# Network Flows in Federated Learning

Alexander Jung (Aalto University)

[linkedin.com/in/aljung](linkedin.com/in/aljung)

@alexjung111

# guiding theme:

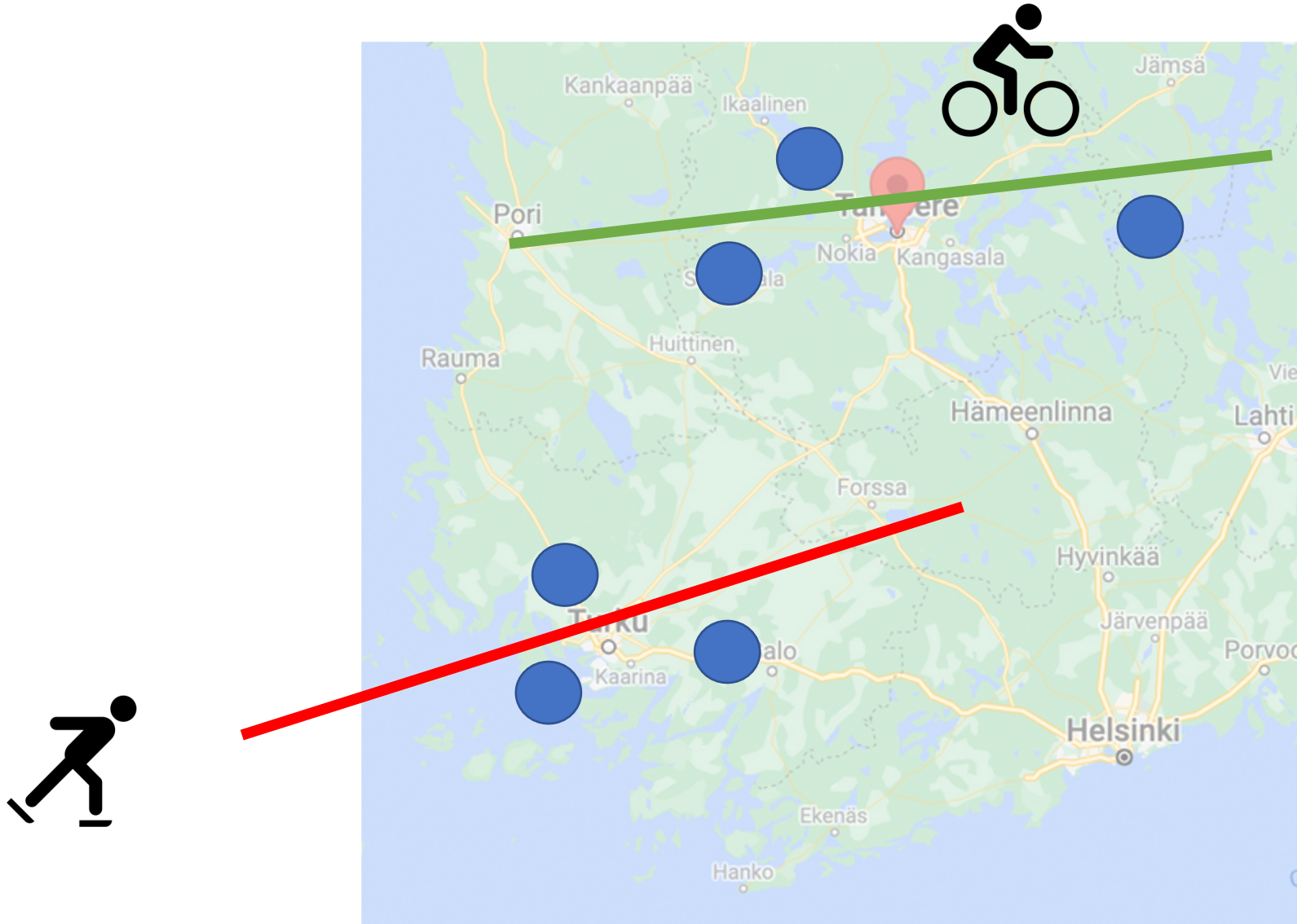organize <span style="color:red">data</span>, <span