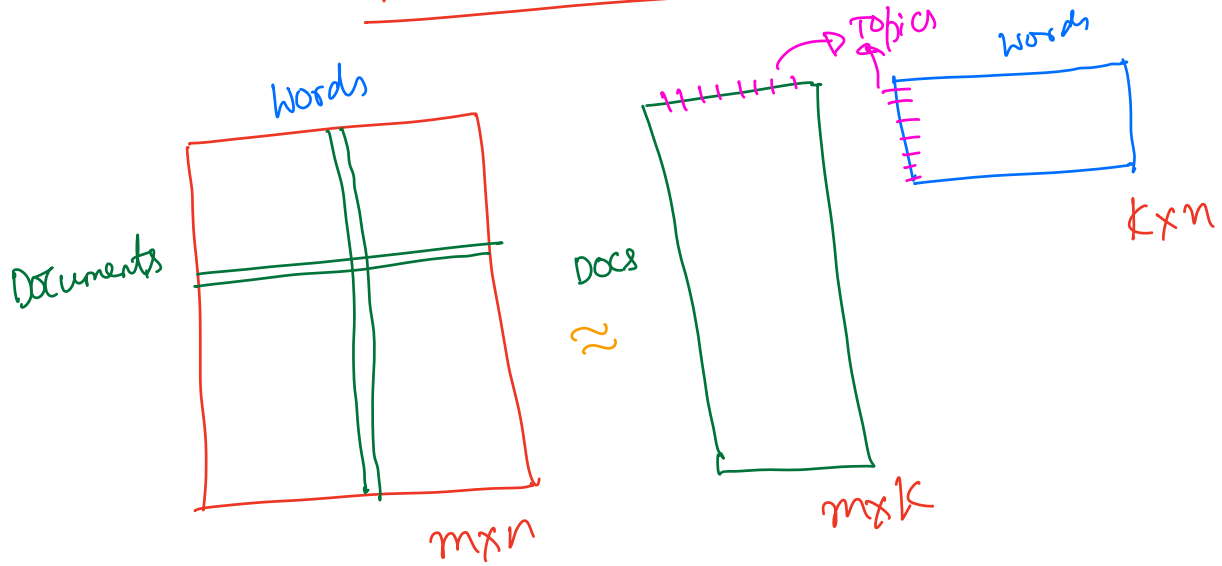
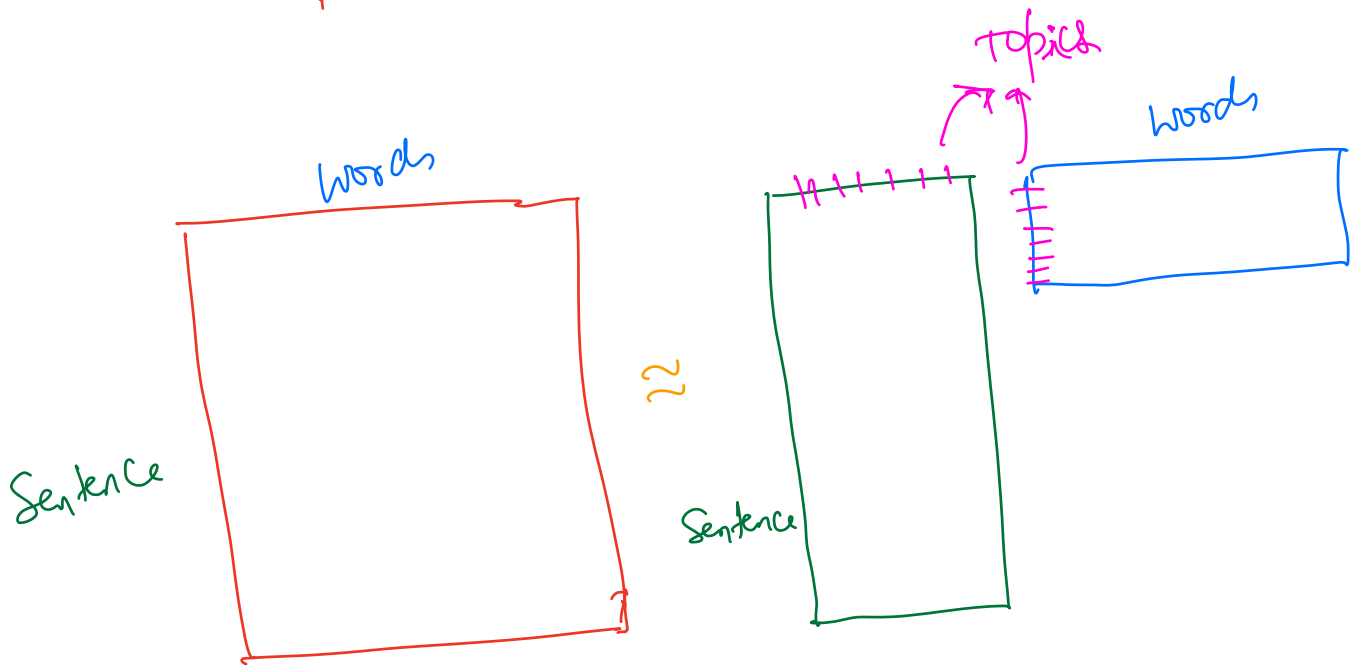
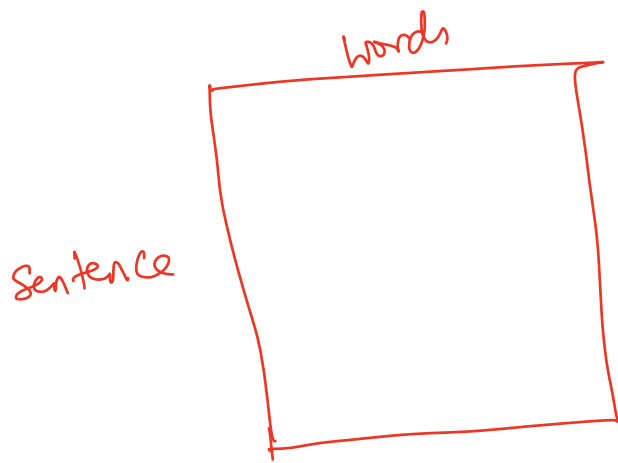


