# CS-E4740 Federated Learning

# "FL Design Principle"

## Dipl.-Ing. Dr.techn. Alexander Jung

# What are the main components of ML and how are they combined?
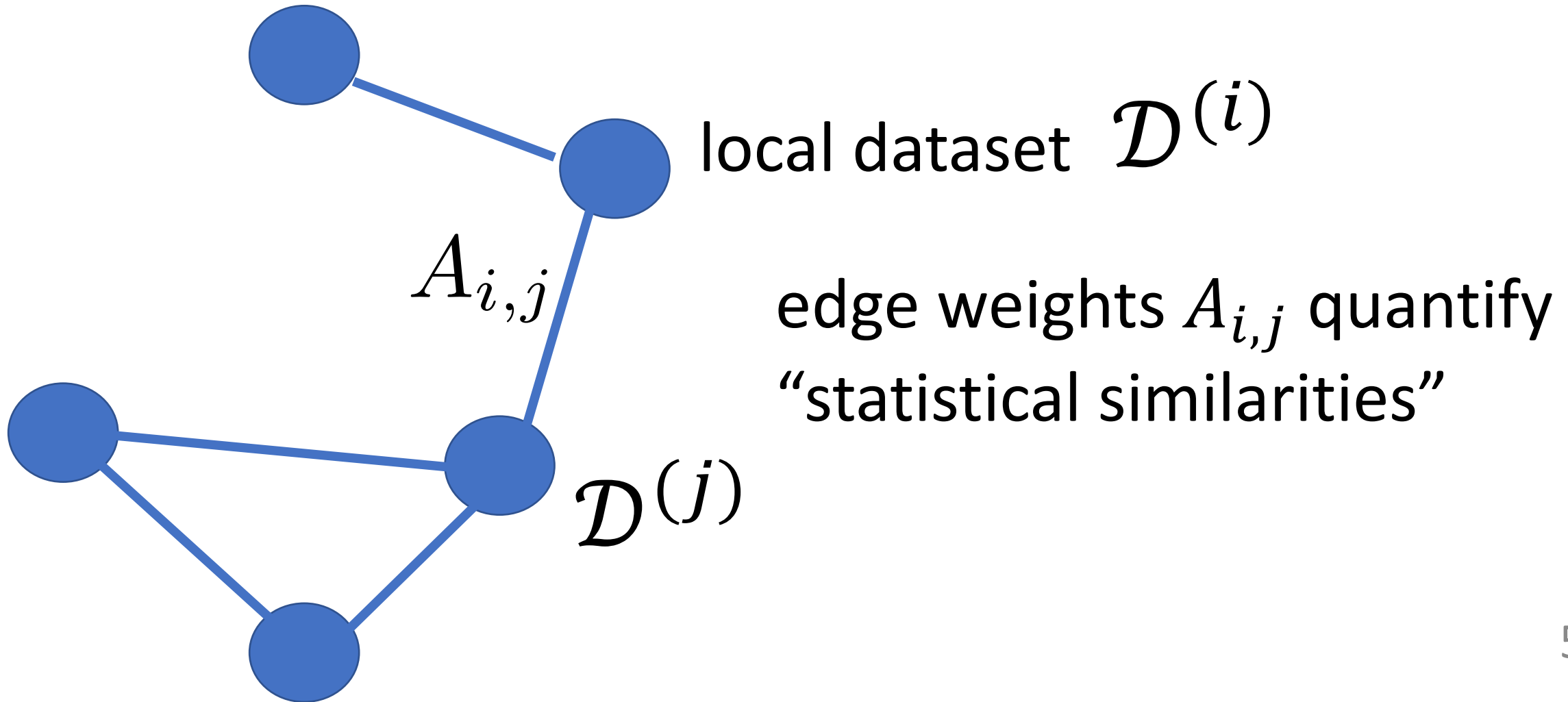
# Previous Lecture:
# Networked Data and Models