style="color:red">models</span> and <span style="color:red">computation</span> for
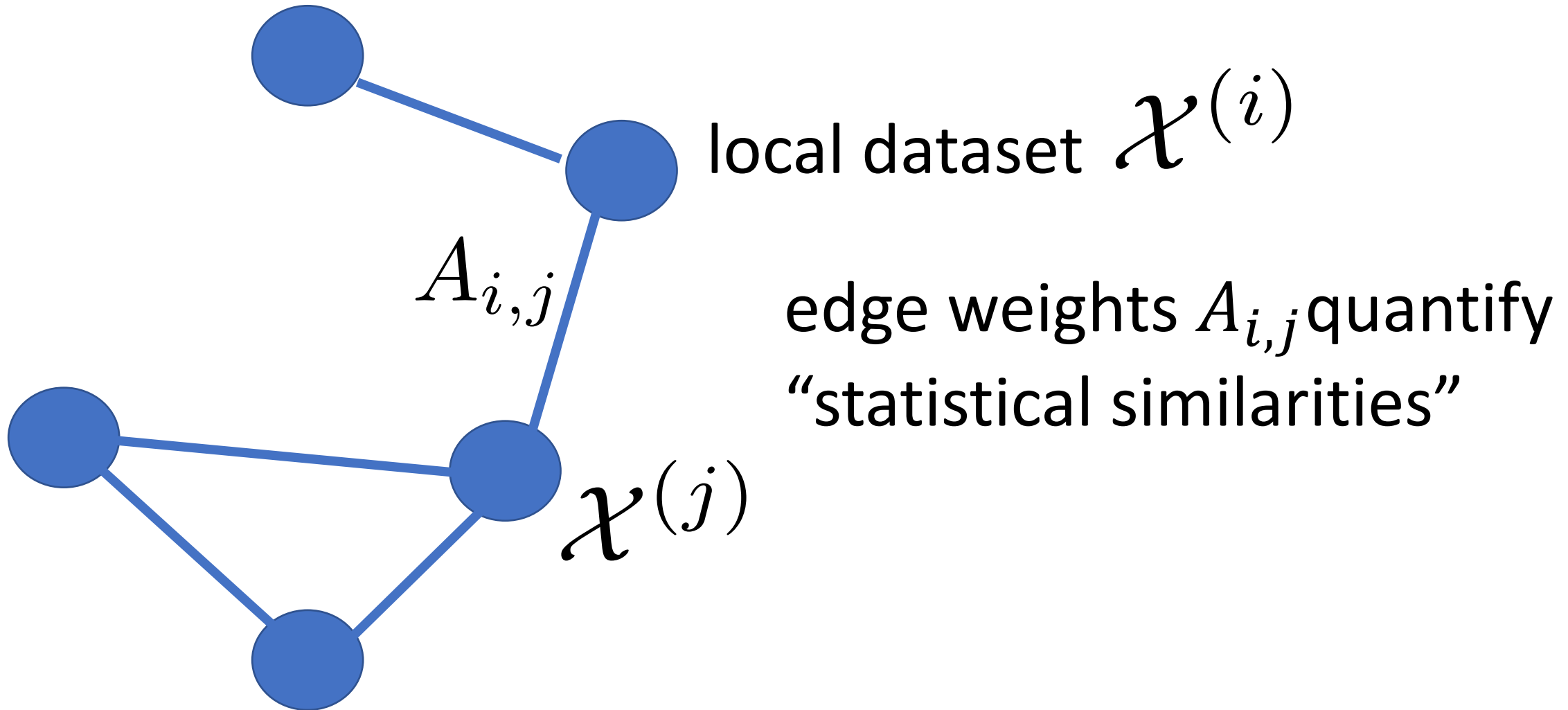
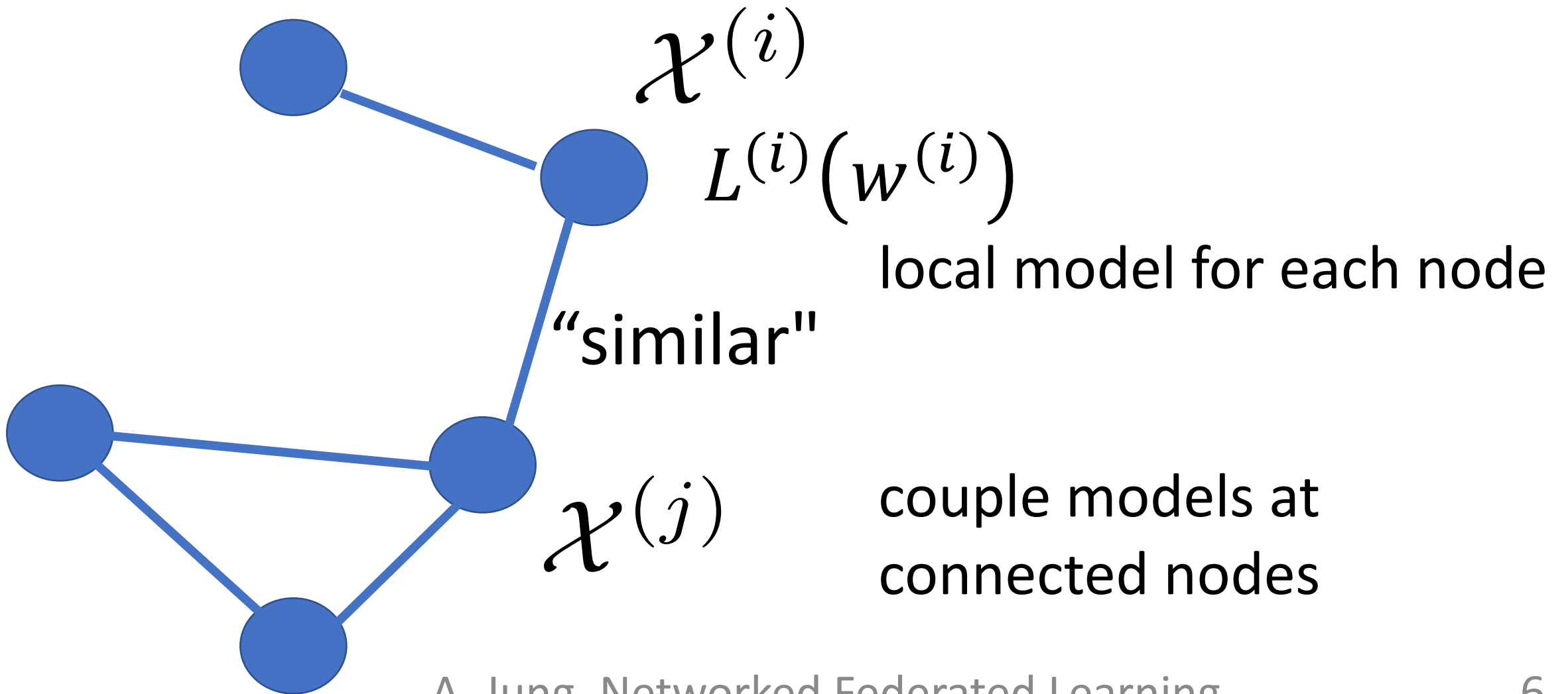machine learning as <span style="color:red">networks</span>.