LECTURE 12
MATRIX FACTORIZATION



$k = \# \text{ Topics}$



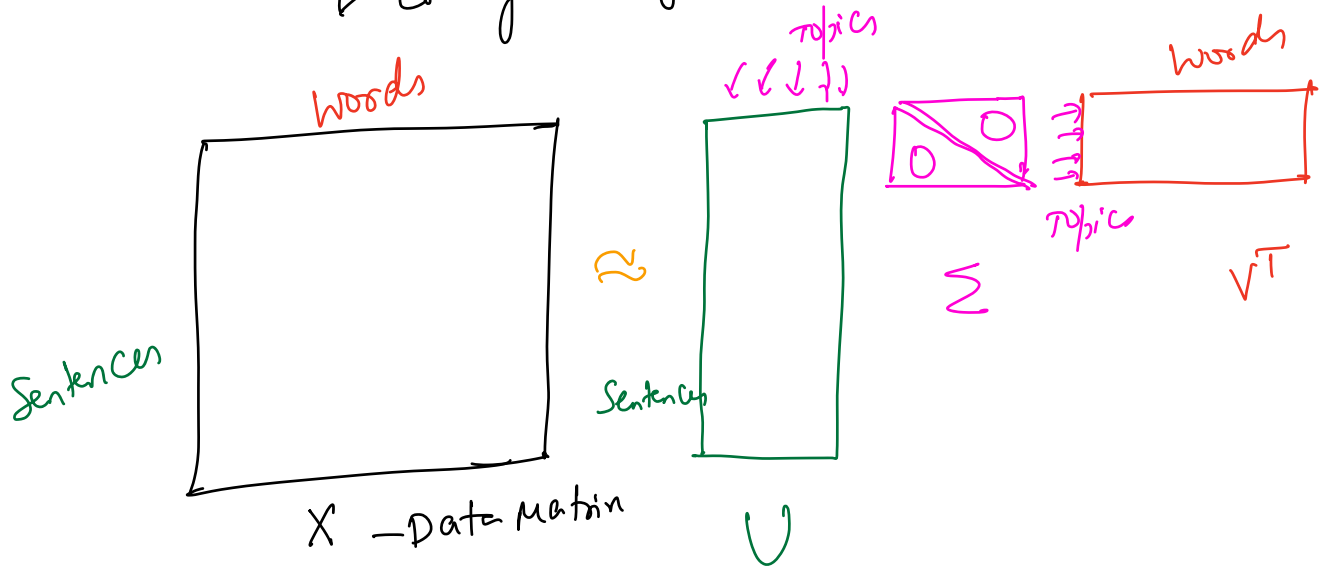
$$k = \# \text{topics}$$

- Give us the weightage of "Latent topics" in a sentence
- For a medical document, need to segment each sentence into a topic!
- How do we assign "one topic" to a sentence?

SVD

→ Singular value decomposition

→ Every single matrix has a SVD!



Σ - Diagonal matrix

ICE #0

Let X be a data matrix and

$X \approx U \Sigma V^T$ be the k th rank SVD of X .

i.e. Σ has k diagonal entries.

If X is $m \times n$ matrix, what is the dimension of

U ?

