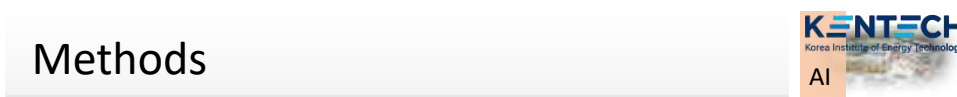


1



- Convolution methods
- Spectral convolution
- Spatial convolution

2

GCN – Math Background



- Graph signal processing
- Normalized graph Laplacian
 - $L = I - D^{-1/2}AD^{-1/2}$
 - L is symmetric and PSD $\rightarrow L = U\Lambda U^T$ where U is matrix of eigenvalues and Λ is diagonal matrix of eigenvalues
 - $UU^T = I$
- Graph Fourier transform
 - $F(x) = U^T x = \hat{x}$
 - Inverse Fourier transform: $x = F^{-1}(\hat{x}) = U\hat{x}$
 $\rightarrow x = \sum_i \hat{x}_i u_i$
- Graph convolution with filter, g
 - $x * g = F^{-1}(F(x) \odot F(g)) = U(U^T x \odot U^T g)$ where \odot is element-wise product
 - Denote the filter as $g_\theta = \text{diag}(U^T g)$
 $\rightarrow x * g_\theta = U g_\theta U^T x$

GCN – Math Background

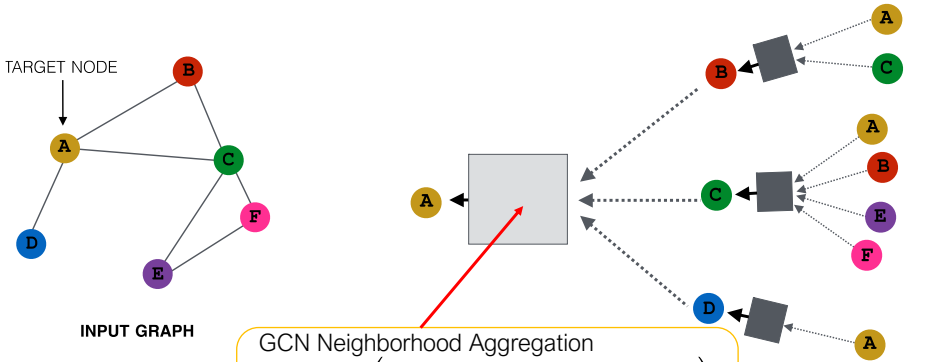


- Chebyshev approximation
 - Approximate g_θ by Chebyshev polynomials of U
 - $g_\theta = \sum_{i=1}^K \theta'_i T_i(\tilde{\Lambda})$ where $\tilde{\Lambda} = \frac{2}{\lambda_{max}} \Lambda - I$ and $\theta'_i \in R^K$ is a vector of Chebyshev coefficients
- $$\begin{aligned} T_0(x) &= 0, T_1(x) = x \\ T_i(x) &= 2xT_{i-1}(x) - T_{i-2}(x) \end{aligned}$$
- $\rightarrow x * g_\theta = U(\sum_{i=1}^K \theta'_i T_i(\tilde{\Lambda}))U^T x = \sum_{i=1}^K \theta'_i T_i(\tilde{L}) x$
- Further simplification
 - $x * g_\theta = \sum_{i=1}^K \theta'_i T_i(\tilde{L}) x$
 - Let $K=1, \lambda_{max} = 2$
 $\rightarrow x * g_\theta = \theta_0 x - \theta_1 D^{-1/2} A D^{-1/2} x$
 - Let $\theta = \theta_0 = -\theta_1$
 $\rightarrow x * g_\theta = \theta(I + D^{-1/2} A D^{-1/2}) x$
- $$\text{Let } \tilde{L} = \frac{2}{\lambda_{max}} L - I$$

$$T_i(\tilde{L}) = U T_i(\tilde{\Lambda}) U^T$$

Proof by induction
- $$\text{Kipf \& Welling, "Semi-supervised classification with graph convolutional networks," ICLR, 2017.}$$
- In case of multi-channel input & output
 $\rightarrow H = X * g_\theta = f(\tilde{A} X \theta)$ where $\tilde{A} = I + D^{-1/2} A D^{-1/2}$, $f(\cdot)$ an activation function

