# Lecture 6: (Probabilistic) Latent Semantic Analysis

## Julia Hockenmaier

*juliahmr@illinois.edu*

3324 Siebel Center
Office hours: by appointment

# Indexing by Latent Semantic Analysis
## (Deerwester et al., 1990)
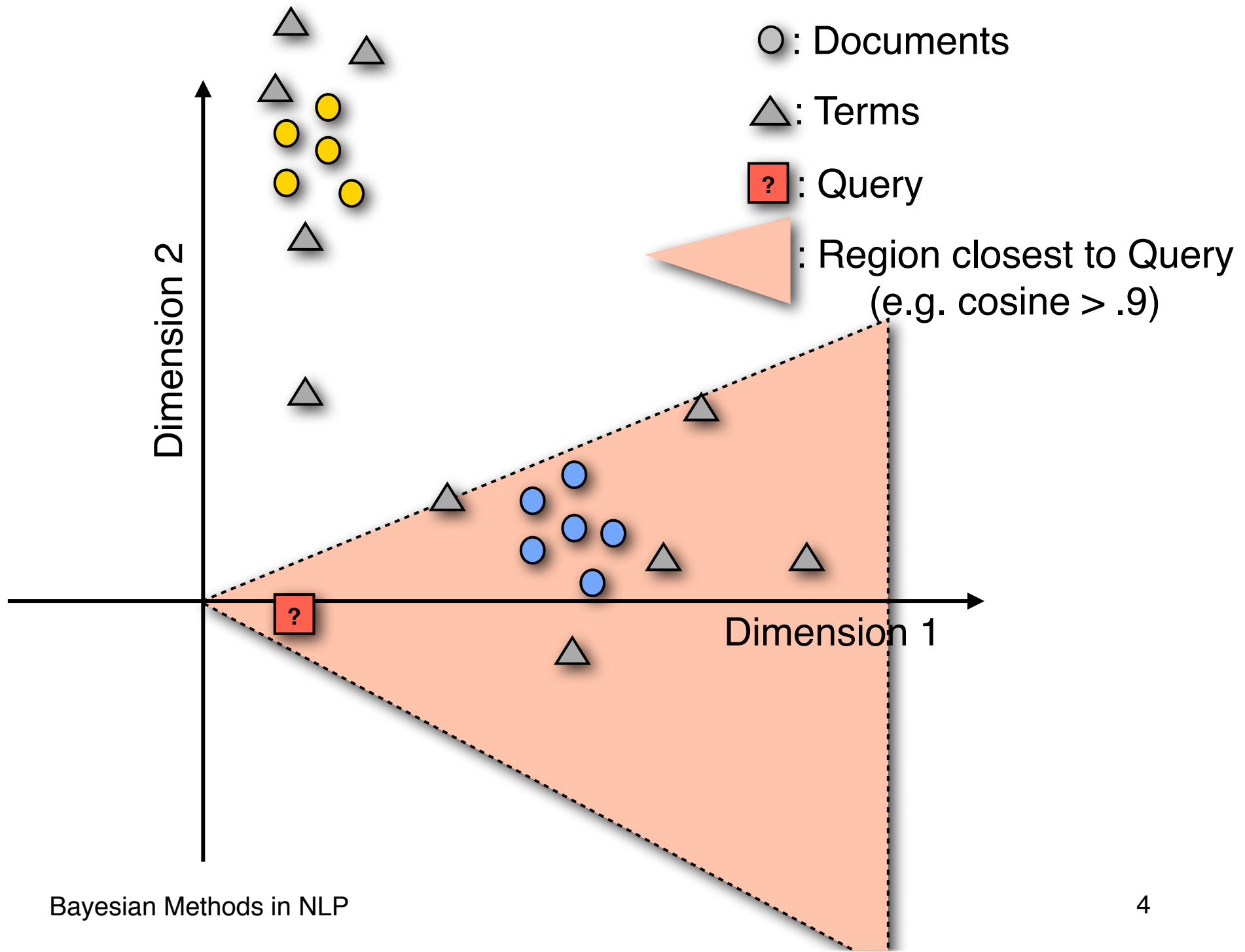
# Latent Semantic Analysis