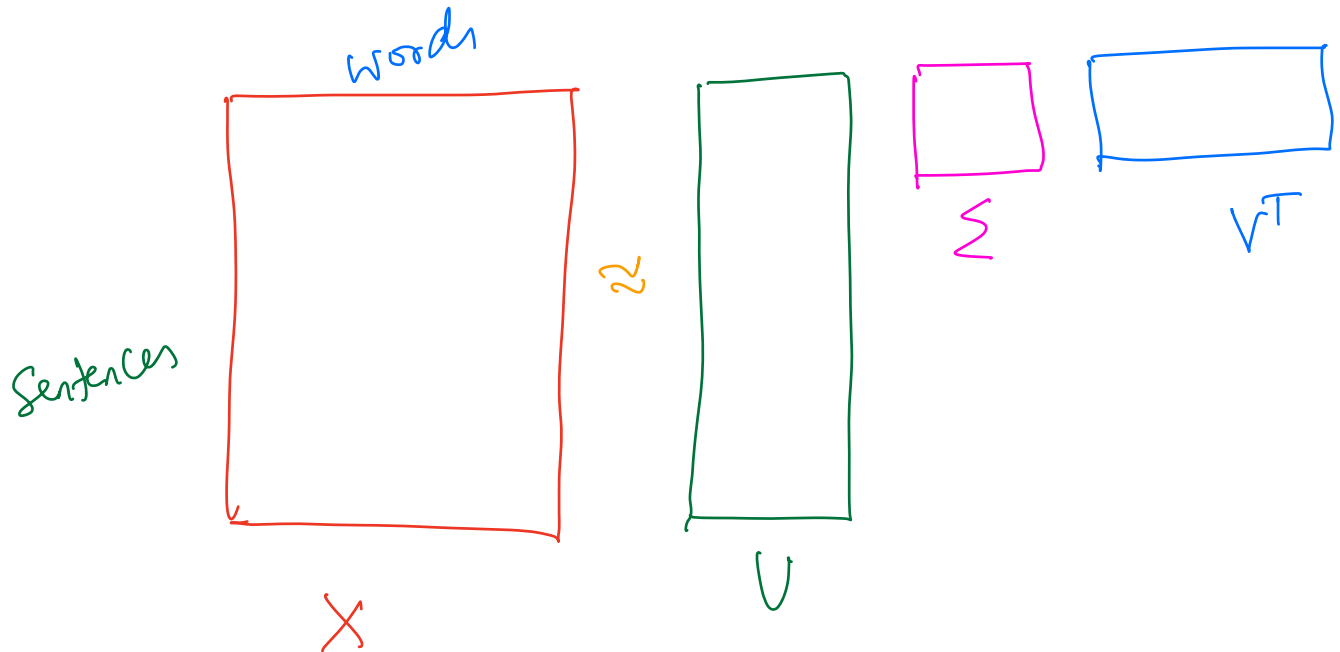
a) $k \times m$

b) $k \times k$

c) $m \times k$

d) $m \times n$

SVD yields a trivial matrix factorization!



$$X \approx U \Sigma V^T = \underbrace{(U \Sigma)}_{\text{Factor 1}} \underbrace{V^T}_{\text{Factor 2}}$$

How do we choose k , the number of ^{latent} topics?

→ SVD based matrix factorization doesn't guarantee that each sentence gets assigned a topic!

Matrix factorization

- An optimization perspective!

X - Data matrix - $m \times n$

U - Left factor

V^T - Right factor

Objective:-

$$\min_{U, V} \|X - UV^T\|_F^2$$

$m \times n$ $m \times k$ $k \times n$

without any constraints, U & V can be obtained from SVD as discussed!

ICE #1

Consider the obj. fn. $\min_{U, V} \|X - UV^T\|_F^2$.

Let's say $k = m$ & X is of dim. $m \times n$.

What is the minimum value the loss fn. can take in this case?