GCN




TARGET NODE

INPUT GRAPH

GCN Neighborhood Aggregation

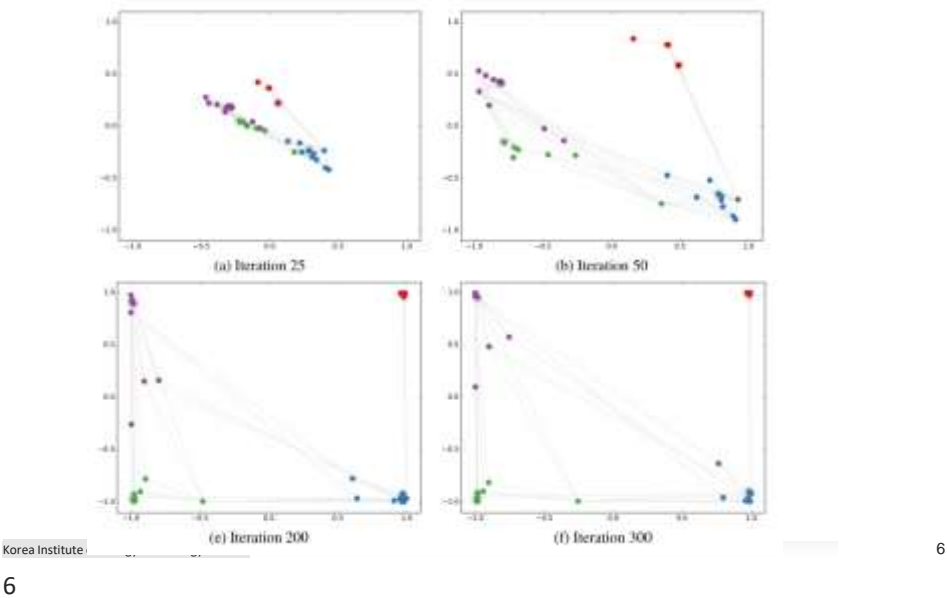
$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)| |N(v)|}} \right)$$



Results

- Node Embedding of GCN





GCN-Results



- Semi-supervised node classification (Accuracy)

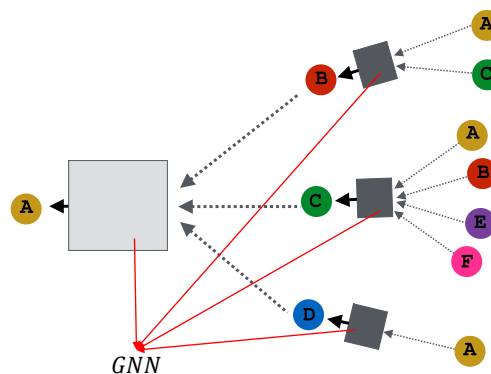
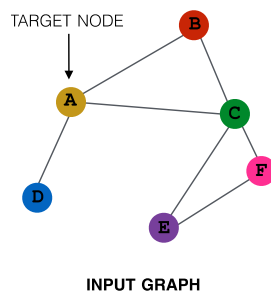
Method	Citeseer	Cora	Pubmed	NELL
ManiReg [3]	60.1	59.5	70.7	21.8
SemiEmb [28]	59.6	59.0	71.1	26.7
LP [32]	45.3	68.0	63.0	26.5
DeepWalk [22]	43.2	67.2	65.3	58.1
ICA [18]	69.1	75.1	73.9	23.1
Planetoid* [29]	64.7 (26s)	75.7 (13s)	77.2 (25s)	61.9 (185s)
GCN (this paper)	70.3 (7s)	81.5 (4s)	79.0 (38s)	66.0 (48s)
GCN (rand. splits)	67.9 ± 0.5	80.1 ± 0.5	78.9 ± 0.7	58.4 ± 1.7

GraphSAGE



Hamilton, Ying, & Leskovec, "Inductive representation learning on large graphs," NIPS, 2017

