

TENSOR FACTORIZATIONS METHODS : Lecture 3.

I) Tensor Methods for Topic models, Multinomial models, GMM. (applications).

II) Matrix Power Method.

III) Tensor Power Method.

I.a) Single topic models for documents.

Following is a simplified model for documents.

A document is a collection of l words x_1, x_2, \dots, x_l

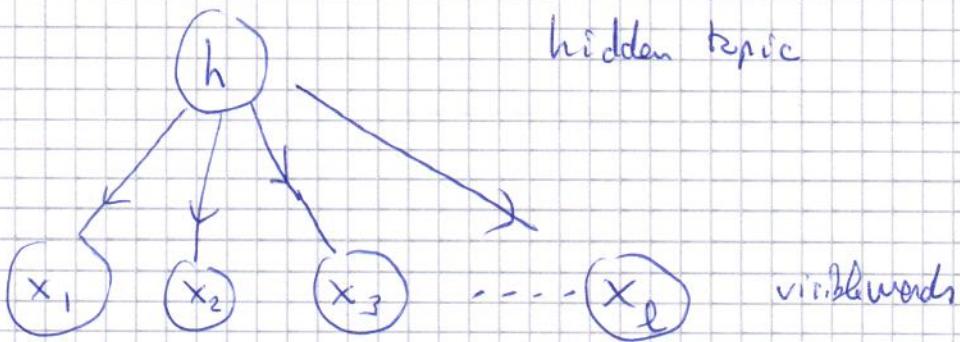
x_1	x_2	x_3
x_4	x_5	x_6
x_7	x_8	x_9
x_{10}	x_{11}	x_{12}
x_{13}	x_{14}	x_{15}

where each word is a random variable taking values in some dictionary of size D . So the fifth word of the

document, say x_5 , can take D possible values.

We assume that each document is about a single topic. There are K possible topics indexed by $h \in \{1, \dots, K\}$ (e.g. sport, music, politics, cinema).

The probabilistic model of documents is a belief network of the form :



where all conditional probabilities are identical $p(x_i|h)$.

So the prob dist for a document is

$$P(x_1, x_2, \dots, x_l, h) = p(h) \prod_{i=1}^l p(x_i|h)$$

Given a Topic words are i.i.d.

We use the following formalism:

- $p(h=i) = w_i$ = weight of topic $i \in \{1 \dots K\}$.
- $\underline{x} \in$ Dictionary of D words = $\{\underline{e}_1, \underline{e}_2, \dots, \underline{e}_D\}$ where each $\underline{e}_i = \begin{pmatrix} 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix}$ i th coordinate. So we

identify words with canonical vectors of basis of \mathbb{R}^D .

We set $p(\underline{x} | i) = \underline{\mu}_i$ = vector of probabilities of words given a topic i .

Remark: words are exchangeable random variables

in the sense that $p(x_1, x_2, \dots, x_l) = \sum_{h=1}^K p(h) \prod_{i=1}^l p(x_i|h)$
is permutation invariant.

Our problem: Given frequencies of occurrence of words in documents

$$\Pr(x_1) ; \Pr(x_1, x_2) ; \Pr(x_1, x_2, x_3) \text{ etc}$$

we want to determine the parameters

$$\left\{ \begin{array}{l} (w_i)_{i=1 \dots k} \text{ weights of Topics} \\ \mu_i \text{ weights of words given a Topic } i \in \{1 \dots k\} \end{array} \right.$$

This is done by matching empirical moments to theoretical moments.

Theoretical Moments of our probabilistic model:

$$\begin{aligned} \mathbb{E}(x_1) &= \sum_{\alpha=1}^D \underbrace{\Pr(x_1 = e_\alpha)}_{\substack{\text{prob of word} \\ e_\alpha}} \cdot \underbrace{e_\alpha}_{\substack{\uparrow \\ \text{possible value of word in Dict.}}} \end{aligned}$$

$$\begin{aligned} &= \sum_{\alpha=1}^D \sum_{i=1}^K \Pr(x_1 = e_\alpha | i) w_i \cdot e_\alpha \\ &= \sum_{i=1}^K \left\{ \sum_{\alpha=1}^D \underbrace{\Pr(x_1 = e_\alpha | i)}_{\mu_i^\alpha} \right\} w_i \end{aligned}$$

$$\Rightarrow \boxed{\mathbb{E}(x_1) = \sum_{i=1}^K \mu_i^\alpha w_i}$$

Note that this formula is the same for the first, second, third etc word of the document.