# Weather Stations.

# Personalized Weather Forecast

# The Empirical Graph



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

$A_{i,j}$

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

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

# Networked Models.



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

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

local model for each node

"similar"

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

couple models at connected nodes

# TV Minimization

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

## Network Lasso: Clustering and Optimization in Large Graphs

by D Hallac · 2015 · Cited by 206 — **Network Lasso: Clustering and Optimization in Large Graphs** ... Keywords: Convex **Optimization**, ADMM, **Network Lasso**. Go to: ... 2013 [**Google Scholar**]. 2.

Abstract · INTRODUCTION · CONVEX PROBLEM... · EXPERIMENTS

# Rewrite GTVMin

$$\widehat{\mathbf{w}} \in \arg\min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}) + g(\mathbf{Dw})$$
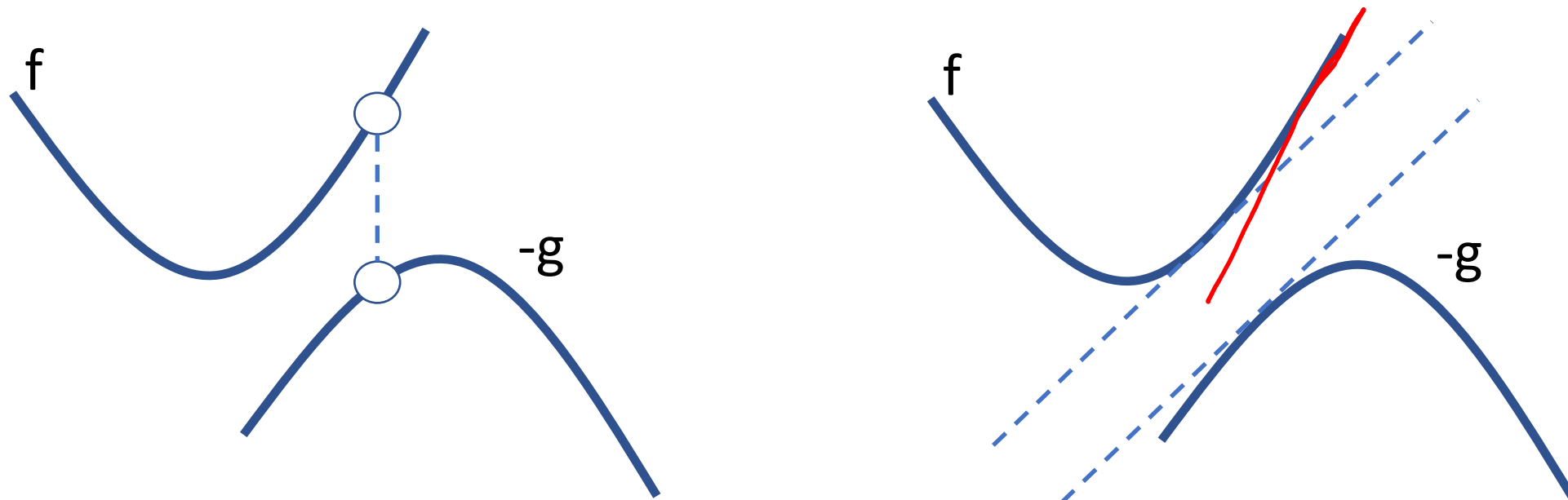
$$\text{with } f(\mathbf{w}) := \sum_{i \in \mathcal{V}} L_i\left(\mathbf{w}^{(i)}\right), \text{ and } g(\mathbf{u}) := \lambda \sum_{e \in \mathcal{E}} A_e \phi\left(\mathbf{u}^{(e)}\right).$$

with incidence matrix/operator

$$\mathbf{D} : \mathcal{W} \to \mathcal{U} : \mathbf{w} \mapsto \mathbf{u} \text{ with } \mathbf{u}^{(e)} = \mathbf{w}^{(e+)} - \mathbf{w}^{(e-)}.$$

# Fenchel's Duality.

$$\min_{\mathbf{w}\in\mathcal{W}} f(\mathbf{w}) + g(\mathbf{Dw}) = \max_{\mathbf{u}\in\mathcal{U}} -g^*(\mathbf{u}) - f^*(-\mathbf{D}^T\mathbf{u}).$$
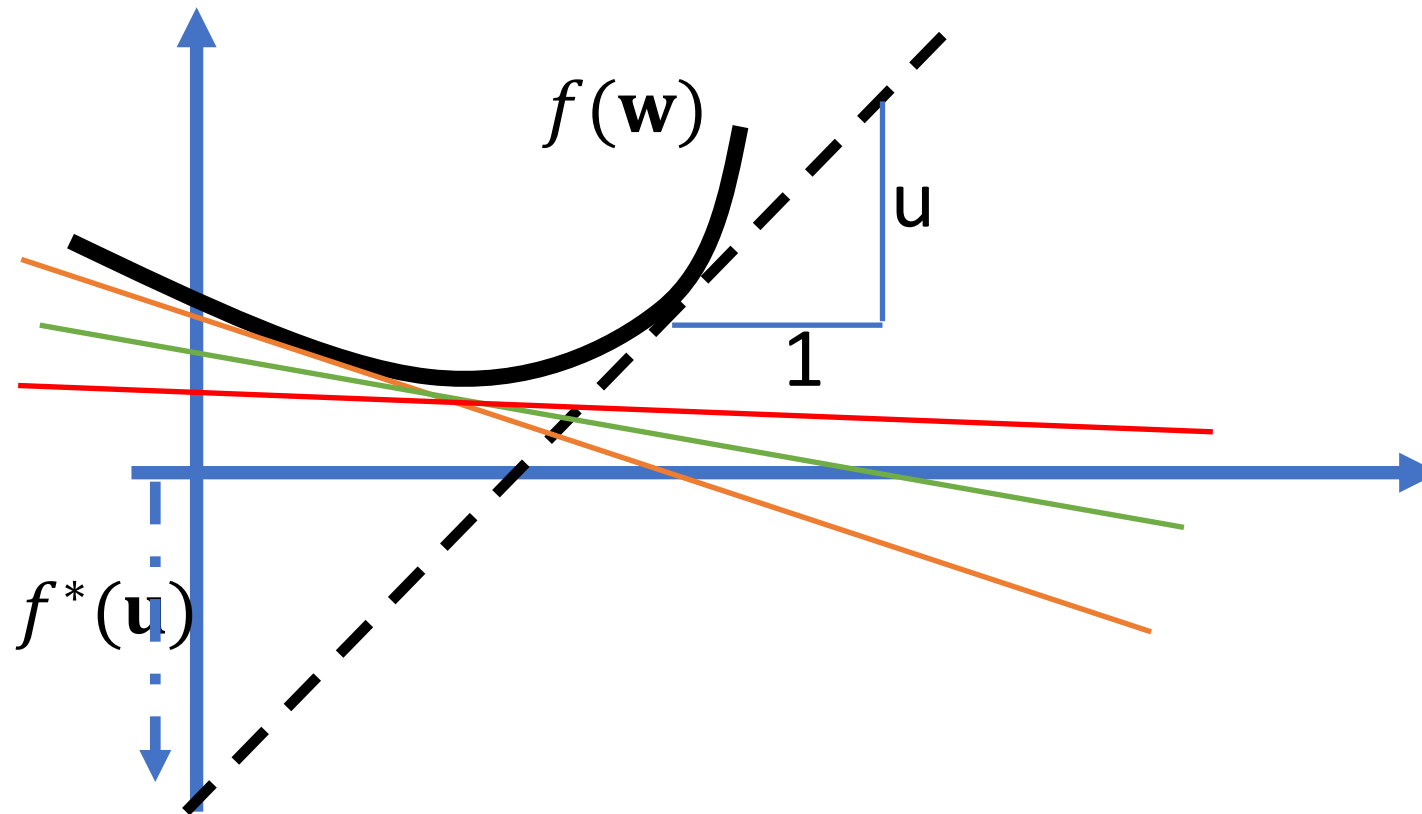