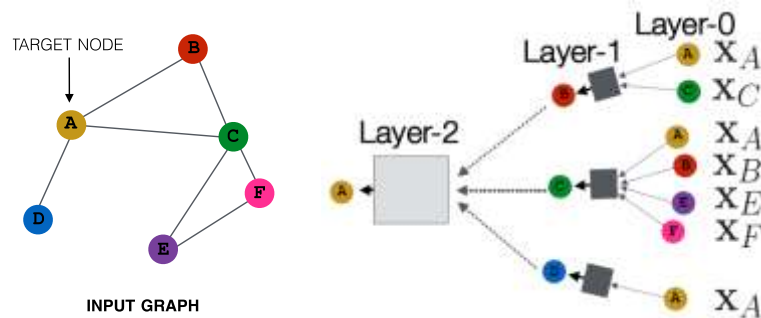
- Compute node embedding based on local network neighborhoods
- Aggregation strategies
 - Max
 - **Average**



Stacking Layers



- Model can be of arbitrary depth (Layers)
 - Nodes have embedding at each layer
 - Layer-0 embedding of node u is its input feature, x_u
 - Layer-K embedding gets information from nodes that are K hops away



Korea Institute of Energy Technology

2022-04-11

9

9

Average Aggregation



- Average information(message) from neighbors

$$\mathbf{h}_v^0 = \mathbf{x}_v \quad (\text{Initialize node embedding})$$

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right), \quad \forall k \in \{1, \dots, K\} \quad (\text{Update})$$

$$\mathbf{z}_v = \mathbf{h}_v^K \quad (\text{After } k - \text{step aggregation})$$

- After K-layers of neighborhood aggregation, we get output embedding for each node (\mathbf{z}_v)
- We can feed these embedding into any loss function
- Run stochastic gradient descent to train

Korea Institute of Energy Technology

2022-04-11

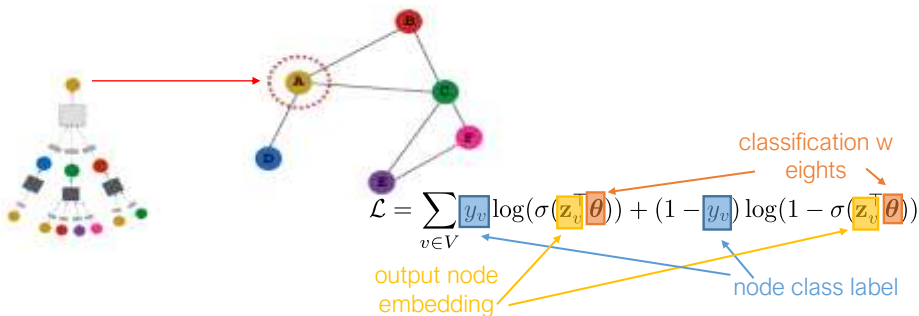
10

10

Training



- Unsupervised training
 - When we do not have label, only graph structure can be used
 - Goal: make similar nodes have similar embeddings
 - Node proximity in the graph
- Supervised training
 - Ground truth or manually assigned labels
 - Directly train the model to make the input has the output of right label



Korea Institute of Energy Technology

2022-04-11

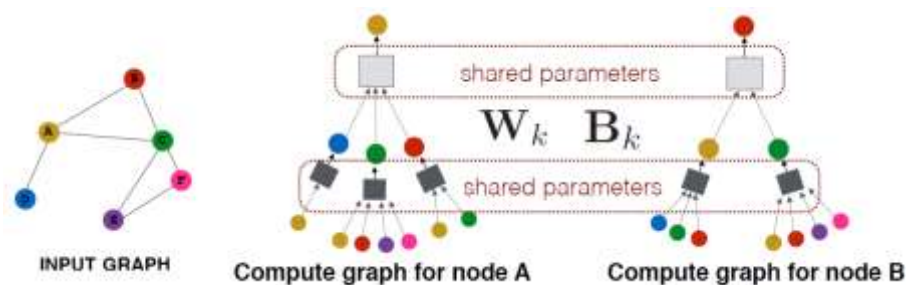
11

11

Inductiveness



- The same aggregation parameters are shared for all nodes:
 - The number of model parameters is sublinear in $|V|$ and we can **generalize to unseen nodes**



Korea Institute of Energy Technology

2022-04-11

12

12

GCN & GraphSAGE



- Graph Convolutional Network (GCN) and GraphSAGE are slightly differ

Basic Neighborhood Aggregation

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|} + \mathbf{B}_k \mathbf{h}_v^{k-1} \right)$$

VS.