Similarly for second and third order moments :

$$\text{IE}(\underline{x}_1 \otimes \underline{x}_2) = \sum_{\alpha=1}^D \sum_{\beta=1}^D \underbrace{\Pr(x_1 = e_\alpha, x_2 = e_\beta)}_{M_{\alpha\beta}^2} e_\alpha \otimes e_\beta$$

$$= \sum_{\alpha, \beta} w_i \underbrace{\Pr(x_1 = e_\alpha, x_2 = e_\beta | i)}_{\Pr(x_1 = e_\alpha | i) \Pr(x_2 = e_\beta | i)} e_\alpha \otimes e_\beta$$

$$= \sum_{i=1}^K w_i \left\{ \sum_{\alpha} p(x_1 = e_\alpha | i) \right\} \left\{ \sum_{\beta} p(x_2 = e_\beta | i) \right\}$$

$$= \sum_{i=1}^K w_i \underline{\mu}_i \otimes \underline{\mu}_i$$

$$\Rightarrow \boxed{\text{IE}(\underline{x}_1 \otimes \underline{x}_2) = \sum_{i=1}^K w_i \underline{\mu}_i \otimes \underline{\mu}_i}$$

Also,

$$\text{IE}(\underline{x}_1 \otimes \underline{x}_2 \otimes \underline{x}_3) = \sum_{i=1}^K w_i \underline{\mu}_i \otimes \underline{\mu}_i \otimes \underline{\mu}_i$$

$$= \sum_{\alpha \beta \gamma} \Pr(x_1 = e_\alpha, x_2 = e_\beta, x_3 = e_\gamma) e_\alpha \otimes e_\beta \otimes e_\gamma$$

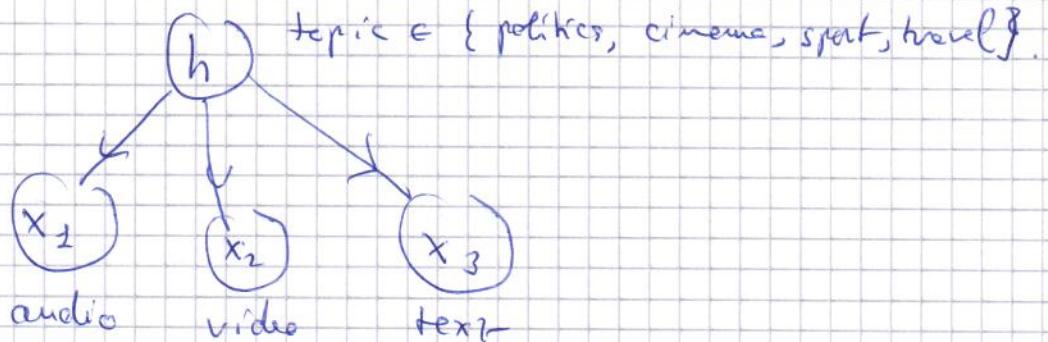
Same formulae are valid for any pair and/or triple of words.

The analysis of documents can be done by making decompositions of empirical Tensors

$$\Pr_{\text{emp}}(x_1 = e_\alpha, x_2 = e_\beta) = M_{\alpha\beta}^2 ; \quad \Pr_{\text{emp}}(x_1 = e_\alpha, x_2 = e_\beta, x_3 = e_\gamma) = \frac{M_{\alpha\beta\gamma}^3}{M_3^3}$$

I, b) Multiview Models,

Very similar but slightly more general than previous model. Suppose we have a hidden Topic $h \in \{1 \dots, K\}$ for a set of K possible Topics (politics, cinema, sport, travel) and suppose we have many views (say: audio, video, text) of a topic :



$$P_1(x_1|h) \quad P_2(x_2|h) \quad P_3(x_3|h)$$

Each $x_i \in \{e_1, e_2, \dots, e_{D_i}\}$ D_i possible values encoded in canonical basis vectors of \mathbb{R}^{D_i} . $i=1, 2, 3$

Joint p.d.:

$$P(x_1, x_2, x_3|h) = P(h) \underbrace{P_1(x_1|h) P_2(x_2|h) P_3(x_3|h)}_{}$$

Now P_1, P_2, P_3 can be unequal (more general than before) so we do not have exchangeability.

