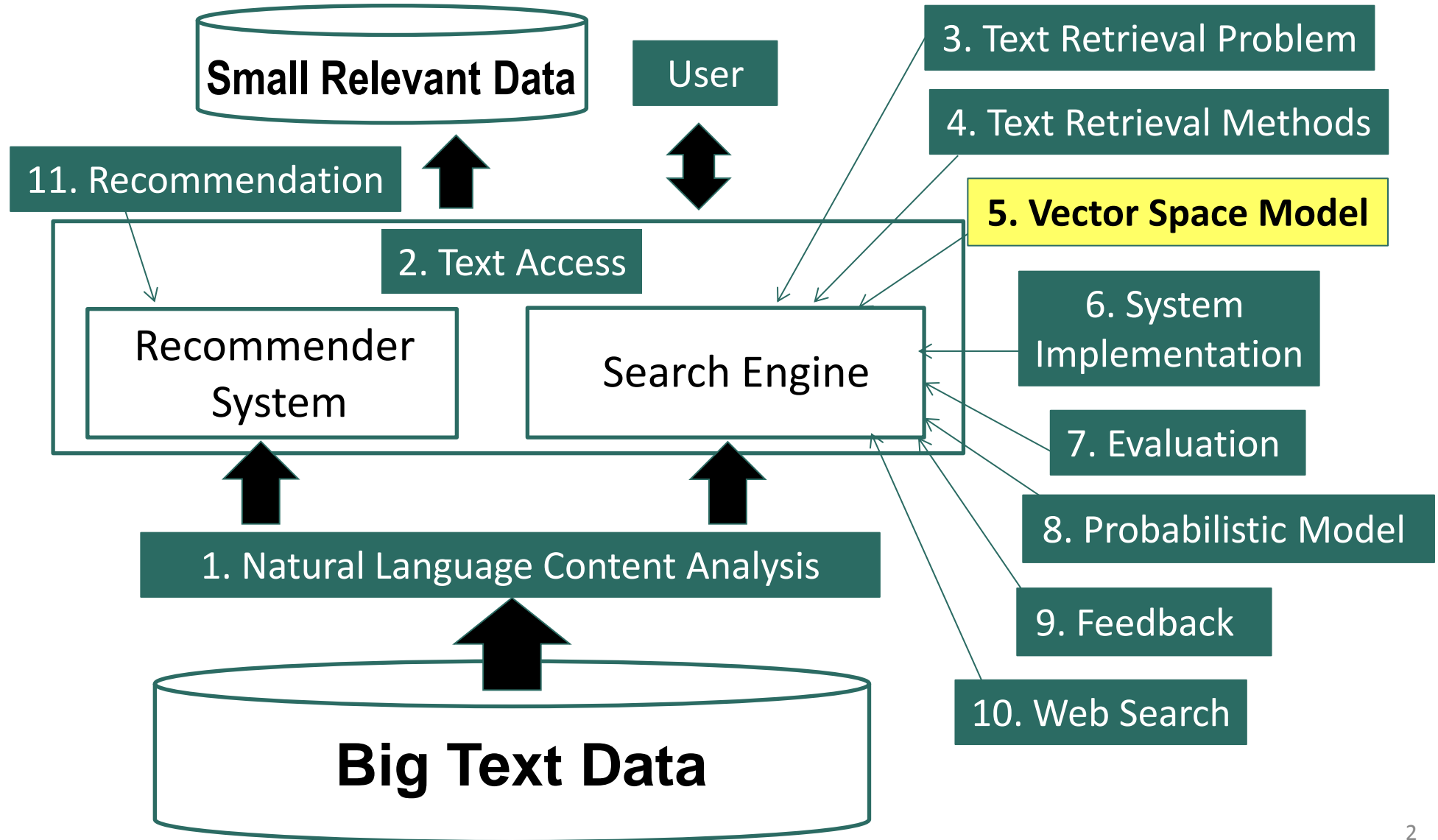


Text Retrieval and Search Engines

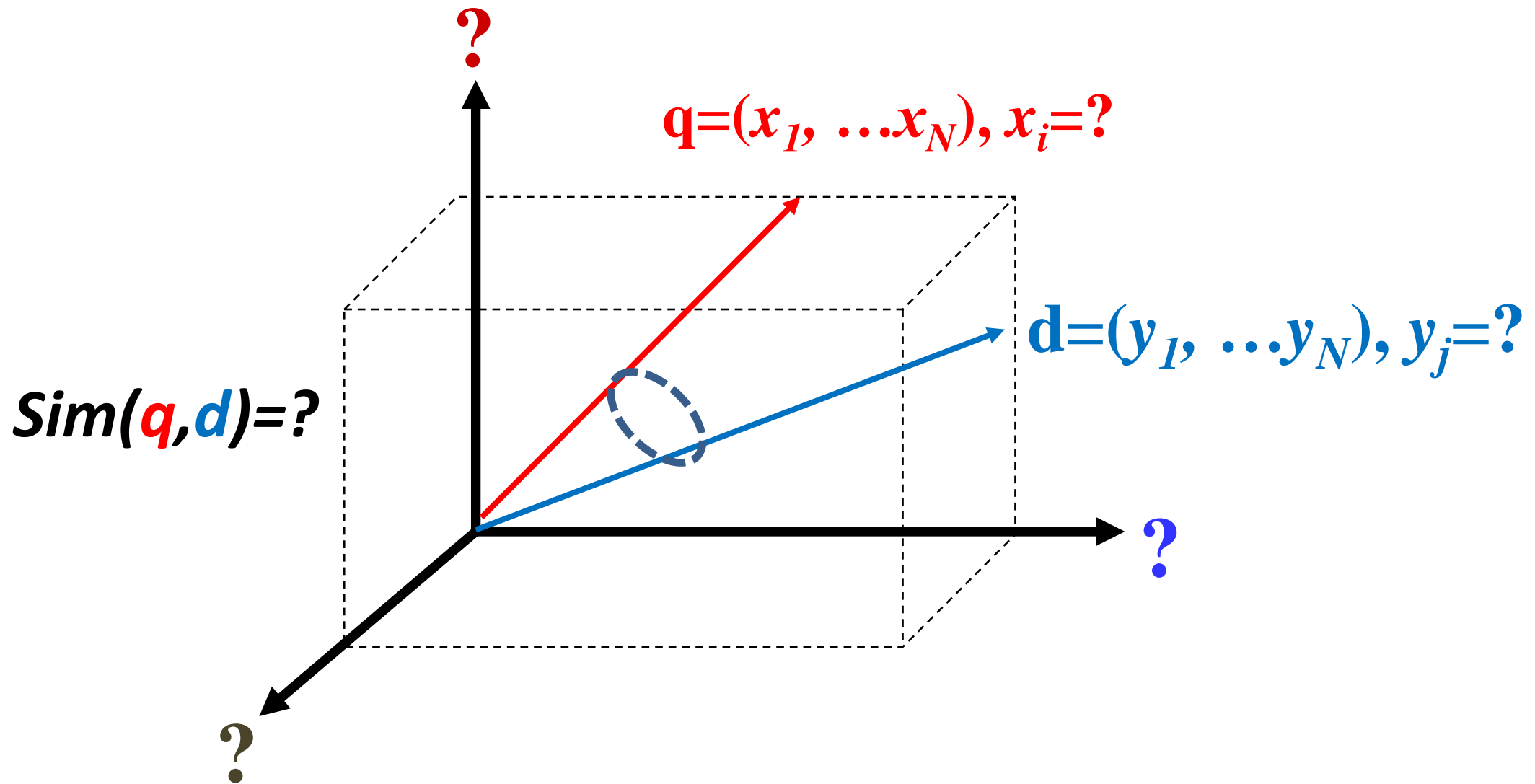
Vector Space Retrieval Model: Simplest Instantiation

ChengXiang “Cheng” Zhai
Department of Computer Science
University of Illinois at Urbana-Champaign