GCN Neighborhood Aggregation

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)| |N(v)|}} \right)$$

same matrix for self and neighbor embeddings

per-neighbor normalization

Graph Convolutional Network



- Empirically, GCN discovers the configuration that births the best results

- More parameter sharing
- Down-weights high degree neighbors

$$\mathbf{h}_v^k = \sigma \left(\mathbf{W}_k \sum_{u \in N(v) \cup v} \frac{\mathbf{h}_u^{k-1}}{\sqrt{|N(u)| |N(v)|}} \right)$$

use the same transformation matrix for self and neighbor embeddings

instead of simple average, normalization varies across neighbors

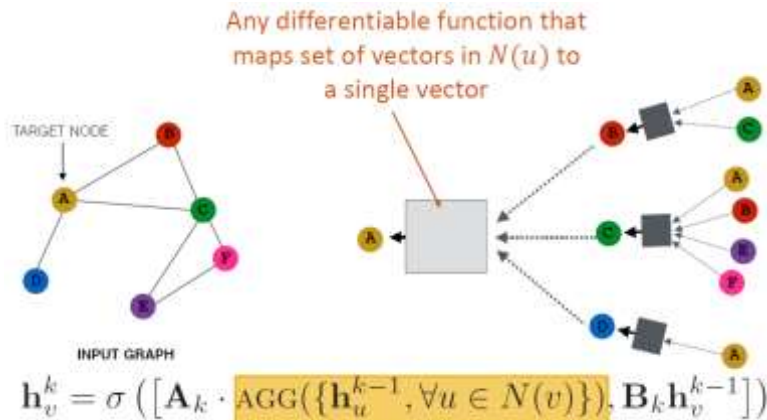
$$\mathbf{H}^{(k+1)} = \sigma \left(\mathbf{D}^{-\frac{1}{2}} \tilde{\mathbf{A}} \mathbf{D}^{-\frac{1}{2}} \mathbf{H}^{(k)} \mathbf{W}_k \right)$$

where $\tilde{\mathbf{A}} = \mathbf{A} + \mathbf{I}$

$$\mathbf{D}_{ii} = \sum_j \mathbf{A}_{i,j}$$

GraphSAGE

- Aggregation function
 - Average



Neighbor Aggregation

- Variants of Aggregation function

- Mean
$$\text{AGG} = \sum_{u \in N(v)} \frac{\mathbf{h}_u^{k-1}}{|N(v)|}$$
- Pool
$$\text{AGG} = \gamma(\{\mathbf{Q}\mathbf{h}_u^{k-1}, \forall u \in N(v)\})$$

element-wise mean/max
- LSTM
$$\text{AGG} = \text{LSTM}([\mathbf{h}_u^{k-1}, \forall u \in \pi(N(v))])$$

- GAT (Graph Attention Network)

- Both GCN & GraphSAGE use fixed weights
- GAT automatically adjusts the importance (weights) of neighbors

Velickovic, et.al., "Graph attention networks," ICLR, 2017.

GCN – appendix



● GCN

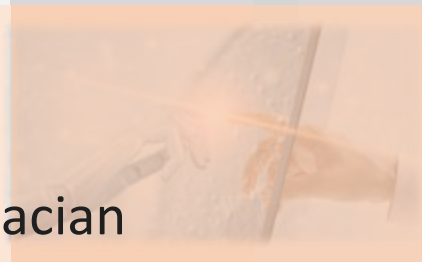
- $H_{:,j}^{(k)} = \sigma \left(\sum_{i=1}^{f_{k-1}} U \theta_{i,j}^{(k)} U^T H_{:,i}^{(k-1)} \right)$, $j = 1, 2, \dots, f_k$ with learnable filter,

$g_\theta = \theta_{i,j}^{(k)}$, σ is an activation function such as ReLU

- k : layer index
- $H^{(k-1)} \in R^{n \times f_{k-1}}$: input signal with f_{k-1} channels
- $\theta_{i,j}^{(k)}$: diagonal matrix