R. T. Rockafellar, *Convex Analysis*.    Princeton, NJ: Princeton Univ. Press, 1970.
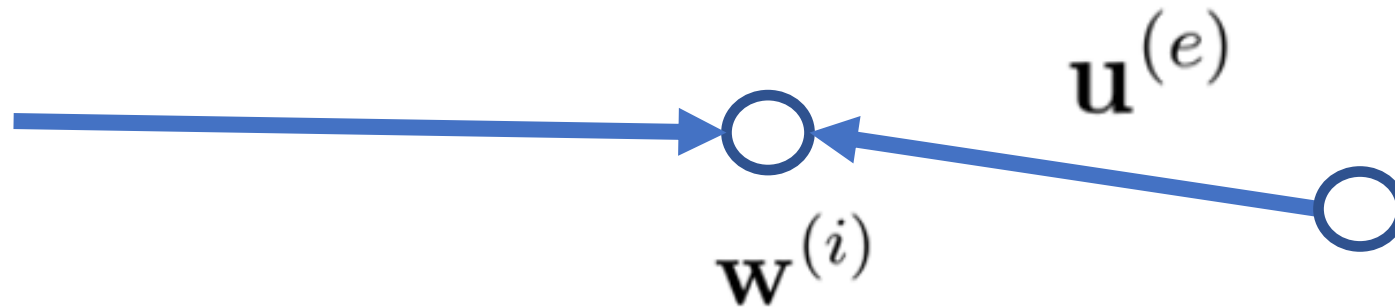https://en.wikipedia.org/wiki/Fenchel%27s_duality_theorem

# Convex Conjugate.

$$f^*(\mathbf{w}) := \sup_{\mathbf{z} \in \mathbb{R}^{n|\mathcal{V}|}} \mathbf{w}^T \mathbf{z} - f(\mathbf{z}) \qquad g^*(\mathbf{u}) := \sup_{\mathbf{z} \in \mathbb{R}^{n|\mathcal{E}|}} \mathbf{u}^T \mathbf{z} - g(\mathbf{z})$$
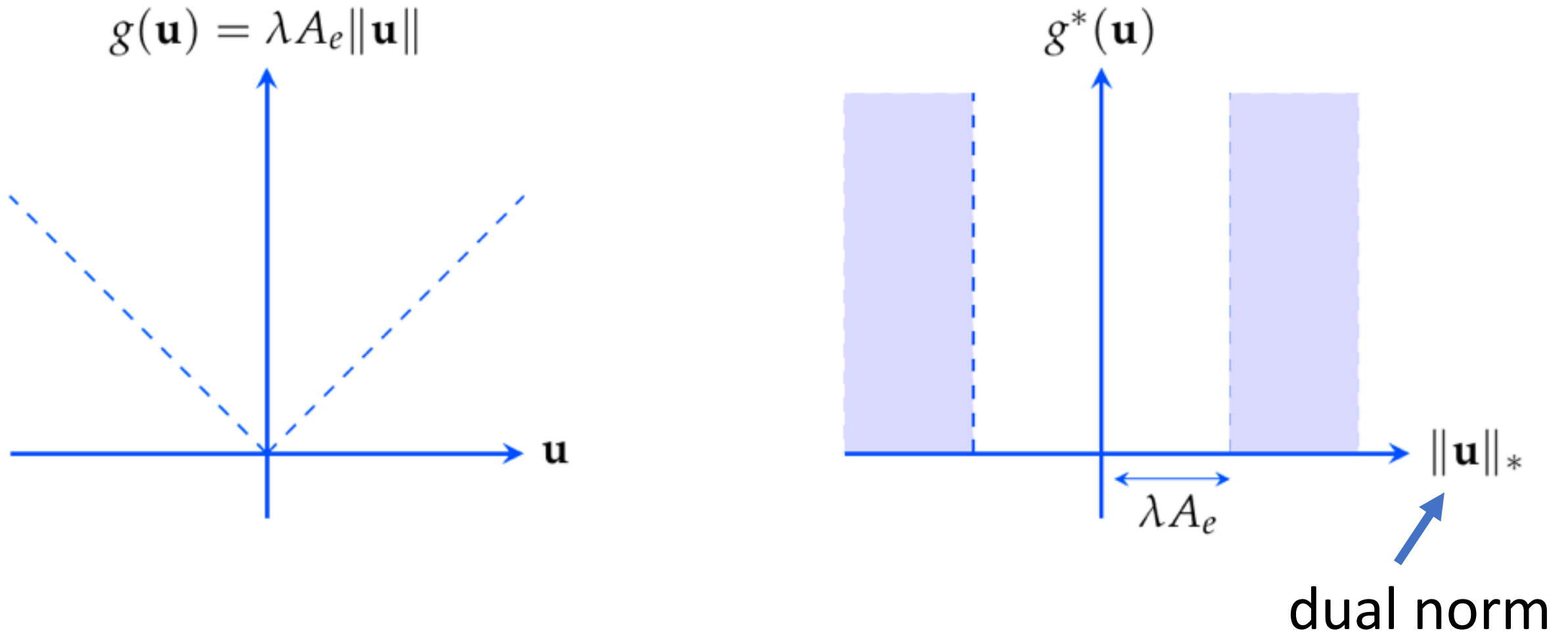
# The Dual of GTVMin.