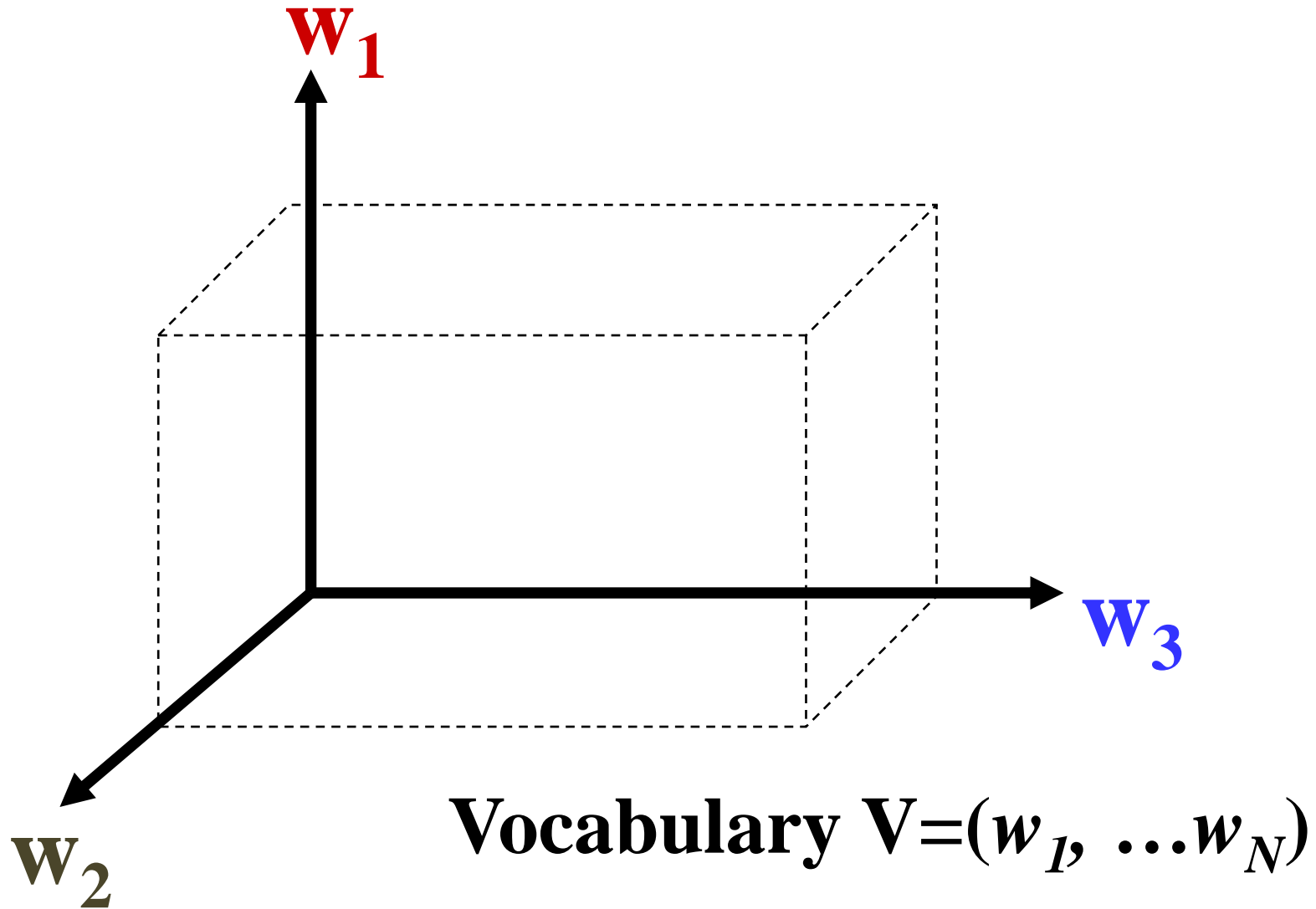
Course Schedule



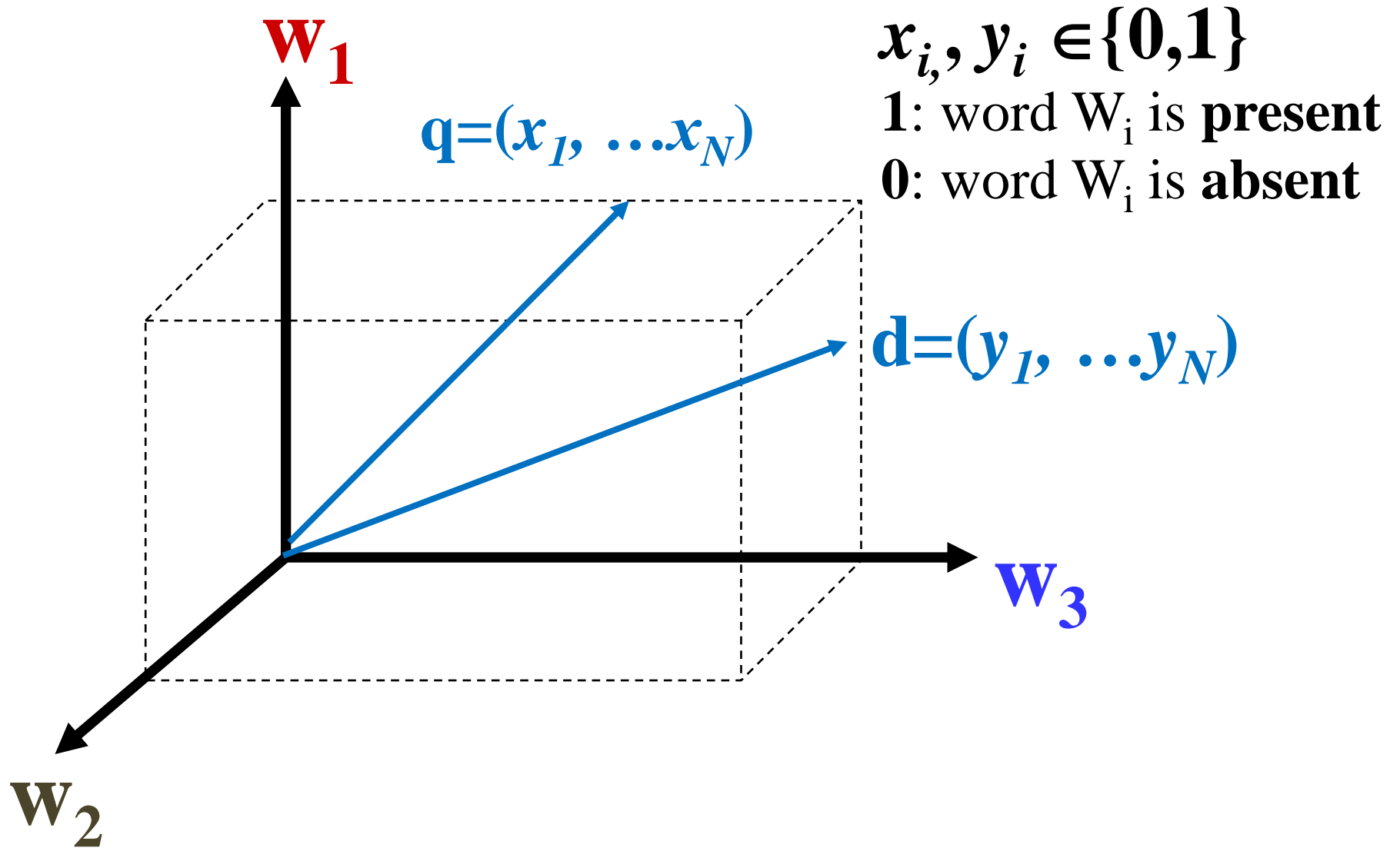
What VSM Doesn't Say



Dimension Instantiation: Bag of Words (BOW)