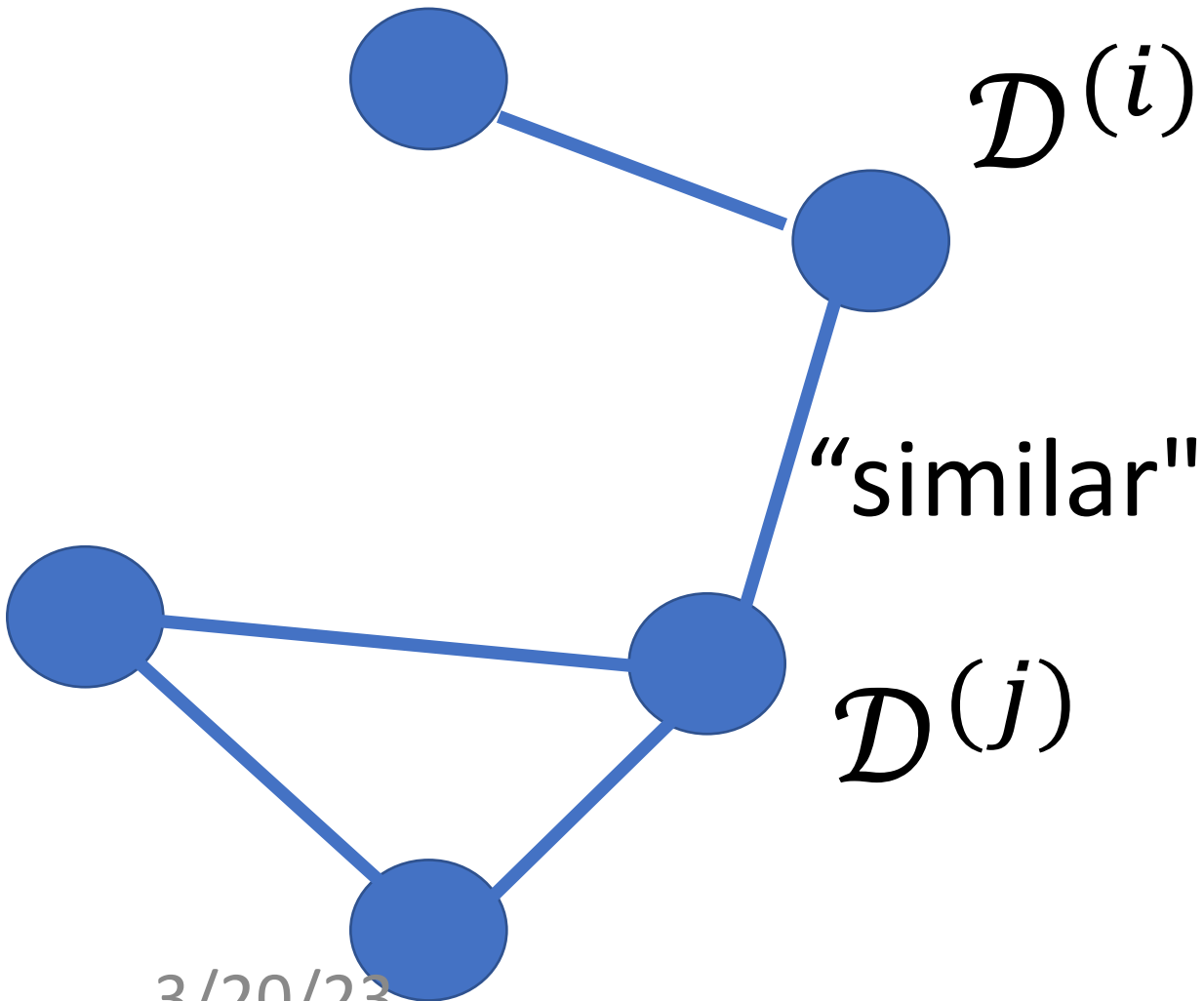# Today:
# Loss and Optimization

# Weather Stations

# The Empirical Graph



local dataset $\mathcal{D}^{(i)}$

$A_{i,j}$

edge weights $A_{i,j}$ quantify "statistical similarities"

$\mathcal{D}^{(j)}$

# Networked Models.



$\mathcal{D}^{(i)}$

"similar"

$\mathcal{D}^{(j)}$

local model for each node

couple models at connected nodes

# Local Parametric Models



$\mathcal{D}^{(i)}$

model params $\mathbf{w}^{(i)}$
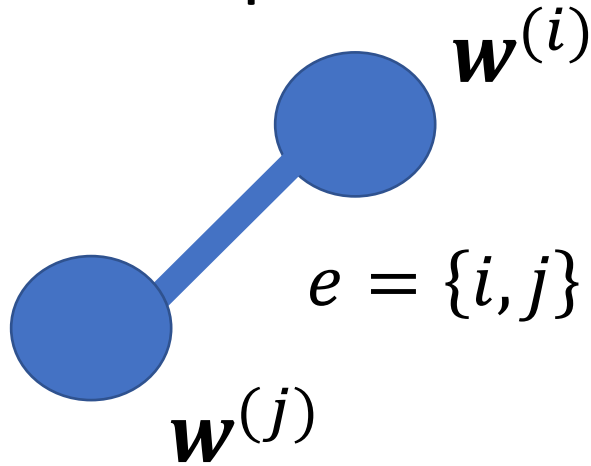
"similar"

$\mathbf{w}^{(j)}$

# Clustering Assumption



the local datasets form clusters

datasets in same can be approximated as realizations of i.i.d. RVs with prob. dist p(x,y;c)
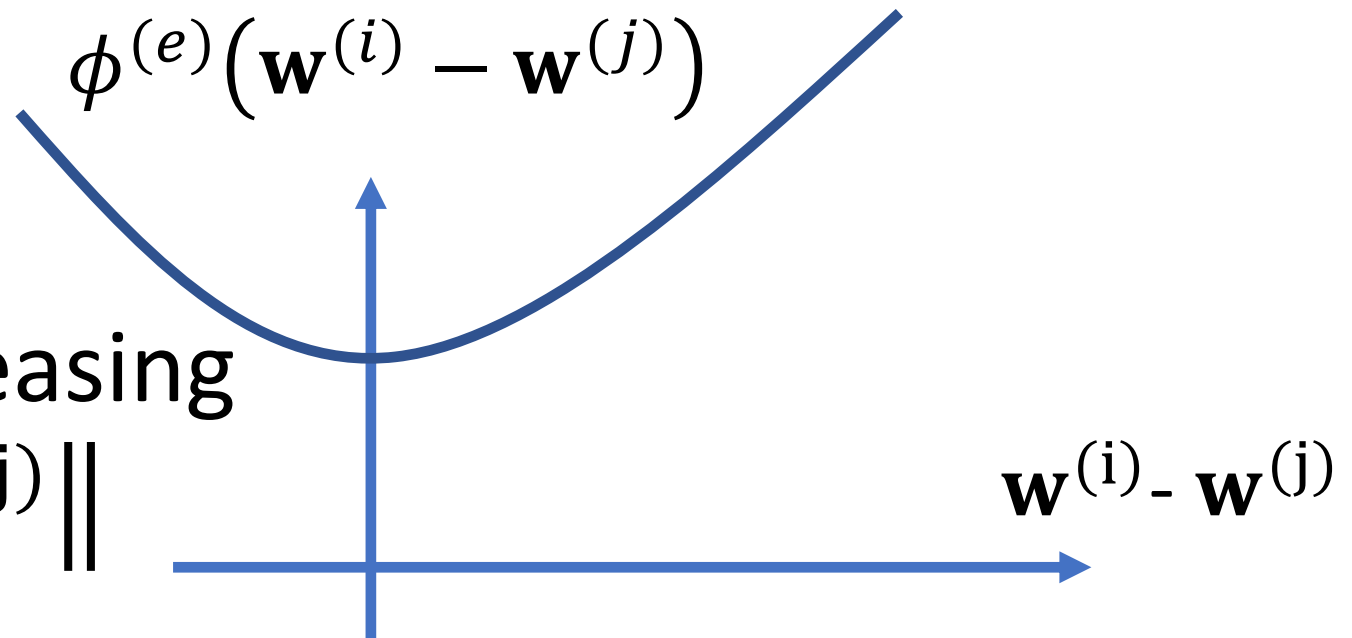
more edges inside clusters