- Again it is easy to see :

$$\mathbb{E}(\underline{x}_i) = \sum_{\alpha=1}^{D_i} \Pr(\underline{x}_i = e_\alpha) e_\alpha .$$

$$= \text{vector} \left(\Pr(\underline{x}_i = e_\alpha) \right)_{\alpha=1}^{D_i} .$$

$$\mathbb{E}(\underline{x}_i \otimes \underline{x}_j) = \sum_{\alpha=1}^{D_i} \sum_{\beta=1}^{D_j} \Pr(\underline{x}_i = e_\alpha, \underline{x}_j = e_\beta) e_\alpha \otimes e_\beta$$

$$= \text{Tensor or Matrix } M^{\alpha\beta} = \Pr(\underline{x}_i = e_\alpha, \underline{x}_j = e_\beta)$$

↑
for each pair $(i, j) \in \{1, 2, 3\}^2$.

$$\mathbb{E}(\underline{x}_1 \otimes \underline{x}_2 \otimes \underline{x}_3) = \sum_{\alpha \beta \gamma}^{D_1 D_2 D_3} \Pr(\underline{x}_1 = e_\alpha, \underline{x}_2 = e_\beta, \underline{x}_3 = e_\gamma) .$$

$e_\alpha \otimes e_\beta \otimes e_\gamma$

= Tensor or 3d array of numbers

$$\Pr(\underline{x}_1 = e_\alpha, \underline{x}_2 = e_\beta, \underline{x}_3 = e_\gamma) .$$

- We have the decomposition:

$$\begin{cases} \mathbb{E}(\underline{x}_1) = \sum_{i=1}^K w_i \underline{\mu}_i^{(1)} & \underline{\mu}_i^{(1)} = \Pr(\underline{x}_1 = i) \\ \mathbb{E}(\underline{x}_1 \otimes \underline{x}_2) = \sum_{i=1}^K w_i \underline{\mu}_i^{(1)} \otimes \underline{\mu}_i^{(2)} \\ \mathbb{E}(\underline{x}_1 \otimes \underline{x}_2 \otimes \underline{x}_3) = \sum_{i=1}^K w_i \underline{\mu}_i^{(1)} \otimes \underline{\mu}_i^{(2)} \otimes \underline{\mu}_i^{(3)} \end{cases}$$

I, C) Gaussian Mixture Models.

$$\text{Let } p(x) = \sum_{i=1}^K w_i \cdot \frac{e^{-\frac{\|x - \underline{\alpha}_i\|^2}{2\sigma^2}}}{(2\pi\sigma^2)^{D/2}}$$

= convex combination of spherical Gaussian clusters around $\underline{\alpha}_1, \dots, \underline{\alpha}_K$.

In the exercises you show that:

$$\mathbb{E}(x) = \sum_{i=1}^K w_i \underline{\alpha}_i \equiv \underline{m}$$

$$\mathbb{E}(x \otimes x) = \sigma^2 I_{D \times D} + \sum_{i=1}^K w_i \underline{\alpha}_i \otimes \underline{\alpha}_i$$

$$\mathbb{E}(x \otimes x \otimes x) = \sum_{i=1}^K w_i \underline{\alpha}_i \otimes \underline{\alpha}_i \otimes \underline{\alpha}_i$$

$$+ \sigma^2 \sum_{i=1}^K (\underline{m} \otimes \underline{\alpha}_i \otimes \underline{\alpha}_i + \underline{\alpha}_i \otimes \underline{m} \otimes \underline{\alpha}_i + \underline{\alpha}_i \otimes \underline{\alpha}_i \otimes \underline{m})$$

If we have a way to estimate \underline{m} & σ^2 then we can determine w_i & $\underline{\alpha}_i$'s by looking at decomposing of the Tensor $\mathbb{E}(x^\alpha x^\beta x^\gamma) = T^{\alpha\beta\gamma}$.

Remark: \underline{m} is estimated by empirical mean. To estimate σ^2 we can compute smallest e.v of $\mathbb{E}(x^\alpha x^\beta)$ for the case when $K < D$ so we know the matrix $\sum_{i=1}^K w_i \underline{\alpha}_i \underline{\alpha}_i^T$ is not full rank and has a zero e.v.