The task:

 Return *relevant* documents for text queries

The problem: relevance is conceptual/semantic
- The index of relevant documents may not contain all query terms (**synonymy** and missing information)
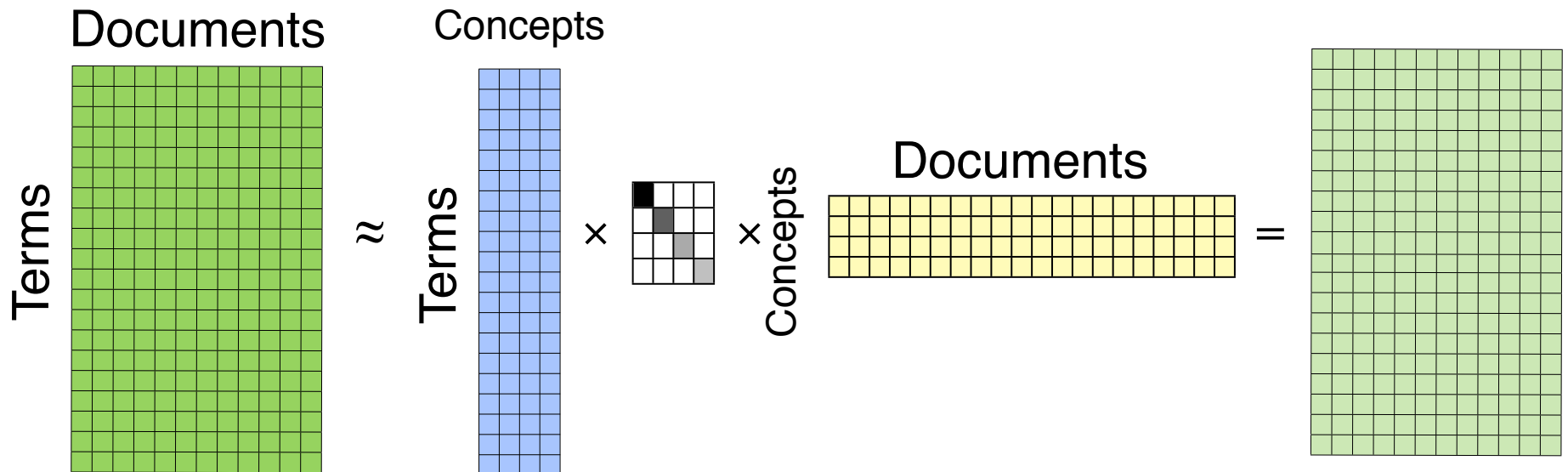- The query terms may be ambiguous (**polysemy**)

Indexing by Latent Semantic Analysis
- Map queries and documents into a new vector space whose $k$ dimensions correspond to independent concepts
- In this space, queries will be near semantically close documents

○ : Documents

△ : Terms

? : Query

: Region closest to Query (e.g. cosine > .9)

Dimension 2

Dimension 1

# Latent Semantic Analysis

Low-rank approximation of Singular Value Decomposition (SVD):



$$X \quad \approx \quad T_0 \quad \times \quad S_0 \quad \times \quad D_0\text{'} \quad = \dot{X}$$

$X$: Term-document matrix (=data): $X_{ij}$ = freq of $w_i$ in $D_j$

$\dot{X} = T_0 S_0 D_0\text{'}$ (k-rank approximation of $X$)

$T_0$: Columns are orthogonal and unit-length $T_0\text{'}T_0 = I$

$S_0$: Diagonal matrix of the $k$ largest singular values

$D_0$: Columns are orthogonal and unit-length $D_0\text{'}D_0 = I$

this should really be $\hat{X}$

# LSA: term similarity

$\mathbf{T_0}$    $\mathbf{\dot{X}}$    $\mathbf{\dot{X}`}$    $= \mathbf{T_0 \ S_0 \ S_0 \ T_0}$

Term $w_i$

dot product of $w_i$, $w_j$
in the new space
$\mathbf{T_0 \ S_0}$

$\mathbf{\dot{X}\dot{X}`} \ = \ \mathbf{T_0 \ S_0 \ S_0 \ T_0}$
(**D** cancels out because **S** is diagonal and **D** orthonormal)