# Measure Clustering via Variation
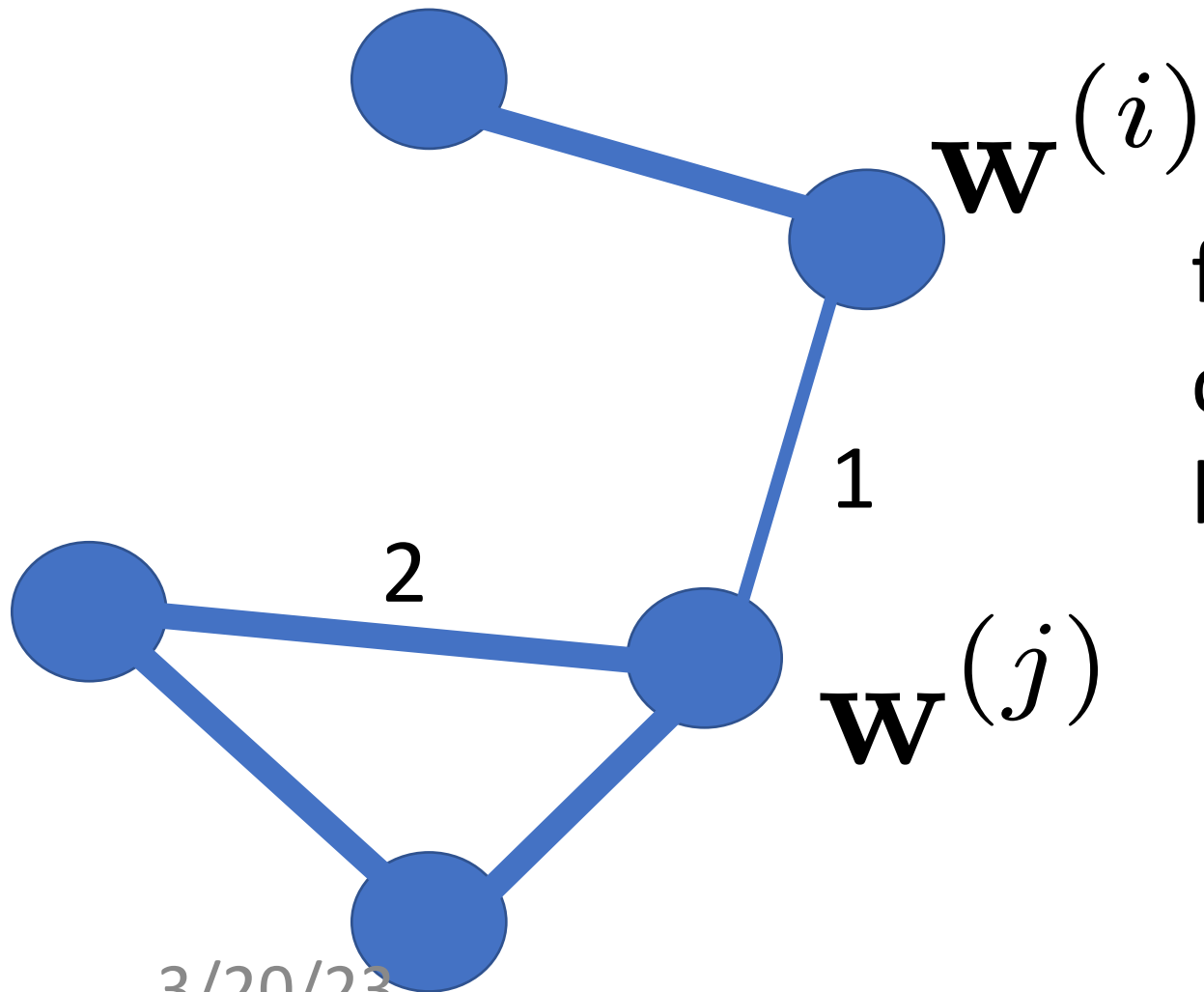
local model params



$$\boldsymbol{w}^{(i)}$$

$$e = \{i, j\}$$

$$\boldsymbol{w}^{(j)}$$

require similar params at ends of edge e

penalty function measures "variation"

$$\phi^{(e)}\big(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\big)$$



$\phi^{(e)}$ convex and increasing with norm $\|\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\|$

$$\mathbf{w}^{(i)} - \mathbf{w}^{(j)}$$

# Generalized Total Variation (GTV)



force model params at well connected nodes to be similar by requiring small GTV

$$\sum_{\{i,j\}} A_{i,j} \phi\left(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\right)$$
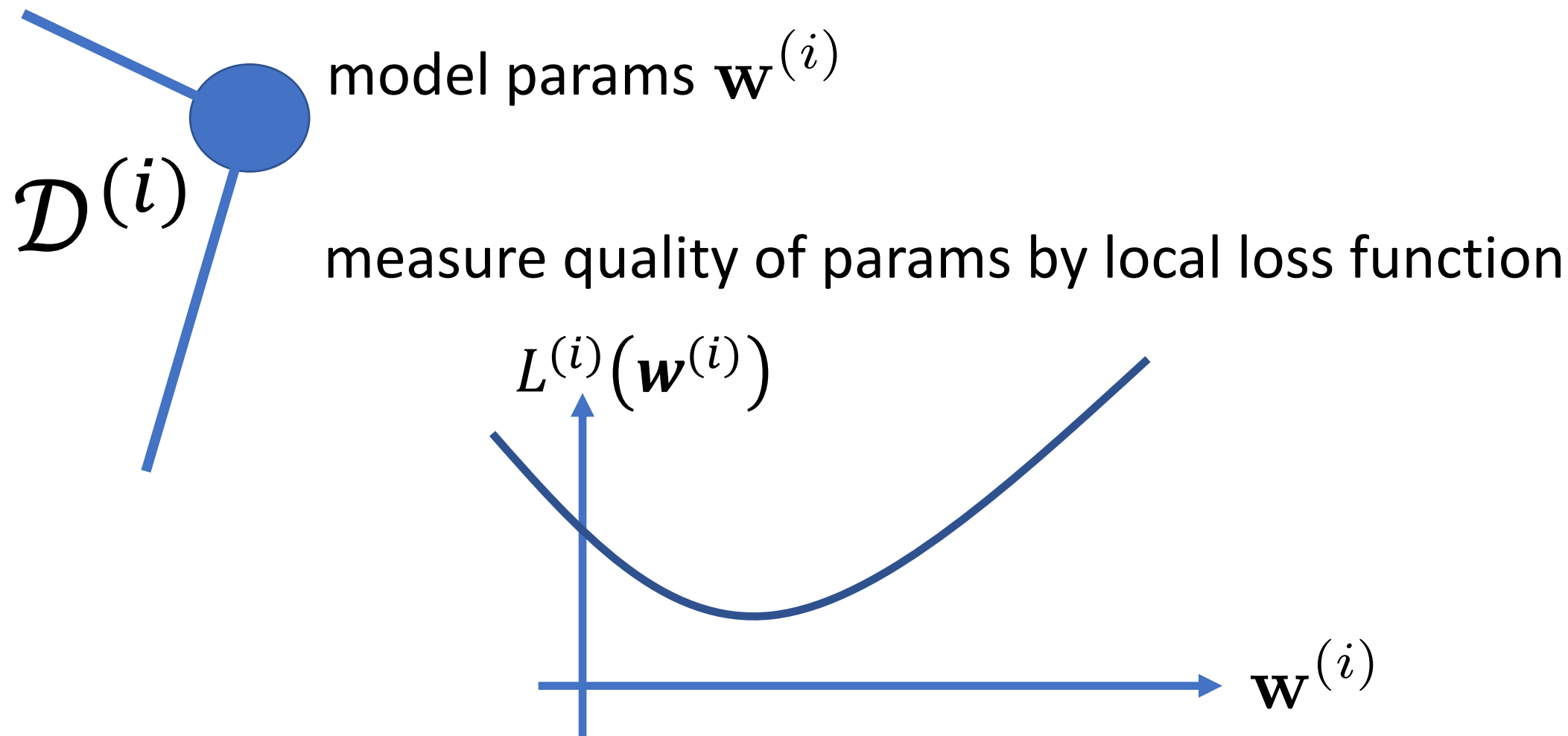
# Two Special Cases of GTV