→ Complexity $O(n^3)$ for eigen-decomposition

Appendix: Graph Laplacian



Spectral

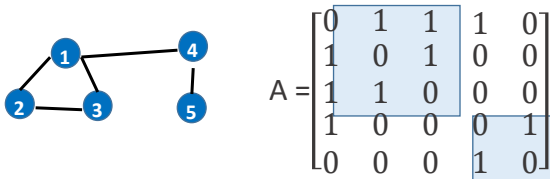


- In signal processing
 - Fourier transform
 - Basis: Elementary sine and cosine waves of different frequencies
 - Represent signal as a sum of basis
- In graph
 - Basis: eigenvectors (values) of graph Laplacian matrix (L, defined later)
 - Eigenvector: Orthogonal components of graph
 - Spectral: Eigen-decomposition of L

Spectral Analysis



- Adjacency matrix, A where $A_{ij} = 1$ if edge $(i, j) \in E$

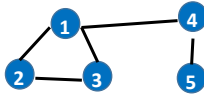


Symmetric
Has n real eigenvalues
 n Eigenvectors are real valued
and orthogonal

- Eigenpair (eigenvector, eigenvalue) of A: $Ax_i = \lambda_i \cdot x_i$
- Spectrum: Eigenvectors of matrix A ordered by the magnitude (strength) of their eigenvalues

$$\Lambda = (\lambda_1, \lambda_2, \lambda_3, \dots, \lambda_n), \lambda_i \leq \lambda_{i+1}$$
- Analyze the spectrum

Adjacency and Degree Matrices



$$D = \begin{bmatrix} 3 & 0 & 0 & 0 & 0 \\ 0 & 2 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 \\ 0 & 0 & 0 & 2 & 0 \\ 0 & 0 & 0 & 0 & 1 \end{bmatrix}$$

Diagonal
Non-singular and positive
definite (if no isolated node)

- N orthogonal eigenpairs

- $\lambda_i = k_i$, $x_i = (0, \dots, 1, 0, \dots, 0)$

i-th element

How about D-A?

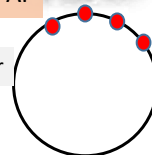
Example: d-regular graph



- d-regular graph

- Every node has d edges

2-regular



- Connected d-regular graph and its adjacency matrix, A

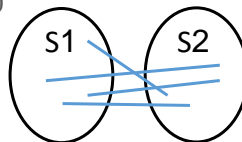
→ $x_n = \mathbf{1} = (1, 1, \dots, 1)$ is the largest (principal) eigenvector with eigenvalue $\lambda_n = d$

- The second largest (secondary) eigenvector, x_{n-1}

- Eigenvectors, x_n and x_{n-1} , are orthogonal

→ $\sum_k x_n[k] \cdot x_{n-1}[k] = 0$

→ $\sum_k x_{n-1}[k] = 0$



Assume $|S1| = |S2|$

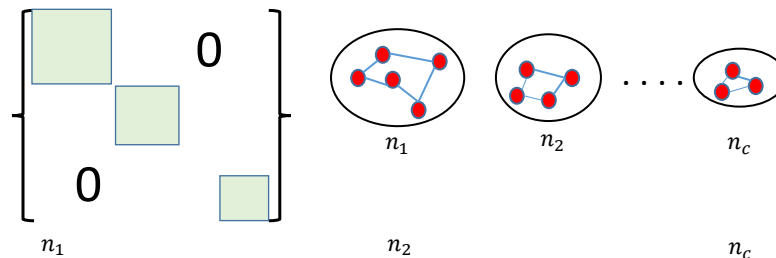
$$x_{n-1}[k] = \begin{cases} 1, & \text{if } k \in S1 \\ -1, & \text{if } k \in S2 \end{cases}$$

Now suppose that S1 and S2 are sparsely connected
Then eigenvector, x_{n-1} , is a good partition
How to find such eigenvector?