For the QMM and single topic document model we see that the tensor decomposition is special because it involves only terms $\underline{\mu}_i \otimes \underline{\mu}_i \otimes \underline{\mu}_i$.

The "power method" has been investigated in detail for this class of tensors, then we give a brief introduction to the basics of the power method.

II. Matrix power method -

Let M a symmetric real matrix $N \times N$: It has a set of real eigenvalues & orthonormal eigenvectors

$$M \underline{v}_i = \lambda_i \underline{v}_i \quad \lambda_i \in \mathbb{R}, \quad \underline{v}_i^T \underline{v}_j = \delta_{ij} = \begin{cases} 1 & i=j \\ 0 & i \neq j \end{cases}$$

The spectral decomposition of M is

$$M = \sum_{i=1}^N \lambda_i \underline{v}_i \underline{v}_i^T = \sum_{i=1}^N \lambda_i \underline{v}_i \otimes \underline{v}_i$$

and we assume $\lambda_1 \geq \lambda_2 \geq \lambda_3 \dots \geq \lambda_N$.

The power method is an iterative method that allows to find eigenvalues and eigenvectors assuming no degeneracy.



Power Method for top eigenvalue / eigenvector:

- Choose initial unit vector $\underline{x}^{(0)}$ at time $t=0$

s.t. $\underline{x}^{(0)}$ is NOT orthogonal to $\underline{\nu}_1$.

- Iterate
$$\underline{x}^{(t)} = \frac{M \underline{x}^{(t-1)}}{\|M \underline{x}^{(t-1)}\|_2} \quad t \geq 1.$$

Lemma,

Assume $\lambda_1 > \lambda_2$ (so in particular

$\underline{\nu}_1$ is unique). Then

$$\begin{cases} \underline{x}^{(t)} \rightarrow \underline{\nu}_1 \text{ in } L_2 \text{ norm as } t \rightarrow +\infty \\ \underline{x}^{(t)} M \underline{x}^{(t)} \rightarrow \lambda_1 \text{ as } t \rightarrow +\infty \end{cases}$$

and the convergence rate is $\approx O\left(\left(\frac{\lambda_2}{\lambda_1}\right)^t\right)$.

exp part.

Remark: In very high dimensions $D \gg 1$ it is not so easy to choose $\underline{x}^{(0)}$ that is "far from orthogonal to $\underline{\nu}_1$ " (even for a random choice!) so there is a big constant C_D in front of $O\left(\left(\frac{\lambda_2}{\lambda_1}\right)^t\right)$ which makes initial convergence rate slow.

Proof of Lemma -

First we remark that

$$\underline{x}^{t-1} = \frac{M \underline{x}^{t-2}}{\|M \underline{x}^{t-2}\|_2} \Rightarrow \underline{x}^t = \frac{M \underline{x}^{t-1}}{\|M \underline{x}^{t-1}\|_2}$$

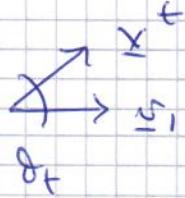
$$= \frac{M^2 \underline{x}^{t-2}}{\|M^2 \underline{x}^{t-2}\|_2} \cdot \frac{1}{\left\| \frac{M^2 \underline{x}^{t-2}}{\|M^2 \underline{x}^{t-2}\|_2} \right\|}$$

$$\Rightarrow \underline{x}^t = \frac{M^2 \underline{x}^{t-2}}{\|M^2 \underline{x}^{t-2}\|_2} \Rightarrow \boxed{\underline{x}^t = \frac{M^t \underline{x}^0}{\|M^t \underline{x}^0\|_2}}$$

From the spectral decomposition we get :

$$\underline{x}^t = \frac{\sum_{i=1}^N \lambda_i^t \underline{w}_i^\top (\underline{w}_i^\top \cdot \underline{x}^0)}{\left\{ \sum_{i=1}^N \lambda_i^{2t} (\underline{w}_i^\top \cdot \underline{x}^0)^2 \right\}^{1/2}}$$