total variation  $\phi(\mathbf{u}) = \|\mathbf{u}\|_2$
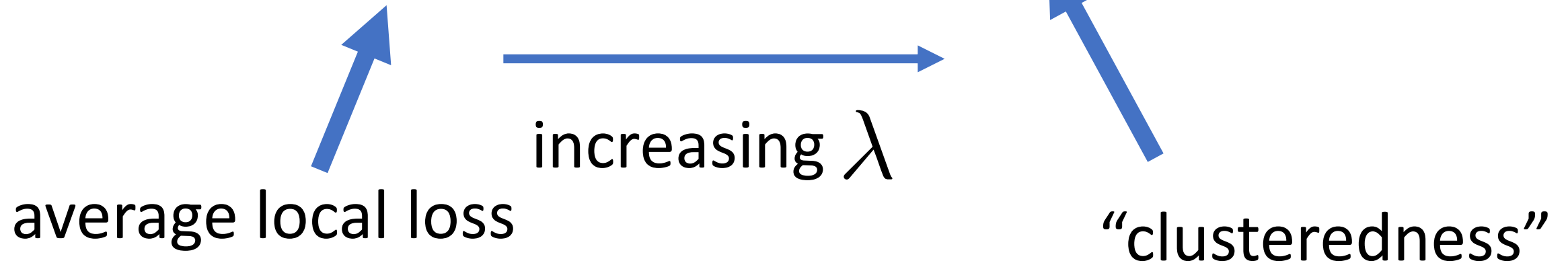
graph Laplacian quadratic from is GTV with
$$\phi(\mathbf{u}) = \|\mathbf{u}\|_2^2$$

# Local Loss Function



model params $\mathbf{w}^{(i)}$

measure quality of params by local loss function

$L^{(i)}(\boldsymbol{w}^{(i)})$

$\mathcal{D}^{(i)}$

$\mathbf{w}^{(i)}$

# GTV Minimization

$$\min_{\mathbf{w}} \sum_i L^{(i)}\big(\mathbf{w}^{(i)}\big) + \lambda \sum_{\{i,j\}} A_{i,j}\phi\big(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\big)$$

increasing $\lambda$

average local loss

"clusteredness"

# Network Lasso

$$\min_{\mathbf{w}} \sum_i L^{(i)}\left(w^{(i)}\right) + \lambda \sum_{\{i,j\}} A_{i,j} \left\| w^{(i)} - w^{(j)} \right\|$$

# Special Case: "MOCHA"

$$\min_{w} \sum_{i} L^{(i)}\left(w^{(i)}\right) + \lambda \sum_{\{i,j\}} A_{i,j} \left\| w^{(i)} - w^{(j)} \right\|^2$$
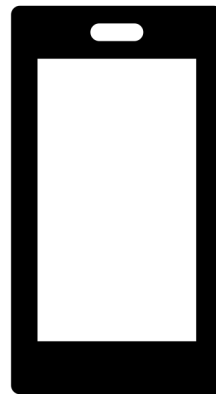
# Two Key Questions of ML

$$\min_{\mathbf{w}} \sum_i L^{(i)}\big(\mathbf{w}^{(i)}\big) + \lambda \sum_{\{i,j\}} A_{i,j}\phi\big(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\big)$$

- computational aspects: how to compute (approximate) solutions efficiently ?

- statistical aspects: are the solutions any good?

# Computational Aspects

# A FL Setting

# Requirements

- run in ad-hoc nets of low-cost devices

- robustness against node/link failures

- robustness against "stragglers"

# Another FL Setting...

https://www.google.com/about/datacenters/







https://en.wikipedia.org/wiki/Optical_fiber

# GTV Min. for Local Lin.Reg.

$$\min_{\mathbf{w}} \sum_i \|\mathbf{X}^{(i)}\mathbf{w}^{(i)} - \mathbf{y}^{(i)}\|^2 + \lambda \sum_{\{i,j\}} A_{i,j} \|\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\|^2$$

using stacked parameters $\mathbf{w} = \left(\mathbf{w}^{(1)}, \cdots, \mathbf{w}^{(n)}\right)^T$,

$$\min_{\mathbf{w}} \mathbf{w}^T \mathbf{Q} \, \mathbf{w} + \mathbf{w}^T \mathbf{q}$$

with psd matrix $\mathbf{Q}$ and vector $\mathbf{q}$ that depend on local datasets, GTVMin parameter $\lambda$ and empirical graph

# GTV Min. for Local Lin.Reg.

$$\min_{\mathbf{w}} \mathbf{w}^T \mathbf{Q} \mathbf{w} + \mathbf{w}^T \mathbf{q}$$

can be solved using gradient methods

$$\boldsymbol{w}^{(k+1)} = \boldsymbol{w}^{(k)} - \boldsymbol{\alpha}_k \left( 2\mathbf{Q}\boldsymbol{w}^{(k)} + \mathbf{q} \right)$$

# Statistical Aspects

# GTV Min. for Local Lin.Reg.

$$\min_{\mathbf{w}} \sum_{i} \|\mathbf{X}^{(i)}\mathbf{w}^{(i)} - \mathbf{y}^{(i)}\|^2 + \lambda \sum_{\{i,j\}} A_{i,j} \|\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\|^2$$

using stacked parameters $\mathbf{w} = \left(\mathbf{w}^{(1)}, \cdots, \mathbf{w}^{(n)}\right)^T$,

$$\Sigma_{\{i,j\}} A_{i,j} \|\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\|^2 = \mathbf{w}^T (\mathbf{L} \otimes \mathbf{I}) \, \mathbf{w}$$