Example: Components



- Consider a graph with c components of n_1, n_2, \dots, n_c nodes
 - Let each component is a d -regular graph
- Adjacency matrix

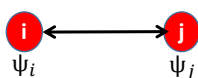


- $x_1 = (1, 1, \dots, 1, 0, \dots, 0)$, $x_2 = (0, \dots, 0, 1, \dots, 1, 0, \dots, 0)$, $x_c = (0, \dots, 0, 1, \dots, 1)$ are eigenvectors with eigenvalue $\lambda_{n-c+1} = \lambda_{n-c+2} = \dots = \lambda_n = d$

Graph Laplacian



- Diffusion (such as pressure, heat, ...) between connected nodes



Diffusion constant

$$\begin{aligned} \frac{d\psi_i(t)}{dt} &= C \cdot \sum_j A_{ij} (\psi_j(t) - \psi_i(t)) \\ &= C \cdot (\sum_j A_{ij} \psi_j(t) - \sum_j A_{ij} \psi_i(t)) \\ &= C \cdot \sum_j (A_{ij} - \delta_{ij} k_i) \psi_j(t) \end{aligned}$$

- In a matrix form, $\frac{d\psi(t)}{dt} = C(\mathbf{A} - \mathbf{D})\psi(t)$

- $\frac{d\psi(t)}{dt} + C(\mathbf{D} - \mathbf{A})\psi(t) = 0$

- $\frac{d\psi(t)}{dt} + C\mathbf{L}\psi(t) = 0$, where

$$L_{ij} = \begin{cases} k_i, & \text{if } i = j \\ -1, & \text{if } i \neq j \text{ and } (i, j) \in E \\ 0, & \text{o.w.} \end{cases}$$

Laplacian operator

$$\text{Diffusion equation: } \frac{d\psi(t)}{dt} + C\nabla^2\psi(t) = 0$$

Graph Laplacian



• Properties of L

- Symmetric
- Singular
 $A\mathbf{1}=\mathbf{0} \rightarrow$ Eigenvector $\mathbf{1}=(1,1,\dots,1)$ with $\lambda=0$
- Positive semidefinite
 \rightarrow all eigenvalues are nonnegative
 \rightarrow Eigenvectors are orthogonal

• Write $\psi(t)$ as a linear combination of the eigenvectors

$$\psi(t) = \sum_i a_i(t) \mathbf{v}_i$$

$$\rightarrow \sum_i \left(\frac{da_i}{dt} + C\lambda_i a_i \right) \mathbf{v}_i = 0$$

$$\rightarrow \frac{da_i}{dt} + C\lambda_i a_i = 0$$

$$\rightarrow a_i(t) = a_i(0)e^{-C\lambda_i t}$$

$$\frac{d\psi}{dt}(t) + C\mathbf{L}\psi(t) = 0$$

Korea Institute of Energy Technology

25

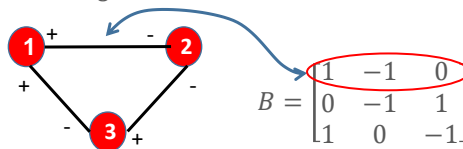
25

Spectrum of Graph Laplacian



- Arbitrarily assign +/- signs to two ends of an edge
- Define a matrix \mathbf{B} (Edge incidence matrix) such that

$$-B_{ij} = \begin{cases} +1, & \text{if } (i,j) \in E \text{ and } i \text{ is a positive end} \\ -1, & \text{if } (i,j) \in E \text{ and } i \text{ is a negative end} \\ 0, & \text{o.w.} \end{cases}$$



• Claim: $\mathbf{L} = \mathbf{B}^T \mathbf{B}$

- $(\mathbf{B}^T \mathbf{B})_{ij} = \sum_k \mathbf{B}_{ki} \cdot \mathbf{B}_{kj}$
- $= \mathbf{L}_{ij}$

• Consider an eigenpair $(\mathbf{v}_i, \lambda_i)$ of \mathbf{L}

$$\rightarrow \mathbf{v}_i^T \mathbf{L} \mathbf{v}_i = \lambda_i \mathbf{v}_i^T \mathbf{v}_i = \lambda_i$$

$$\mathbf{v}_i^T \mathbf{B}^T \mathbf{B} \mathbf{v}_i = (\mathbf{B} \mathbf{v}_i)^T \mathbf{B} \mathbf{v}_i$$

Inner product of a real vector \rightarrow Nonnegative

Korea Institute of Energy Technology

26

26