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

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



dual variables $\mathbf{u}^{(e)}$ for each (oriented) edge $e = (j, i)$
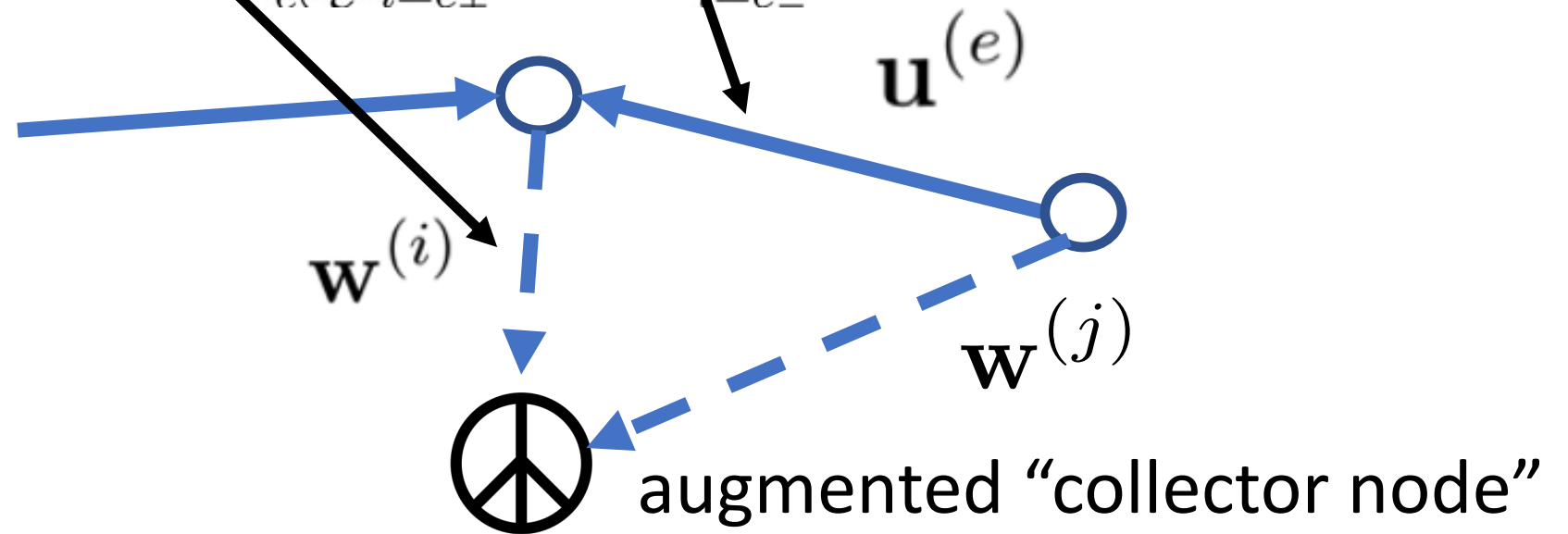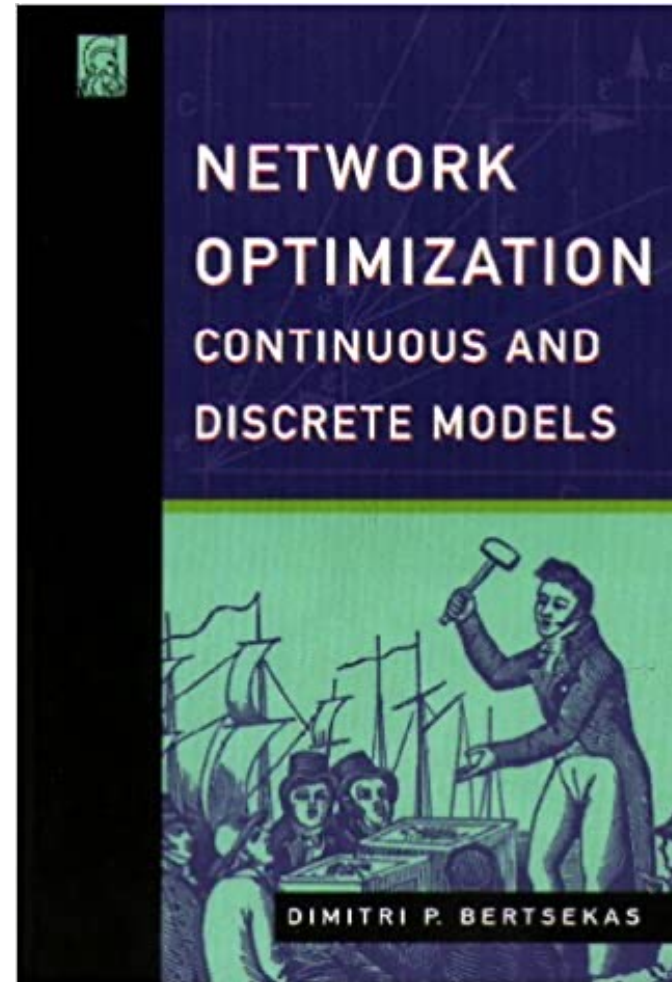
# Convex Conjugate of Norm

$$g(\mathbf{u}) = \lambda A_e \|\mathbf{u}\|$$

$$g^*(\mathbf{u})$$

$\mathbf{u}$

$\lambda A_e$

$\|\mathbf{u}\|_*$

dual norm

# Non-Linear Min-Cost-Flow

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

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

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

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

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

augmented "collector node"

# Non-Linear Min-Cost-Flow

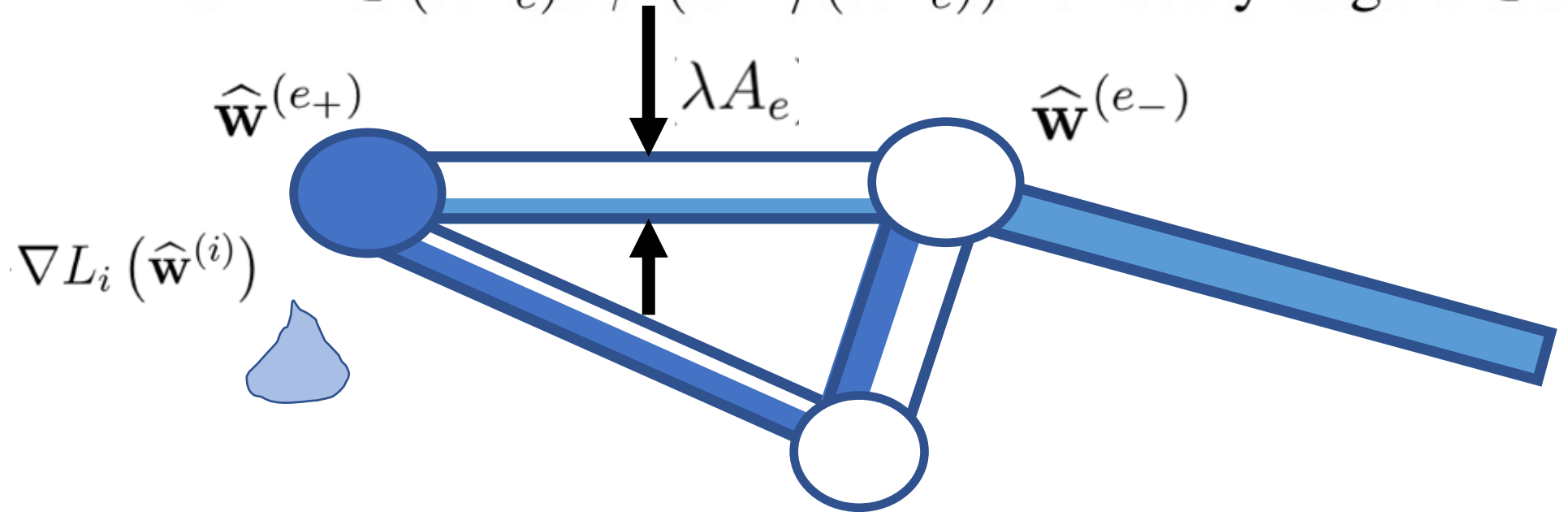# Primal and Dual Optimality.