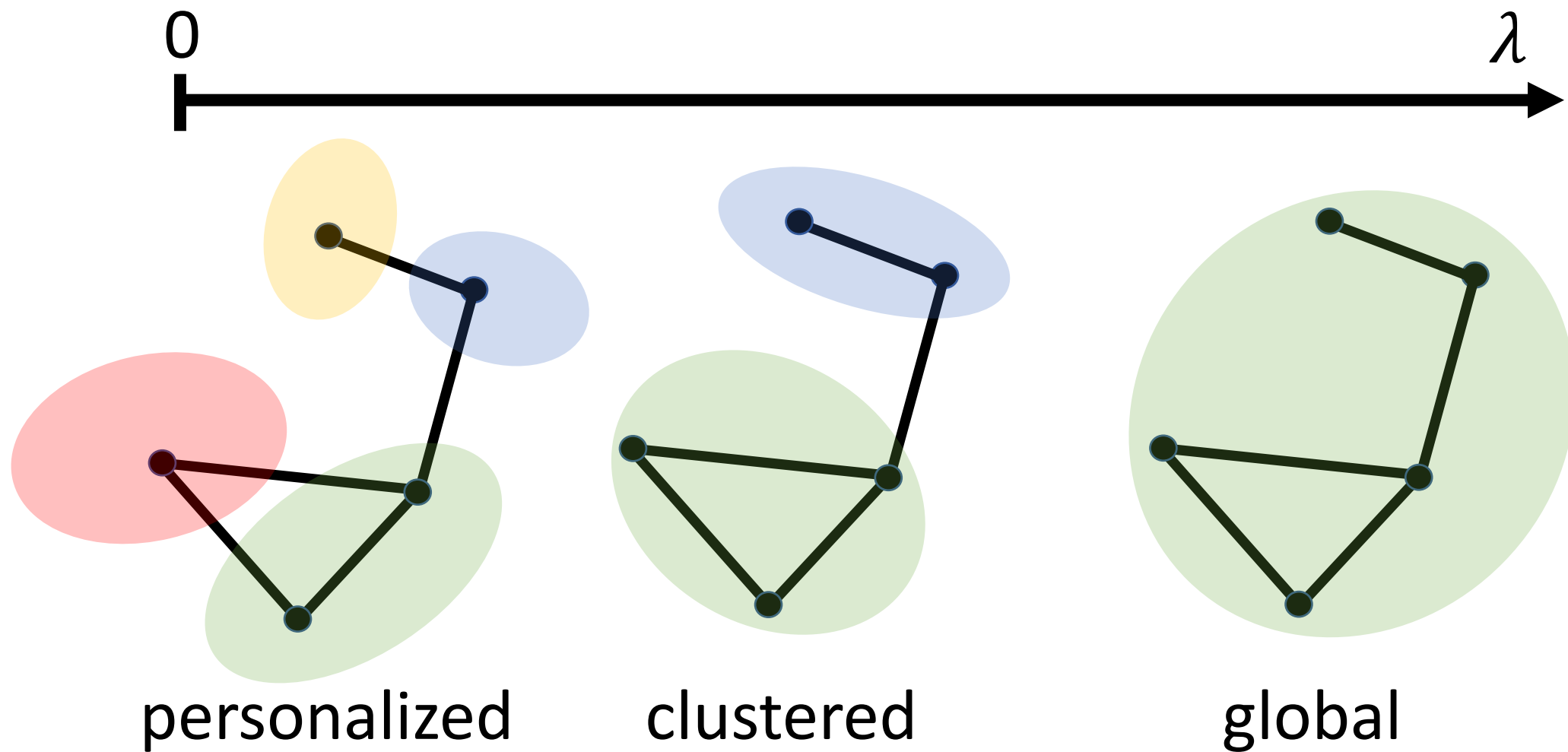with the graph Laplacian $\mathbf{L}$

# Spectral Clustering

for large $\lambda$, GTVMin is to minimize

$$\sum_{\{i,j\}} A_{i,j} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(j)} \right\|^2 = \mathbf{w}^T (\mathbf{L} \otimes \mathbf{I}) \, \mathbf{w}$$