We look at the angle



$$(\cos \theta_t)^2 = \underline{w}_1^\top \underline{x}^t = \frac{\lambda_1^{2t} (\underline{w}_1^\top \cdot \underline{x}^0)^2}{\sum_{i=1}^N \lambda_i^{2t} (\underline{w}_i^\top \cdot \underline{x}^0)^2}$$

$$\Rightarrow (\cos \theta_t)^2 = \frac{1}{1 + \sum_{i=2}^N \left(\frac{\lambda_i}{\lambda_1} \right)^{2t} \frac{(\underline{w}_i^\top \cdot \underline{x}^0)^2}{(\underline{w}_1^\top \cdot \underline{x}^0)^2}}$$

Since $\left(\frac{\lambda_i}{\lambda_1}\right)^{2t} \leq \left(\frac{\lambda_2}{\lambda_1}\right)^{2t}$ for $i \geq 2$ we have

$$\sum_{i=2}^N \left(\frac{\lambda_i}{\lambda_1}\right)^{2t} \left(\frac{\mathbf{n}_i^T \mathbf{x}^0}{\mathbf{n}_1^T \mathbf{x}^0}\right)^2 \leq \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} \frac{\sum_{i=2}^N (\mathbf{n}_i^T \mathbf{x}^0)^2}{(\mathbf{n}_1^T \mathbf{x}^0)^2}$$

and since $\sum_{i=1}^N (\mathbf{n}_i^T \mathbf{x}^0)^2 = 1 \Rightarrow \sum_{i=2}^N (\mathbf{n}_i^T \mathbf{x}^0)^2 = 1 - \mathbf{n}_1^T \mathbf{x}^0$

we finally get

$$(\cos \delta_t)^2 \geq \frac{1}{1 + \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} \frac{1 - (\cos \delta_0)^2}{(\cos \delta_0)^2}}$$

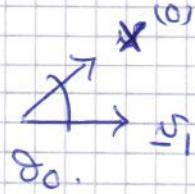
$$\Rightarrow 1 - (\cos \delta_t)^2 \leq \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} (\log \delta_0)^2 \cdot \frac{1}{1 + \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} (\log \delta_0)^2} \\ \leq \left(\frac{\lambda_2}{\lambda_1}\right)^{2t} (\log \delta_0)^2$$

end of proof.

If δ_0 is away from $\frac{\pi}{2}$ this tends to

zero expect as $\left(\frac{\lambda_2}{\lambda_1}\right)^{2t}$. Note that

$$\|\mathbf{x}^t - \mathbf{n}_1\|^2 = 1 + 1 - 2 \underline{x}^t \cdot \underline{\mathbf{n}}_1 = 2(1 - \cos \delta_t)$$



For $\mathbf{x}^t \perp \mathbf{x}^t$ we show that it tends to $\sum_i \lambda_i \underline{\mathbf{n}}_i = \lambda_1$.

(details left as exercise).



Power Method for other eigenvalues / eigenvectors
of the Matrix!

- "Deflate" the matrix $M' \leftarrow M - \lambda_1 \underline{v}_1 \underline{v}_1^T$.

Since $M' = \sum_{i=2}^N \lambda_i \underline{v}_i \underline{v}_i^T$ we can

apply the same power method to M' to get

$\lambda_2, \underline{v}_2$ as long as $\lambda_2 > \lambda_3$.

- and so on for other $\lambda_i, \underline{v}_i$'s as long as we have always a gap between eigenvalues.

Remark!

Stability of the algorithm with respect to noise will require good gaps between eigenvalues [see for example Golub & van Loan 1996, Sec 8.2-3].

III Tensor Power Method.

III, a) Method for "orthogonal tensors".

First we look at tensors that have a decomposition of the form

$$T = \sum_{i=1}^k \lambda_i \underline{v}_i \otimes \underline{v}_i \otimes \underline{v}_i$$

where $\lambda_i \in \mathbb{R}$ and $[\underline{v}_1 \dots \underline{v}_k]$ is an array of orthonormal vectors in \mathbb{R}^D . Note that necessarily we must have $k \leq D$ and that Tammrich's theorem applies. So this polyadic decomposition is unique.

We will generalize from matrix case. The main idea is as follows:

input T^{xy} ; output λ_i & \underline{v}_i .

1) Take $\underline{x}^{(0)}$ unit vector in \mathbb{R}^D at random.