$$\sum_{e \in \mathcal{E}} \sum_{i=e_+} \widehat{\mathbf{u}}^{(e)} - \sum_{i=e_-} \widehat{\mathbf{u}}^{(e)} = -\nabla L_i \left( \widehat{\mathbf{w}}^{(i)} \right) \quad \text{for all nodes } i \in \mathcal{V}$$

$$\widehat{\mathbf{w}}^{(e_+)} - \widehat{\mathbf{w}}^{(e_-)} \in (\lambda A_e) \partial \phi^* (\widehat{\mathbf{u}}^{(e)} / (\lambda A_e)) \quad \text{for every edge } e \in \mathcal{E}.$$



$\widehat{\mathbf{w}}^{(e_+)}$

$\lambda A_e$

$\widehat{\mathbf{w}}^{(e_-)}$

$\nabla L_i \left( \widehat{\mathbf{w}}^{(i)} \right)$

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

**Kirchhoff's Current Law**

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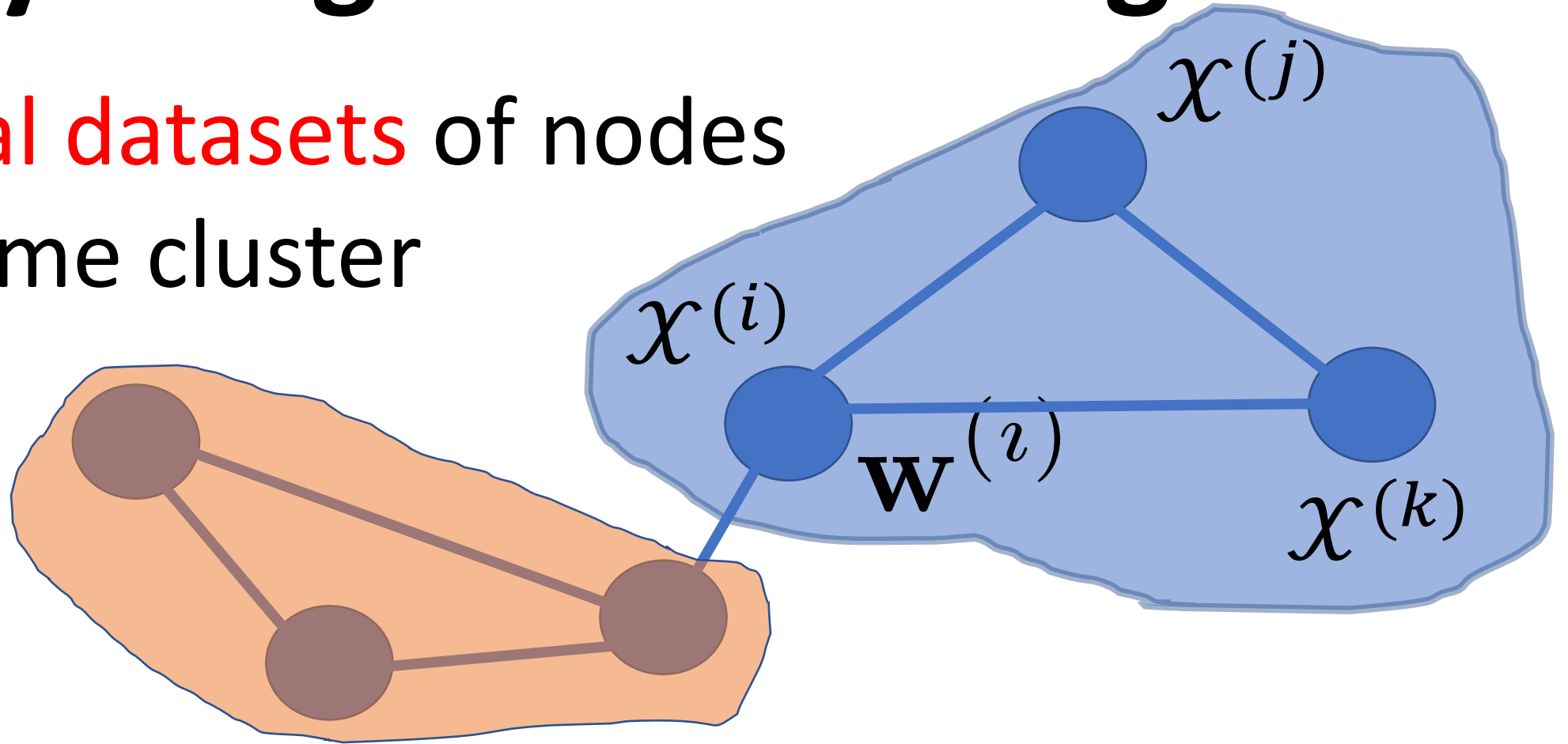
$$\widehat{\mathbf{w}}^{(e_+)} - \widehat{\mathbf{w}}^{(e_-)} \in (\lambda A_e)\partial\phi^*(\widehat{\mathbf{u}}^{(e)}/(\lambda A_e)) \text{ for every edge } e \in \mathcal{E}.$$

**Generalized Ohm Law**

# Locally Weighted Learning

pool local datasets of nodes in the same cluster



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

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

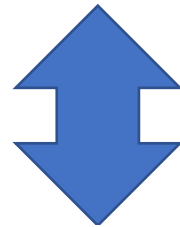$\mathbf{w}^{(i)}$

$\mathcal{X}^{(k)}$

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.