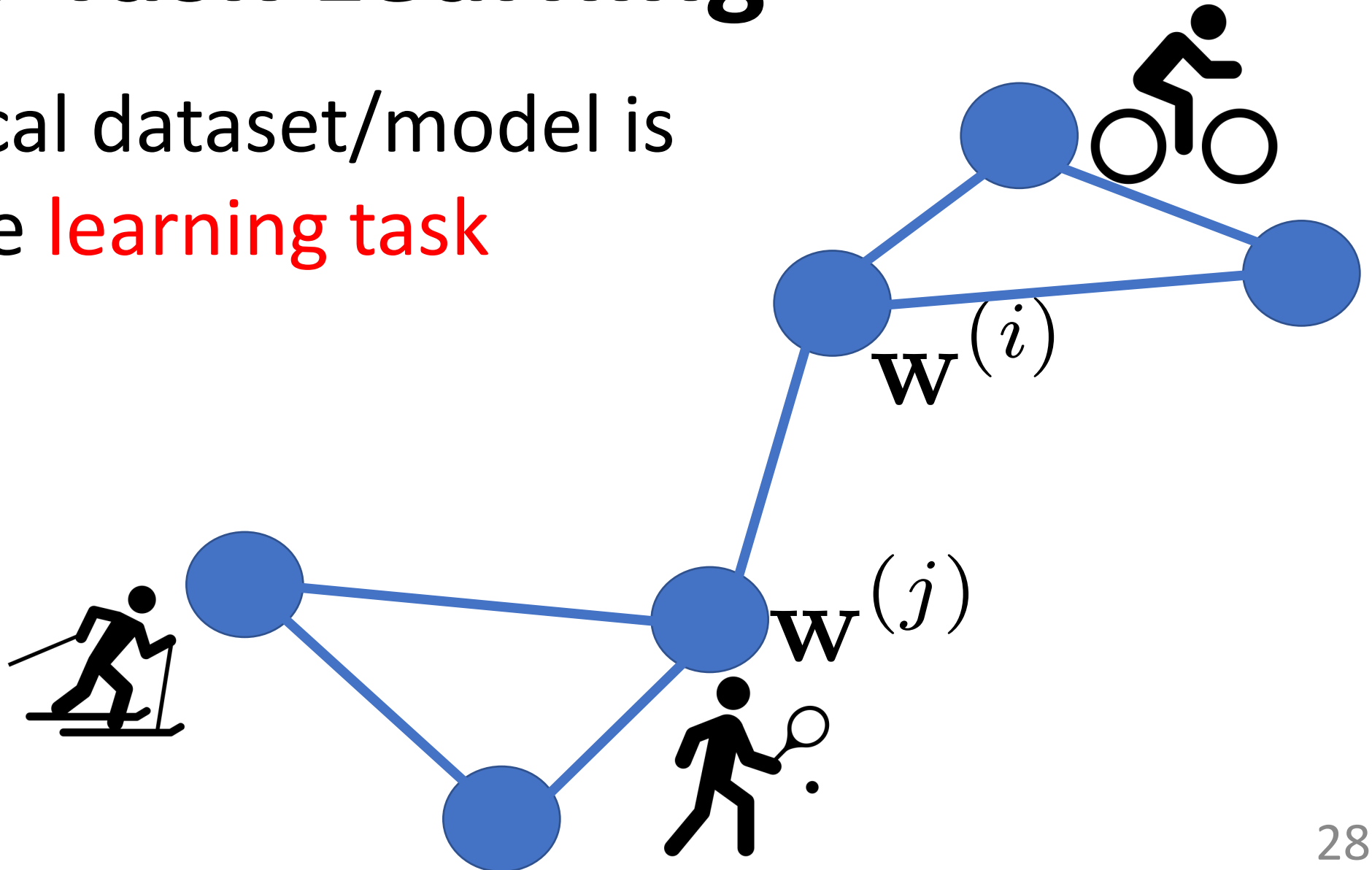$\Rightarrow$ local model parameters composed of eigvecs. of L corresponding to smallest eig.vals
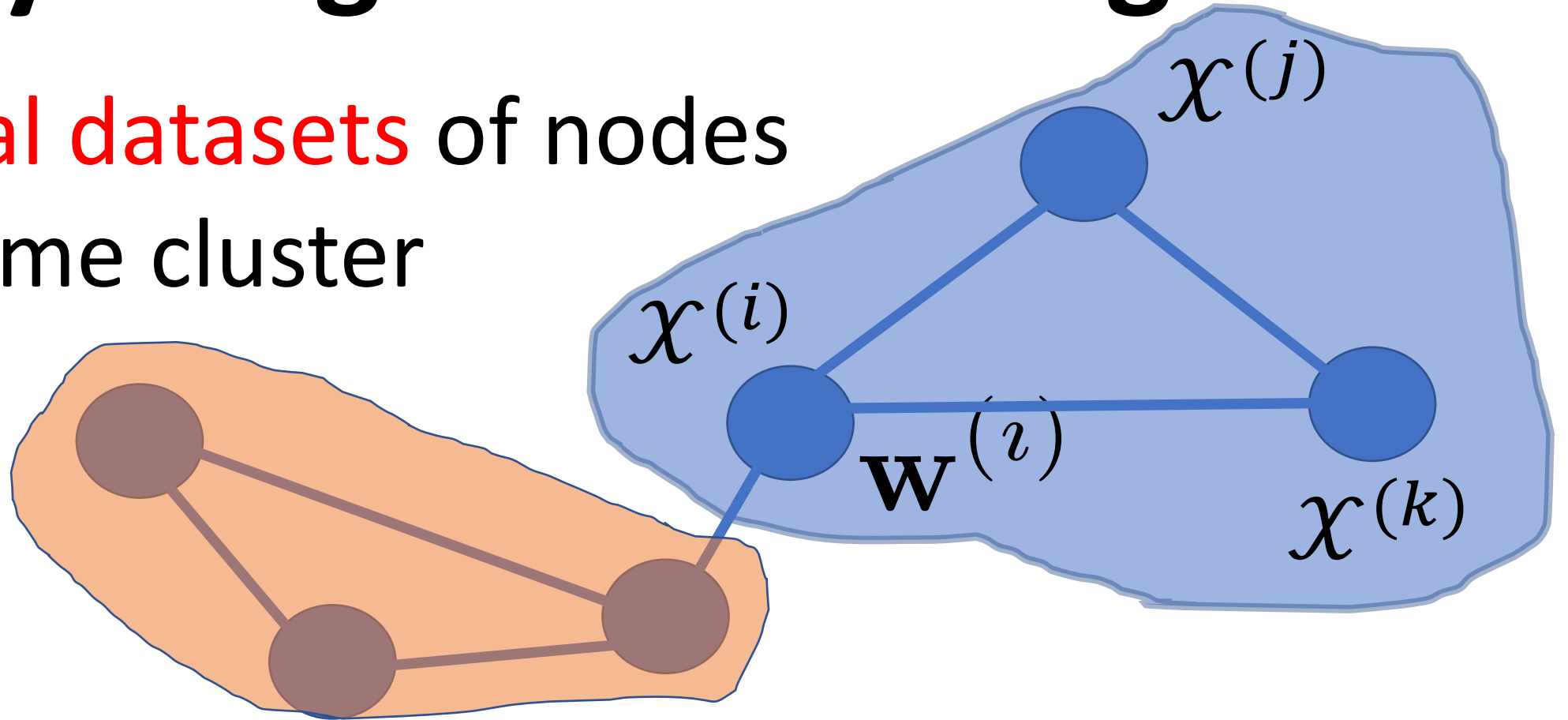
# Clustering of GTVMin Solutions

# Interpretations

# Multi-Task Learning

each local dataset/model is
separate learning task

$\mathbf{w}^{(i)}$

$\mathbf{w}^{(j)}$

# Locally Weighted Learning

<span style="color:red">pool local datasets</span> of nodes
in the same cluster



William S. Cleveland, Susan J. Devlin, Eric Grosse,
"Regression by local fitting: Methods, properties, and computational algorithms,"
Journal of Econometrics, Volume 37, Issue 1, 1988.

# Generalized Convex Clustering

$$\min_{\mathbf{w}} \sum_{i} \left\| w^{(i)} - a^{(i)} \right\|^2 + \lambda \sum_{\{i,j\}} A_{i,j} \left\| w^{(i)} - w^{(j)} \right\|_p$$
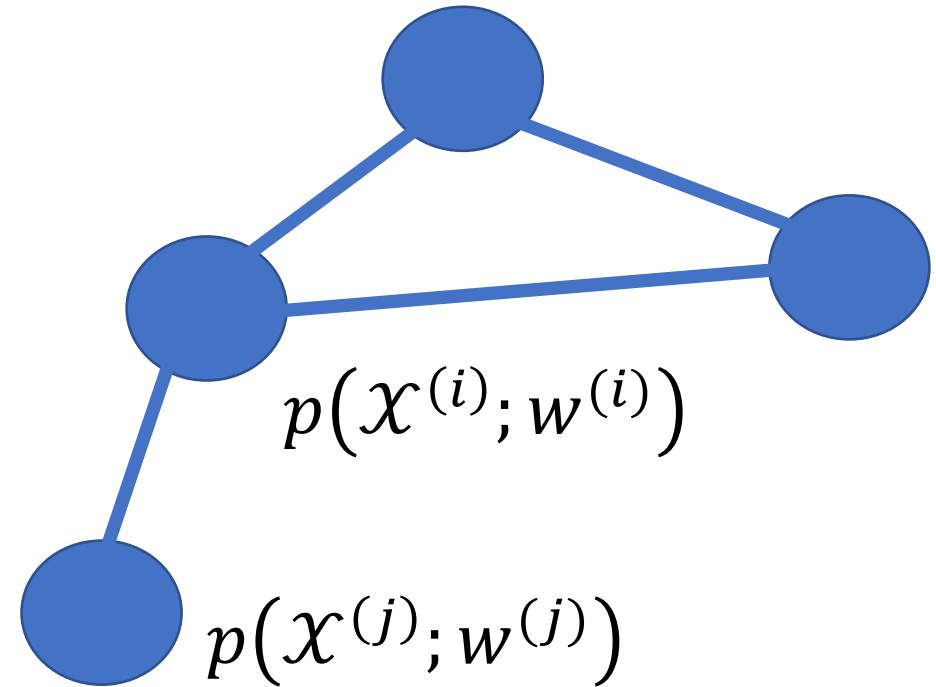
D. Sun, K.-C. Toh, Y. Yuan;
**Convex Clustering: Model, Theoretical Guarantee and Efficient Algorithm**, JMLR, 22(9):1–32, 2021

# (Probabilistic) Graphical Model

separate prob. space for each local dataset



$$p(\mathcal{X}^{(i)}; w^{(i)})$$

traditionally, PGMs use a common prob. space for all local datasets

$$p(\mathcal{X}^{(j)}; w^{(j)})$$

AJ, **"Networked Exponential Families for Big Data Over Networks,"** in *IEEE Access*, vol. 8, pp. 202897-202909, 2020, doi: 10.1109/ACCESS.2020.3033817.