2) Iterate $\underline{x}^t = \frac{T(I, \underline{x}^t, \underline{x}^t)}{\|\underline{T}(I, \underline{x}^t, \underline{x}^t)\|_2}$ where by

definition $\underline{T}(I, \underline{x}^t, \underline{x}^t) = \sum_{i=1}^k \lambda_i \underline{v}_i (\underline{v}_i^T \cdot \underline{x}^t)^2$

3) Under suitable conditions: $\underline{x}^t \rightarrow$ some \underline{v}^* and $\lambda^t \rightarrow$ some λ^* (see later).

4) Deflate the Tensor

$$\underline{T}' \leftarrow \underline{T} - \lambda^* \underline{v}^* \otimes \underline{v}^* \otimes \underline{v}^*$$

5) and so on.

Remarks:

- a) Convergence to \underline{v}^* , λ^* depends on the initial choice among other things and one does not converge to "top" "eigenvalue" "eigenvector".
- b) Steps 3 and 4 work out if a certain non-degeneracy condition is met (see them later).
- c) This method can be applied to tensors s.t. the \underline{v}_i 's are NOT orthogonal after a suitable orthonormalization procedure called "whitening" (see last paragraph).

Theorem on Tensor Power Method.

Let $T = \sum_{i=1}^k \lambda_i \underline{v}_i \otimes \underline{v}_i \otimes \underline{v}_i$ with $[\underline{v}_1, \dots, \underline{v}_k] \in \mathbb{R}^{D \times k}$ an orthonormal set of vectors and $\lambda_1, \dots, \lambda_k \in \mathbb{R}$.

Suppose we initialize the method with $\underline{x}^{(0)} \in \mathbb{R}^D$ s.t. the set of numbers

$$|\lambda_1 \underline{v}_1^T \cdot \underline{x}^{(0)}|, \dots, |\lambda_k \underline{v}_k^T \cdot \underline{x}^{(0)}|$$

has a unique largest element. Without loss of generality suppose

$$|\lambda_1 \underline{v}_1^T \cdot \underline{x}^{(0)}| > |\lambda_2 \underline{v}_2^T \cdot \underline{x}^{(0)}| \geq \dots \geq \dots$$

\uparrow strict for this one

$$\text{Set } \underline{x}^t = \frac{T(I, \underline{x}^{t-1}, \underline{x}^{t-1})}{\|T(I, \underline{x}^{t-1}, \underline{x}^{t-1})\|} \quad t \geq 1$$

$$\text{Then } \|\underline{x}^t - \underline{v}_1\|^2 \leq \left(2\lambda_1^2 \sum_{i=2}^k \lambda_i^{-2}\right) \left| \frac{\lambda_2 \underline{v}_2^T \cdot \underline{x}^0}{\lambda_1 \underline{v}_1^T \cdot \underline{x}^0} \right|^{2^{t+1}}.$$

and similarly for

$$\lambda^t = T(\underline{x}^t, \underline{x}^t, \underline{x}^t) \rightarrow \lambda_1 \text{ fast.}$$

Proof.

First note that by definition:

$$T(I, x^t, x^t) = \sum_{i=1}^k \lambda_i \underline{w}_i (\underline{w}_i^T, x^t)^2$$

By looking at iteration t & $t-1$ we check (as far as I can) that

$$\underline{x}^t = \frac{\sum_{i=1}^k \lambda_i^{\frac{3}{2}} (\underline{w}_i^T \cdot x^{t-1})^2 \underline{w}_i}{\left\| \sum_{i=1}^k \lambda_i^{\frac{3}{2}} (\underline{w}_i^T \cdot x^{t-1})^2 \underline{w}_i \right\|_2}$$

$$= \frac{\sum_{i=1}^k \lambda_i^{\frac{3}{2}} (\underline{w}_i^T \cdot x^{t-2})^4 \underline{w}_i}{\left\| \sum_{i=1}^k \lambda_i^{\frac{3}{2}} (\underline{w}_i^T \cdot x^{t-2})^4 \underline{w}_i \right\|_2}$$

$$= \frac{\sum_{i=1}^k \lambda_i^{\frac{2t-1}{2}} (\underline{w}_i^T \cdot x^{(0)})^{2^t} \underline{w}_i}{\left\| \sum_{i=1}^k \lambda_i^{\frac{2t-1}{2}} (\underline{w}_i^T \cdot x^{(0)})^{2^t} \underline{w}_i \right\|_2}$$