# Primal-Dual Optimality Conditions.

(assuming convexity of loss functions and GTV penalty)

$$\mathbf{M}^{-1} \begin{pmatrix} \partial f & \mathbf{D}^T \\ -\mathbf{D} & \partial g^* \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{w}} \\ \widehat{\mathbf{u}} \end{pmatrix} \ni \mathbf{0} \text{ with } \mathbf{M} := \begin{pmatrix} \mathbf{T}^{-1} & -\mathbf{D}^T \\ -\mathbf{D} & \mathbf{\Sigma}^{-1} \end{pmatrix}$$

$$\begin{pmatrix} \widehat{\mathbf{w}} \\ \widehat{\mathbf{u}} \end{pmatrix} = \left( \mathbf{I} + \mathbf{M}^{-1} \begin{pmatrix} \partial f & \mathbf{D}^T \\ -\mathbf{D} & \partial g^* \end{pmatrix} \right)^{-1} \begin{pmatrix} \widehat{\mathbf{w}} \\ \widehat{\mathbf{u}} \end{pmatrix}$$

## this is again a fixed-point problem !

# Proximal Point Algorithm.

primal and dual variables $\widehat{w}, \widehat{u}$ optimal if and only if

$$\mathbf{M}^{-1} \begin{pmatrix} \partial f & \mathbf{D}^T \\ -\mathbf{D} & \partial g^* \end{pmatrix} \begin{pmatrix} \widehat{\mathbf{w}} \\ \widehat{\mathbf{u}} \end{pmatrix} \ni \mathbf{0} \text{ with } \mathbf{M} := \begin{pmatrix} \mathbf{T}^{-1} & -\mathbf{D}^T \\ -\mathbf{D} & \mathbf{\Sigma}^{-1} \end{pmatrix}$$

solve iteratively by <span style="color:red">proximal point algorithm</span>

$$\begin{pmatrix} \widehat{\mathbf{w}}^{(k+1)} \\ \widehat{\mathbf{u}}^{(k+1)} \end{pmatrix} = \left( \mathbf{I} + \mathbf{M}^{-1} \begin{pmatrix} \partial f & \mathbf{D}^T \\ -\mathbf{D} & \partial g^* \end{pmatrix} \right)^{-1} \begin{pmatrix} \widehat{\mathbf{w}}^{(k)} \\ \widehat{\mathbf{u}}^{(k)} \end{pmatrix}$$

A. Chambolle, T. Pock. An introduction to continuous optimization for imaging. Acta Numerica, 2016.
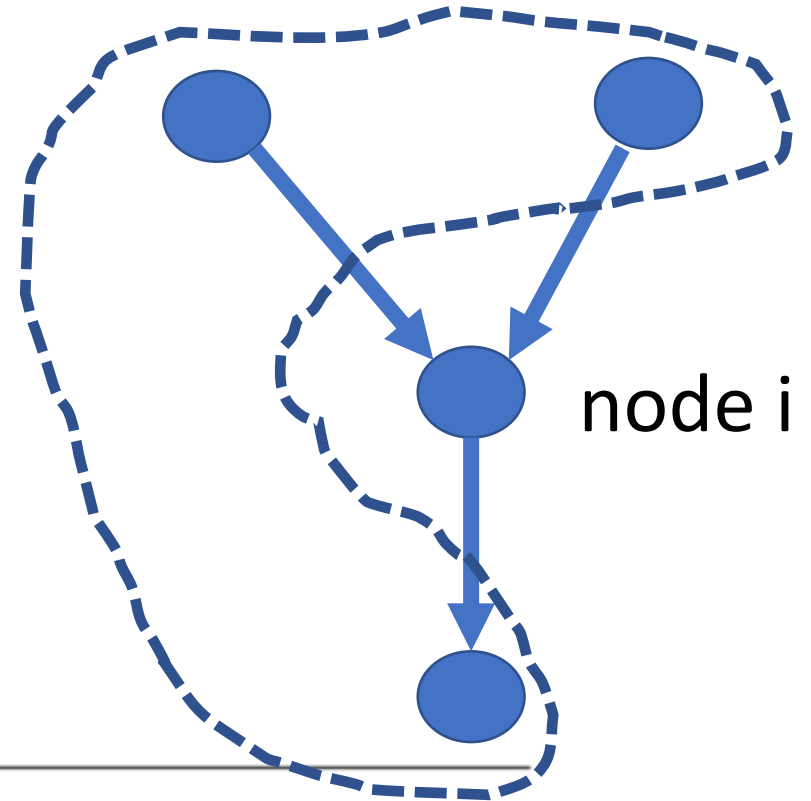
# After Some Manipulations.

**Algorithm 1** Primal-Dual Method for Networked FL

**Input**: empirical graph $\mathcal{G} = (\mathcal{V}, \mathcal{E}, \mathbf{A})$; training set $\{\mathbf{X}^{(i)}\}_{i \in \mathcal{M}}$; regularization parameter $\lambda$; loss $\mathcal{L}$; GTV penalty $\phi$

**Initialize**: $k := 0; \widehat{\mathbf{w}}_0 := \mathbf{0}; \widehat{\mathbf{u}}_0 := \mathbf{0}; \sigma_e = 1/2$ and $\tau_i = 1/|\mathcal{N}_i|$
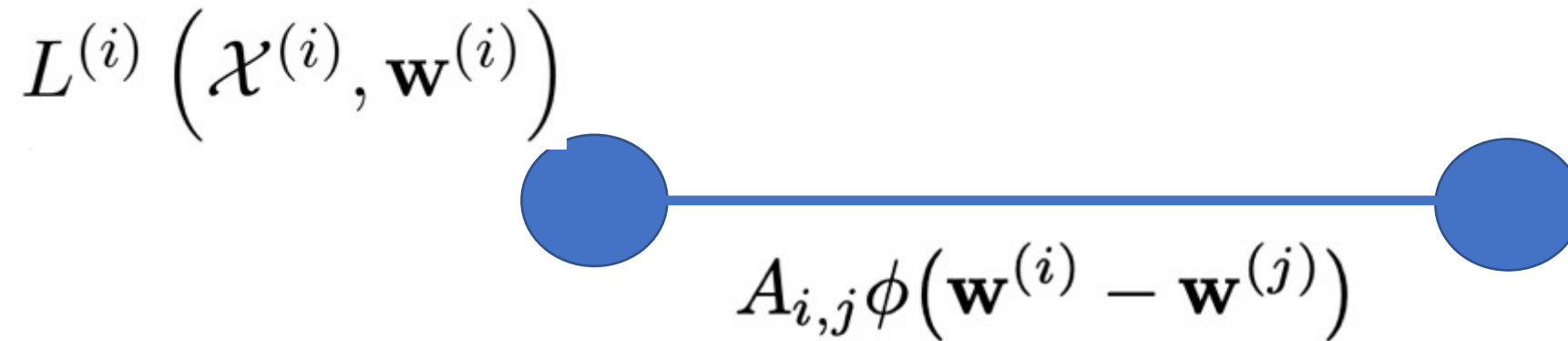
1:   **while** stopping criterion is not satisfied **do**
2:       **for** all nodes $i \in \mathcal{V}$ **do**
3:           $\widehat{\mathbf{w}}_{k+1}^{(i)} := \widehat{\mathbf{w}}_k^{(i)} - \tau_i \sum_{e \in \mathcal{E}} D_{e,i} \widehat{\mathbf{u}}_k^{(e)}$
4:       **end for**
5:       **for** nodes in the training set $i \in \mathcal{M}$ **do**
6:           $\widehat{\mathbf{w}}_{k+1}^{(i)} := \mathcal{P}\mathcal{U}^{(i)}\{\widehat{\mathbf{w}}_{k+1}^{(i)}\}$
7:       **end for**
8:       **for** all edges $e \in \mathcal{E}$ **do**
9:           $\widehat{\mathbf{u}}_{k+1}^{(e)} := \widehat{\mathbf{u}}_k^{(e)} + \sigma_e\big(2\big(\widehat{\mathbf{w}}_{k+1}^{(e_+)} - \widehat{\mathbf{w}}_{k+1}^{(e_-)}\big) - \big(\widehat{\mathbf{w}}_k^{(e_+)} - \widehat{\mathbf{w}}_k^{(e_-)}\big)\big)$
10:          $\widehat{\mathbf{u}}_{k+1}^{(e)} := \mathcal{D}\mathcal{U}^{(e)}\{\widehat{\mathbf{u}}_{k+1}^{(e)}\}$
11:      **end for**
12:      $k := k+1$
13: **end while**