(a) -0.5

(b) 0.5

(c) 0

(d) Some value > 0

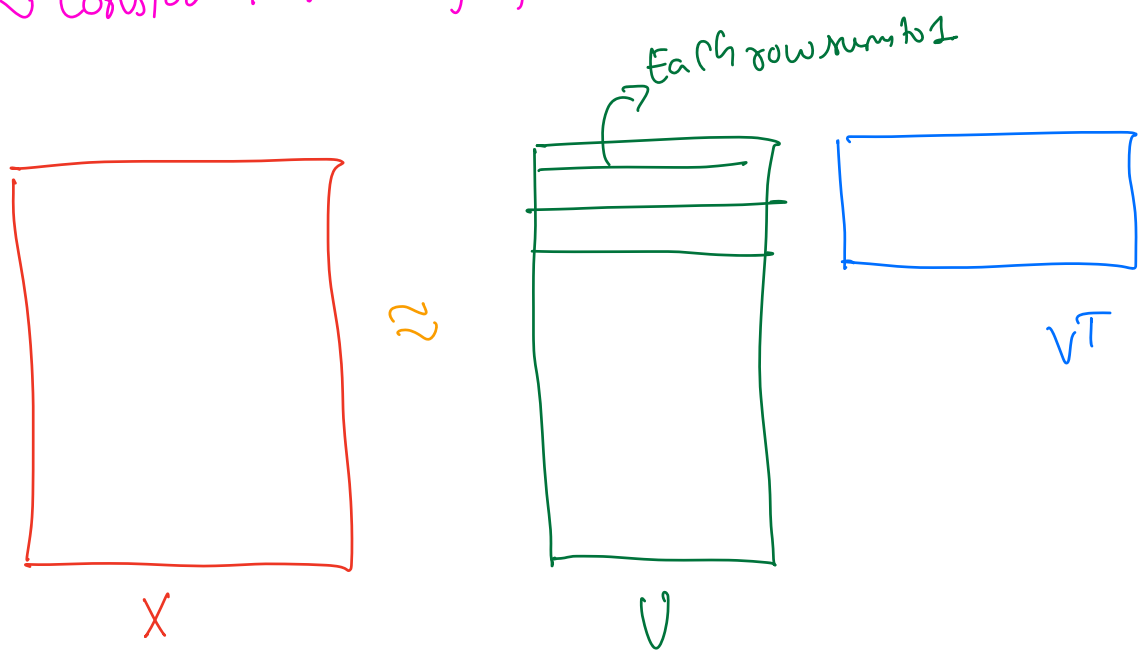
Simply fix $U=I$ & $v = x^T$.
What's the 1011 fn. value now?

Segmenting each sentence into a "Latent Topic"
↳ Unsupervised Topic Segmentation

~ SVD or simple matrix factorization doesn't assign a single topic to a sentence

~ work around!

~ Constrain the left factor to be a probability matrix



$U \mathbf{1} = \mathbf{1}$ (Each row of U is a probability vector!)

$$\min \|X - UV^T\|_F^2$$

$$U^T = 1,$$

$$U \geq 0, V$$

This turns out to be exactly what we think it is!!

U - Assignment Matrix

V - Centroid Matrix

Algorithm:- (k-means)

1. Initialize V to be k random sentences from X

2. Iterate:-

a) For fixed V , optimize for U

- Assign sentences to a representative Centroid

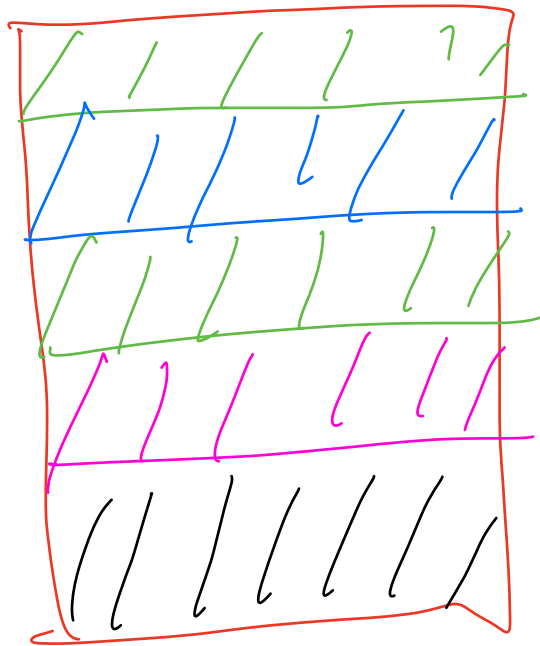
b) For fixed U , optimize for V

- Average sentence re cls in a cluster

k-means = Algorithm that optimizes a
Cost Matrix Factorization Loss
with constraints!

Probabilistic/soft k-means
→ Instead of hard assignment
of a sentence to a topic,
soft assign!

Topic Segmentation



Two problems in Health Care

1) Digital scribe notes to EHR

2) Given EHR's \rightarrow Info. retrieval!



E.S.
1) Get all patient records with a particular strain of disease?

2) Get % of patients that got treatment y

Challenges to Info Retrieval

1) Subject Identification

E.g. "Patient's father had diabetes"

b) Temporality

E.g. "Patient had diabetes until last year."

NLP Can help. 1) Graph of entities & relationships

2) query graph

Explainability of AI in health care.

Why?

→ Informed Consent

→ Liability Issues!

→ Decision based on meta data

—E.S. Huskies vs Wolves!

→ Avoid medical malpractice Law suits (Legal Issue)

Note:- Next Lecture!