Similarity of terms $w_i$, $w_j$ in the new space: $\mathbf{(\dot{X}\dot{X}`)}_{ij}$

# LSA: document similarity

$$\dot{X}' \quad \dot{X} \quad = \mathbf{D_0}\, \mathbf{S_0}\, \mathbf{S_0}\, \mathbf{D_0}$$

$\mathbf{D_0}$

Doc. $D_j$

$\dot{X}\dot{X}'$

dot product of $D_i$, $D_j$
in the new space
$\mathbf{D_0}\, \mathbf{S_0}$

$\dot{X}'\dot{X} \;=\; \mathbf{D_0}\, \mathbf{S_0}\, \mathbf{S_0}\, \mathbf{D_0}$
(**T** cancels out because **S** is diagonal and **T** orthonormal)

Similarity of documents $d_i$, $d_j$ in the new space: $(\dot{X}'\dot{X})_{ij}$

# LSA: term-document similarity

The elements of $\dot{X}$ give the similarity of terms and documents.

Now, terms are projected to $\mathbf{T}\mathbf{S}^{1/2}$, documents to $\mathbf{D}\mathbf{S}^{1/2}$

# LSA: query-document similarity

Queries $q$ are 'pseudo-documents':
they don't appear in $\mathbf{X}$

Construct their term vector $\mathbf{X}_q$
Define their document vector $\mathbf{D}_q = \mathbf{X'}_q \mathbf{T} \mathbf{S}^{-1}$

# Probabilistic Latent Semantic Indexing
## (Hofmann 1999)

# The aspect model

Observations are document-word pairs $(d, w)$

Assume there are *k* aspects $z_1...z_k$
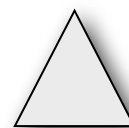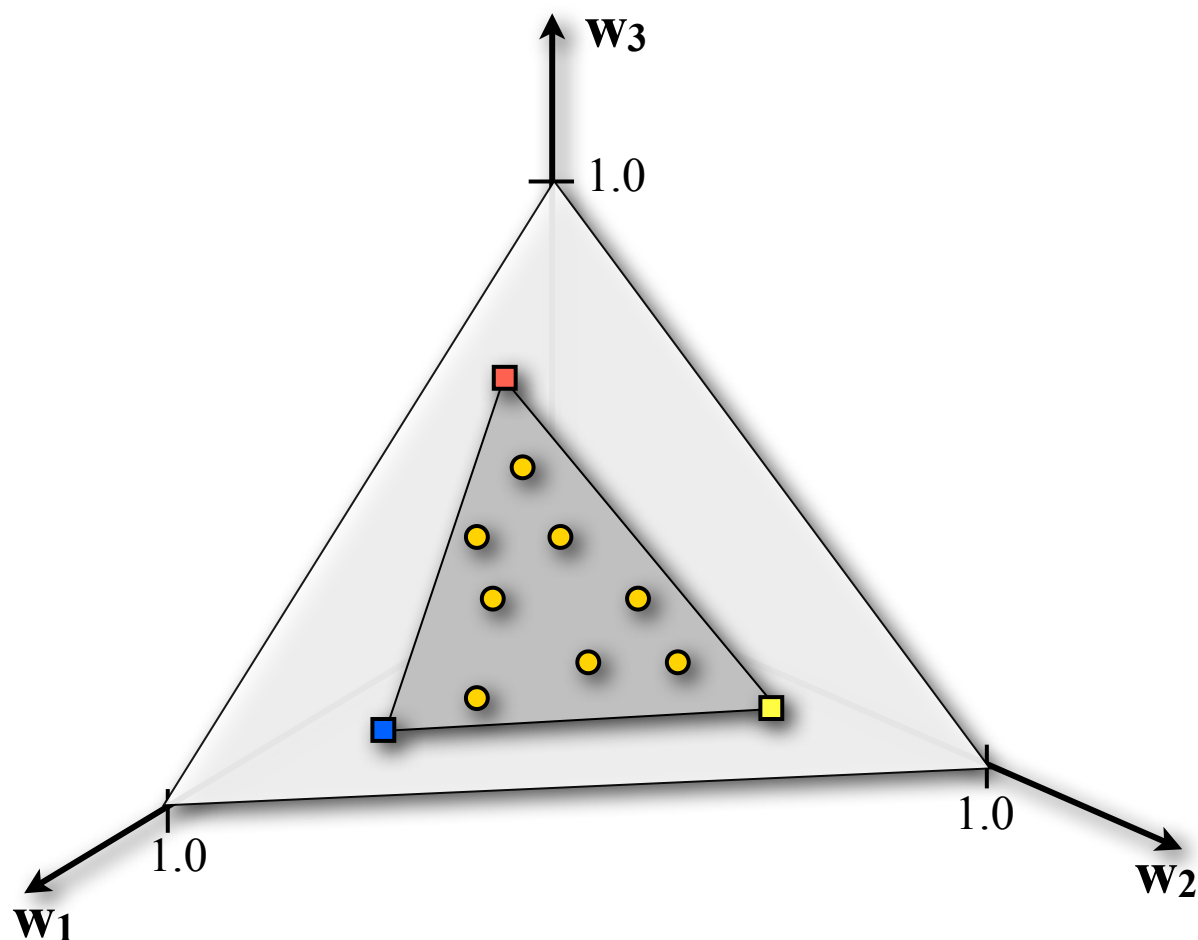Each observation is associated with a hidden aspect $z$