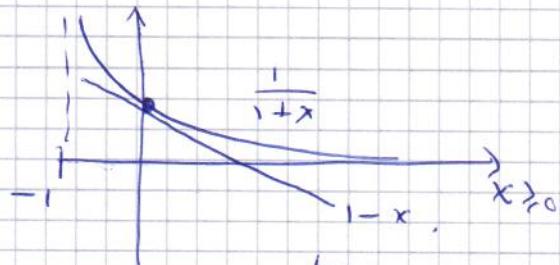
Multiplying by \underline{w}_i^T we find

$$(\underline{w}_i^T \cdot \underline{x}^t)^2 = \frac{\lambda_i^{\frac{2t+1}{2}} \cdot (\underline{w}_i^T \cdot \underline{x}^{(0)})^{2^{t+1}}}{\sum_{i=1}^k \lambda_i^{\frac{2t+1}{2}} (\underline{w}_i^T \cdot \underline{x}^{(0)})^{2^{t+1}}}.$$

$$\Rightarrow (\underline{v}_1^T \cdot \underline{x}^t)^2 = \frac{1}{1 + \sum_{i=2}^K \left(\frac{\lambda_i \underline{v}_i^T \cdot \underline{x}^0}{\lambda_1 \underline{v}_1^T \cdot \underline{x}^0} \right)^2 \left(\frac{\lambda_1}{\lambda_i} \right)^{2^{t+1}}}$$

$$\geq 1 - \lambda_1^2 \sum_{i=2}^K \lambda_i^{-2} \left(\frac{\lambda_i \underline{v}_i^T \cdot \underline{x}^0}{\lambda_1 \underline{v}_1^T \cdot \underline{x}^0} \right)^{2^{t+1}}$$

where we used $x \geq 0 \Rightarrow \frac{1}{1+x} \geq 1-x$.



$$\Rightarrow 1 - (\underline{v}_1^T \cdot \underline{x}^t)^2 \leq \lambda_1^2 \left(\sum_{i=2}^K \lambda_i^{-2} \right) \cdot \left(\frac{\lambda_2 \underline{v}_2^T \cdot \underline{x}^0}{\lambda_1 \underline{v}_1^T \cdot \underline{x}^0} \right)^{2^{t+1}}$$

$$1 - (\cos \theta_t)^2 \geq 1 - \cos \theta_t \quad (\text{for } \cos \theta_t < 1).$$

$$\Rightarrow 2(1 - \cos \theta_t) \leq 2 \cdot \lambda_1^2 \left(\sum_{i=2}^K \lambda_i^{-2} \right) \left(\frac{\lambda_2 \underline{v}_2^T \cdot \underline{x}^0}{\lambda_1 \underline{v}_1^T \cdot \underline{x}^0} \right)^{2^{t+1}}$$

$\underbrace{\quad}_{\|\underline{x}^t - \underline{v}_1\|_2^2}$

For $\lambda_t = T(\underline{x}^t, \underline{x}^t, \underline{x}^t)$ we have

$$\lambda_t = \sum_{i=1}^K \lambda_i (\underline{v}_i \cdot \underline{x}^t)^3$$

$$= \lambda_1 (\underline{v}_1^T \cdot \underline{x}^t)^3 \left\{ 1 + \sum_{i=2}^K \frac{\lambda_i (\underline{v}_i^T \cdot \underline{x}^t)^3}{\lambda_1 (\underline{v}_1^T \cdot \underline{x}^t)^3} \right\}$$

$\rightarrow 0$ for $\underline{v}_i^T \cdot \underline{v}_1 = 0$
 $i \geq 2$

$$\rightarrow \lambda_1.$$

III. b) Whitening procedure

Often in applications we have

$$\mathbf{T} = \sum_{i=1}^k \lambda_i \underline{u}_i \otimes \underline{u}_i \otimes \underline{u}_i$$

where \underline{u}_i 's ~~are~~ are not necessarily orthogonal.

Before applying the power method we would like to change basis in order to transform \mathbf{T} to an "orthogonal tensor".

This is possible when

$$\mathbf{M} = \sum_{i=1}^k \lambda_i \underline{u}_i \otimes \underline{u}_i \otimes \underline{u}_i = \sum_{i=1}^k \lambda_i \underline{u}_i \underline{u}_i^\top$$

is available (e.g. Topic Models, LDA),

Assumption: $[\underline{u}_1, \dots, \underline{u}_K]$ are linearly independent and $\lambda_1, \dots, \lambda_K$ are strictly positive.