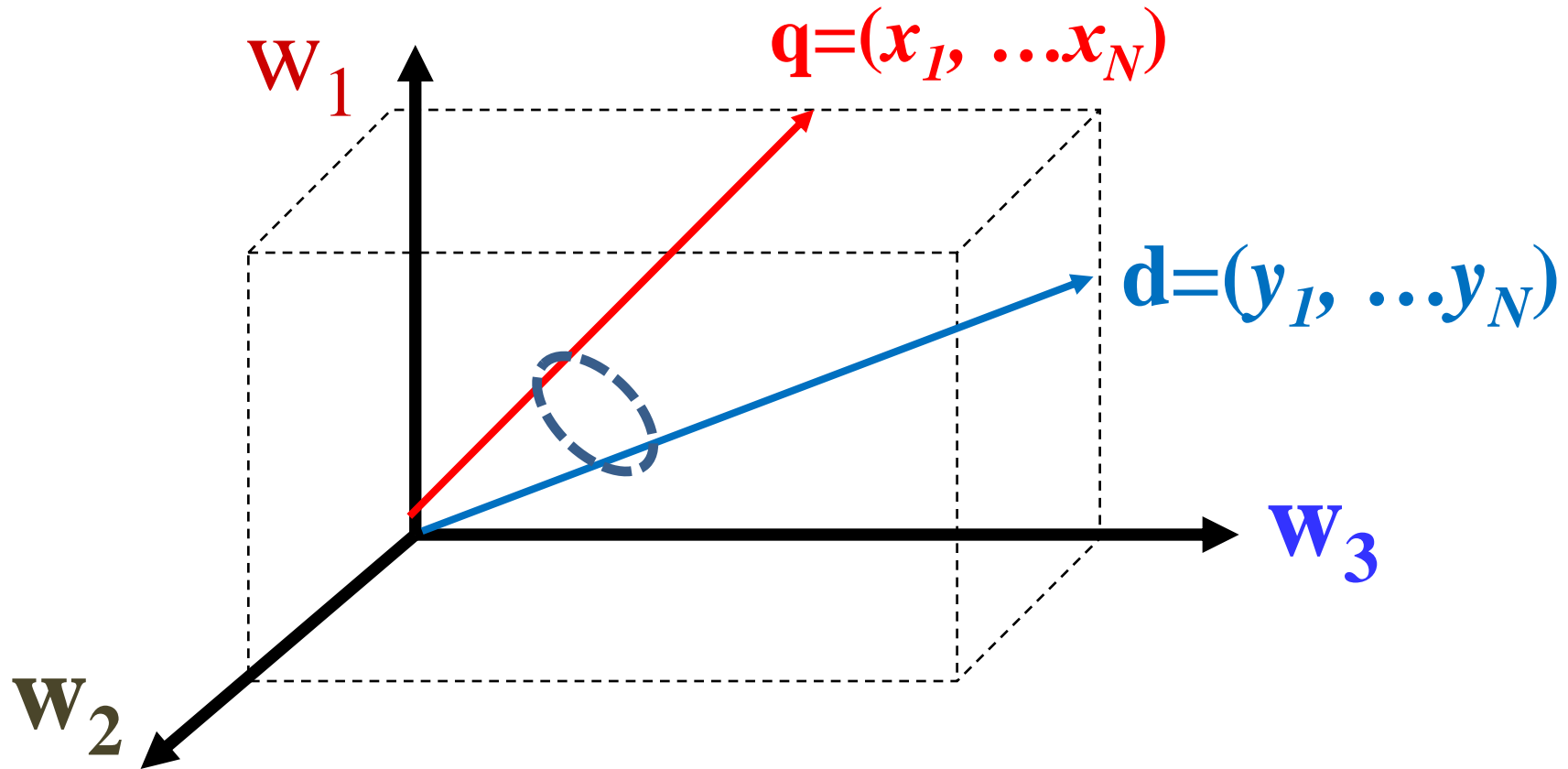


Vector Placement: Bit Vector



Similarity Instantiation: Dot Product

$$\text{Sim}(\mathbf{q}, \mathbf{d}) = \mathbf{q} \cdot \mathbf{d} = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$



Simplest VSM= Bit-Vector + Dot-Product + BOW

$$\begin{aligned} \mathbf{q} &= (x_1, \dots, x_N) & x_i, y_i &\in \{0, 1\} \\ \mathbf{d} &= (y_1, \dots, y_N) & 1: \text{word } W_i \text{ is present} \\ & & 0: \text{word } W_i \text{ is absent} \end{aligned}$$

$$\text{Sim}(\mathbf{q}, \mathbf{d}) = \mathbf{q} \cdot \mathbf{d} = x_1 y_1 + \dots + x_N y_N = \sum_{i=1}^N x_i y_i$$

What does this ranking function intuitively capture?
Is this a good ranking function?

An Example: How Would You Rank These Documents?

Query = “**news about presidential campaign**”

Ideal Ranking?

d1	... news about ...	d4 +
		d3 +
d2	... news about organic food campaign ...	
d3	... news of presidential campaign ...	
d4	... news of presidential campaign presidential candidate ...	d1 -
		d2 -
d5	... news of organic food campaign ... campaign ... campaign ... campaign ...	d5 -

Ranking Using the Simplest VSM

Query = “**news about presidential campaign**”

d1 ... **news about** ...

d3 ... **news** of **presidential campaign** ...

$V = \{\text{news, about, presidential, campaign, food ...}\}$

$q = (1, 1, 1, 1, 0, \dots)$

$d1 = (1, 1, 0, 0, 0, \dots)$

$f(q, d1) = 1*1 + 1*1 + 1*0 + 1*0 + 0*0 + \dots = 2$

$d3 = (1, 0, 1, 1, 0, \dots)$

$f(q, d3) = 1*1 + 1*0 + 1*1 + 1*1 + 0*0 + \dots = 3$

Is the Simplest VSM Effective?

Query = “news about presidential campaign”

d1	... news about ...	$f(q, d1)=2$
d2	... news about organic food campaign ...	$f(q, d2)=3$
d3	... news of presidential campaign ...	$f(q, d3)=3$
d4	... news of presidential campaign presidential candidate ...	$f(q, d4)=3$
d5	... news of organic food campaign ... campaign ... campaign ... campaign ...	$f(q, d5)=2$

Summary

- VSM instantiation: dimension, vector placement, similarity
- Simplest VSM
 - Dimension = word
 - Vector = 0-1 bit vector (word presence/absence)
 - Similarity = dot product
 - $f(q,d)$ = number of **distinct** query words matched in d