# Approx. Hierarch. Bayes' Model



$$p(\mathbf{w}) \xrightarrow{\text{i.i.d. draw}} \mathbf{w}^{(i)}$$

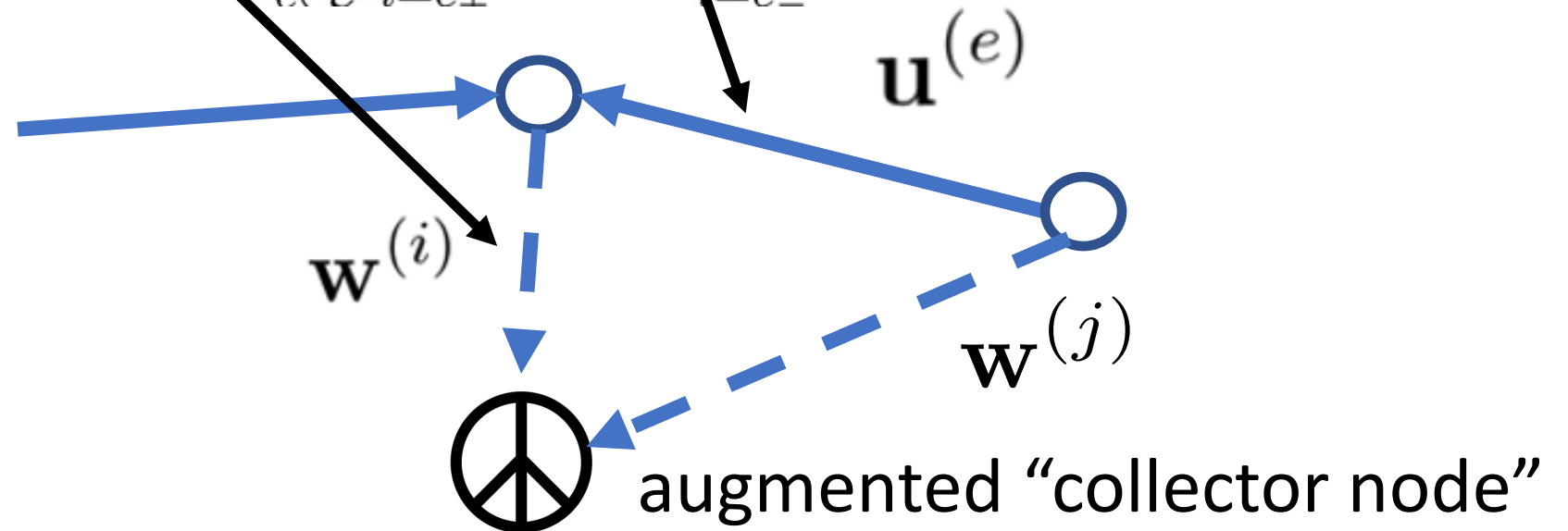i.i.d. draw

$$\mathbf{w}^{(j)}$$

Lyu, B., Hanzely, F., and Kolar, M., "Personalized Federated Learning with Multiple Known Clusters", *arXiv e-prints*, 2022. doi:10.48550/arXiv.2204.13619.
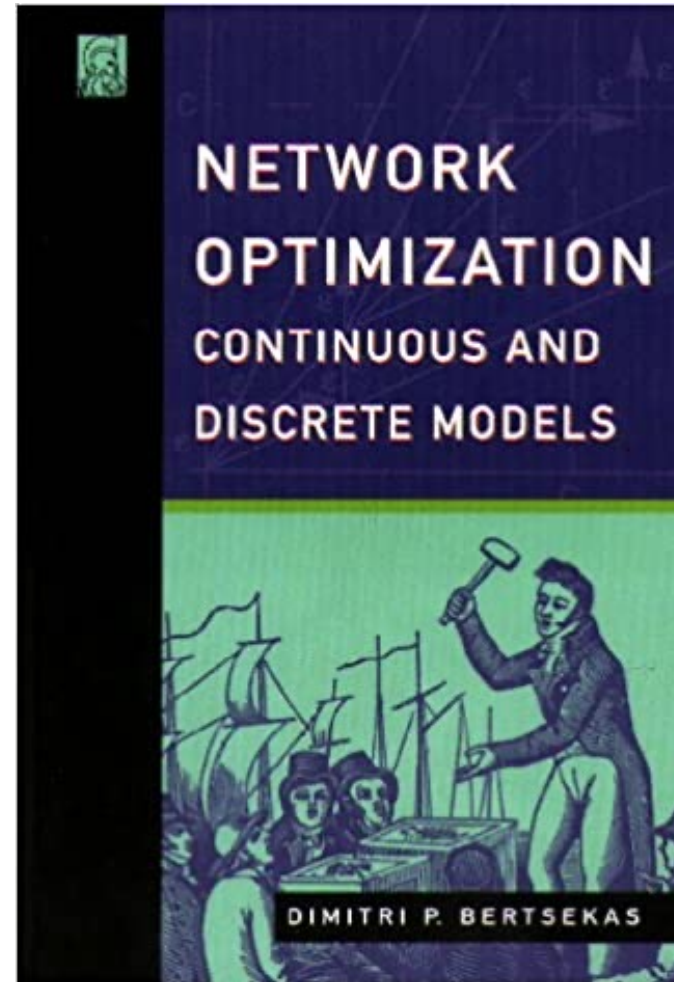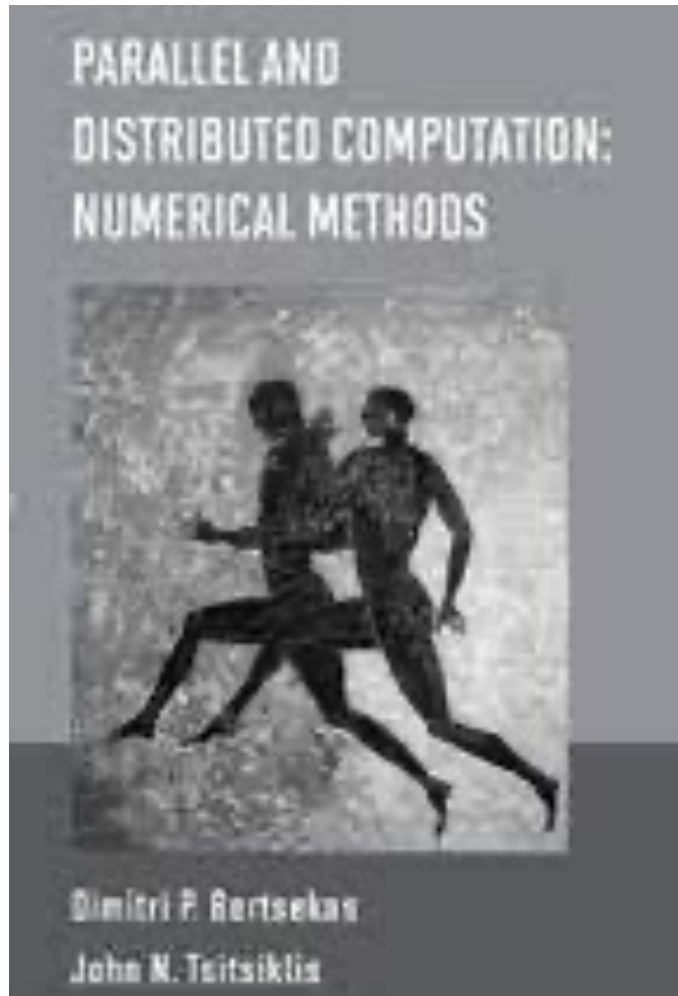
# Non-Linear Min-Cost-Flow