Then \mathbf{M} is a symmetric semi-definite positive (square) matrix and we have

$$\mathbf{M}_2 = \mathbf{U} \mathbf{D} \mathbf{U}^\top$$

with $\mathbf{D} = \text{diag}(\mathbf{d})$ $d_i > 0$. Since \mathbf{M}_2 is rank K we have $d_1 > d_2 > \dots > d_K > 0$ (with other possible eigenvalues zero)

$$\Rightarrow M_2 = U_{N \times K} \text{Diag}(d_1 \dots d_K) U^T_{K \times N}$$

where $U_{N \times K} = [\underline{u}_1 \dots \underline{u}_K]$ are orthonormal eigen vectors,

Define : $W = U_{N \times K} \text{Diag}(d_1^{-1/2} \dots d_K^{-1/2})$.

We claim that the vectors

$$\underline{\nu}_i = \sqrt{d_i} W^T \underline{u}_i \quad i=1 \dots K \quad (*)$$

form an orthonormal system of vectors. (Proof below).

Whitening of T is the operation:

$$\star T(W, W, W) = \sum_{i=1}^K \lambda_i W^T \underline{u}_i \otimes W^T \underline{u}_i \otimes W^T \underline{u}_i$$

(a multilinear transformation to a new tensor)

$$= \sum_{i=1}^K \lambda_i^{-1/2} \underline{\nu}_i \otimes \underline{\nu}_i \otimes \underline{\nu}_i$$

After applying e.g. the Tensor Power method to the whitened tensor we retrieve $\lambda_i^{-1/2}$ and $\underline{\nu}_i$. Then we get \underline{u}_i from $\underline{\nu}_i$ as follows :

$$\underline{u}_i = \frac{1}{\sqrt{d_i}} U_{N \times K} \text{diag}(d_1^{1/2} \dots d_K^{1/2}) \underline{\nu}_i$$

as can be checked from (*) above.

Proof of Claim

$$\text{Recall } M_2 = U_{N \times K} \text{Diag}(d_1, \dots, d_K) U_{K \times N}^T$$

$$\text{and by definition } W = U_{N \times K} \text{Diag}(d_1^{-1/2}, \dots, d_K^{-1/2})$$

Thus:

$$\begin{aligned} W^T M_2 W &= \text{Diag}(d_1^{-1/2}, \dots, d_K^{-1/2}) U_{K \times N}^T U_{N \times K} \\ &\quad \cdot \text{Diag}(d_1, \dots, d_K) U_{K \times N}^T U_{N \times K} \\ &\quad \cdot \text{Diag}(d_1^{-1/2}, \dots, d_K^{-1/2}) \end{aligned}$$

$$\Rightarrow W^T M_2 W = I_{K \times K}.$$

which means:

$$\sum_{i=1}^K \underline{\mu}_i (W^T \underline{\mu}_i) (\cancel{W^T \mu_i} W) = I$$

$$\Rightarrow \sum_{i=1}^K (\sqrt{d_i} W^T \underline{\mu}_i) (\sqrt{d_i} \cancel{W^T \mu_i} W) = I$$

$$\Rightarrow \sum_{i=1}^K (\sqrt{d_i} W^T \underline{\mu}_i) (\sqrt{d_i} W^T \underline{\mu}_i)^T = I_{K \times K}.$$

Since $\{\underline{\mu}_1, \dots, \underline{\mu}_K\}$ are lin indep and W^T is full rank

$[W^T \underline{\mu}_1, \dots, W^T \underline{\mu}_K]$ are also lin indep. Thus this equation implies they must be orthogonal.

For this last point in more details:

Let $\underline{w}^T \underline{\mu}_i = \underline{w}_i$, Since $\underline{\mu}_1 \dots \underline{\mu}_K$

form a basis of K -dim space and $\text{Rank}(W) = K$

The \underline{w}_i also form a basis of K -dim space.

Moreover $\sum_{i=1}^K \underline{w}_i \underline{w}_i^T = I_{K \times K}$ implies

$\sum_{i=1}^K \underline{w}_i (\underline{w}_i^T \underline{w}_l) = \underline{w}_l$, Thus we must

have $\underline{w}_i \cdot \underline{w}_l = \delta_{il}$, i.e they are orthonormal.