$$P(d, w) = P(d)P(w \mid d)$$
$$\text{with} \quad P(w \mid d) = \sum_{z \in Z} P(w \mid z)P(z \mid d)$$

Or, equivalently:
$$P(d, w) = \sum_{z \in Z} P(z)P(d \mid z)P(w \mid z)$$
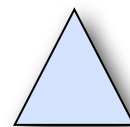
# A geometric interpretation



Word simplex
Any point in this simplex defines a multinomial over words

○ Documents $P(w \mid d)$
Each document corresponds to one multinomial over words

■ Topics $P(w \mid z)$
Each topic is a multinomial over words

Topic simplex
The topics define the corners of a (sub)simplex.
All training documents lie inside this topic simplex.

$$P(w \mid d) = \lambda_1 P(w \mid z_1) + \lambda_2 P(w \mid z_2) + \lambda_3 P(w \mid z_3)$$
$$= P(z_1 \mid d)P(w \mid z_1) + P(z_2 \mid d)P(w \mid z_2) + P(z_3 \mid d)P(w \mid z_3)$$

Bayesian Methods in NLP

# PLSA is a mixture model

## Mixture models:
- K mixture components and N observations $x_{1...}\ x_N$
- Mixing weights $(\theta_1 ... \theta_K)$: $P(\ k\ ) = \theta_K$
- Each observation $x_n$ is generated by mixture component $z_n$
  $$P(\ x_n\ ) = P(\ z_n\ )\ P(\ x_n\ |\ z_n\ )$$

## PLSI:
- Mixture components = topics
- Mixing weights are specific to each document $\theta_d = (\theta_{d1}...\theta_{dK})$
- Each observation (word) $w_{d,n}$ is a sample
  from the document-specific mixture model.
  It is drawn from one of the components $z_{d,n}$
  $$P(\ w_{d,n}\ ) = P(\ z_{d,n}\ |\ \theta_d\ )\ P(\ w_{d,n}\ |\ z_{d,n}\ )$$

# Estimation: EM algorithm

**E-step:** Recompute

$P(z \mid d, w) = P(z, d, w) / \sum_{z'} P(z', d, w)$

with $\quad P(z, d, w) = P(z)P(d \mid z)P(w \mid z)$

**M-step:** Recompute

$P(w \mid z) \propto \sum_d \text{freq}(d, w) \, P(z \mid d, w)$

$P(d \mid z) \propto \sum_w \text{freq}(d, w) \, P(z \mid d, w)$

$P(z) \quad \propto \sum_d \sum_w \text{freq}(d, w) \, P(z \mid d, w)$