$$\max_{\mathbf{u} \in \mathcal{U}} - \sum_{i \in \mathcal{V}} L_i^* \left( \mathbf{w}^{(i)} \right) - \lambda \sum_{e \in \mathcal{E}} A_e \phi^* \left( \mathbf{u}^{(e)} / (\lambda A_e) \right)$$

subject to $-\mathbf{w}^{(i)} = \sum_{e \in \mathcal{E}} \sum_{i=e_+} \mathbf{u}^{(e)} - \sum_{i=e_-} \mathbf{u}^{(e)}$ for all nodes $i \in \mathcal{V}$.

$\mathbf{u}^{(e)}$

$\mathbf{w}^{(i)}$

$\mathbf{w}^{(j)}$

augmented "collector node"

# Non-Linear Min-Cost-Flow

# Electrical Network.
# ("AI is new Electricity!")

**Kirchhoff's Current Law**

$$\sum_{e \in \mathcal{E}} \sum_{i=e_+} \widehat{\mathbf{u}}^{(e)} - \sum_{i=e_-} \widehat{\mathbf{u}}^{(e)} = -\nabla L_i\left(\widehat{\mathbf{w}}^{(i)}\right) \text{ for all nodes } i \in \mathcal{V}$$
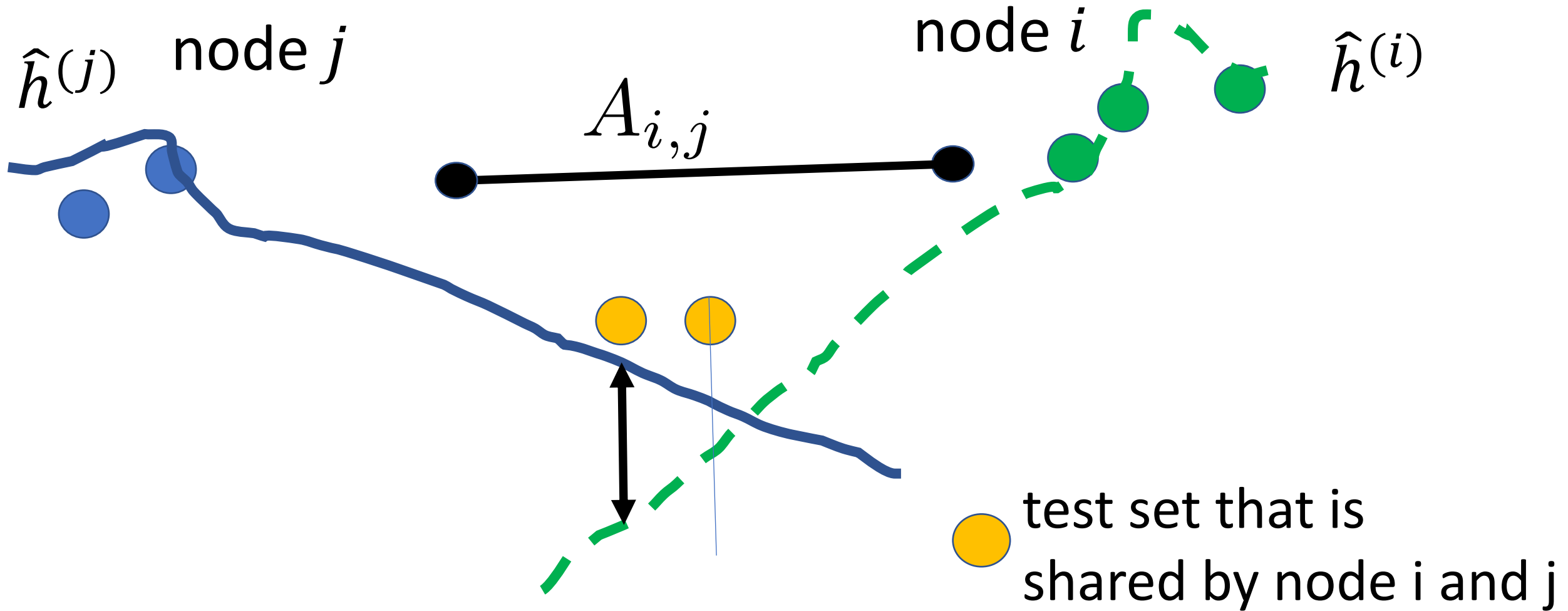
$$\widehat{\mathbf{w}}^{(e_+)} - \widehat{\mathbf{w}}^{(e_-)} \in (\lambda A_e)\partial\phi^*(\widehat{\mathbf{u}}^{(e)}/(\lambda A_e)) \text{ for every edge } e \in \mathcal{E}.$$

**Generalized Ohm Law**

# GTVMin for Non-Param. Models

Variation of Non-Param. Models

$\hat{h}^{(j)}$  node $j$  node $i$  $\hat{h}^{(i)}$

$A_{i,j}$

test set that is shared by node i and j

# Wrap Up.

- couple local model training via regularization

- regularizer obtained via GTV (over empirical graph)

- FL algorithms = optimization methods for GTV min

- GTVmin pools local datasets into clusters

- cluster structure depends on emp.graph <span style="color:red">and</span> local data!

# Thank you for your attention!