node i

# Algorithm 1 is Attractive for NFL...

➢ decentralized implementation (mess. pass.)

➢ robust against various imperfections

  ➢ approximate primal/dual updates
  ➢ node/link failures

➢ privacy friendly; no raw data exchanged

# Local Computations in Algorithm 1.

$$L^{(i)}\left(\mathcal{X}^{(i)}, \mathbf{w}^{(i)}\right)$$



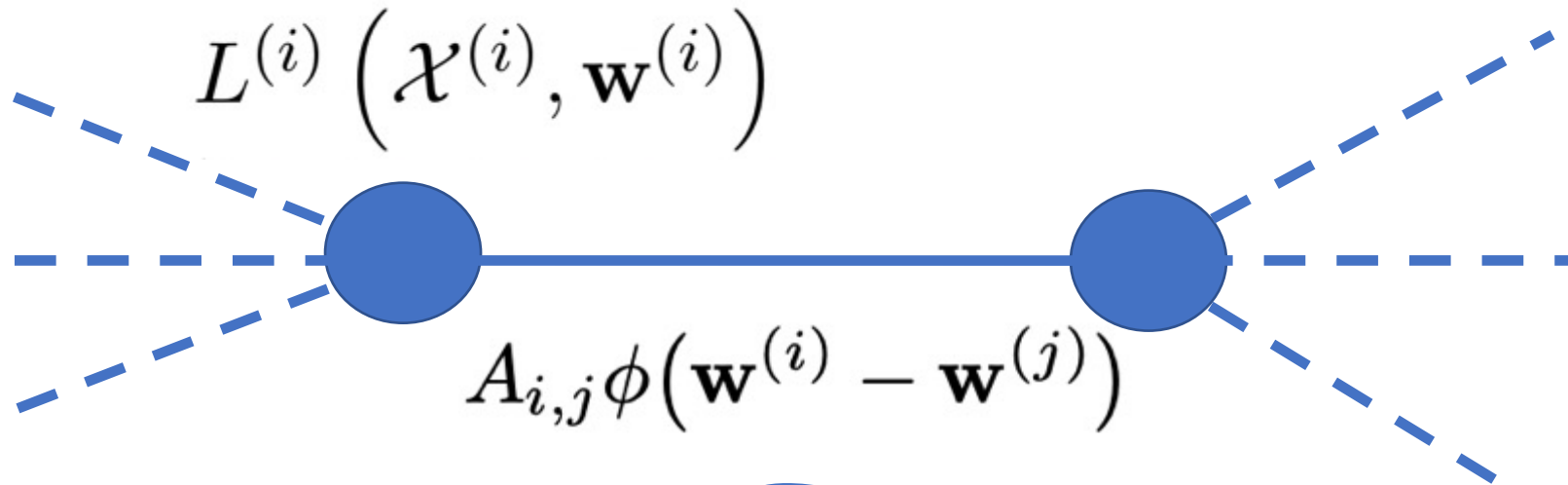$$A_{i,j}\phi\big(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\big)$$

node-wise primal update:

$$\mathcal{PU}^{(i)}\{\mathbf{v}\} := \operatorname*{argmin}_{\mathbf{z}\in\mathbb{R}^n} L^{(i)}(\mathbf{z}) + (1/2\tau_i)\|\mathbf{v} - \mathbf{z}\|^2.$$

edge-wise dual update:

$$\mathcal{DU}^{(e)}\{\mathbf{v}\} := \operatorname*{argmin}_{\mathbf{z}\in\mathbb{R}^n} \lambda A_e \phi^*\big(\mathbf{z}/(\lambda A_e)\big) + (1/2\sigma_e)\|\mathbf{v} - \mathbf{z}\|^2.$$

# Spreading Local Results.

$$L^{(i)}\left(\mathcal{X}^{(i)}, \mathbf{w}^{(i)}\right)$$

$$A_{i,j}\phi\left(\mathbf{w}^{(i)} - \mathbf{w}^{(j)}\right)$$

2:      **for** all nodes $i \in \mathcal{V}$ **do**

3:          $\widehat{\mathbf{w}}^{(i)}_{k+1} := \widehat{\mathbf{w}}^{(i)}_k - \tau_i \sum_{e \in \mathcal{E}} D_{e,i} \widehat{\mathbf{u}}^{(e)}_k$

4:      **end for**

8:      **for** all edges $e \in \mathcal{E}$ **do**

9:          $\widehat{\mathbf{u}}^{(e)}_{k+1} := \widehat{\mathbf{u}}^{(e)}_k + \sigma_e\left(2\left(\widehat{\mathbf{w}}^{(e_+)}_{k+1} - \widehat{\mathbf{w}}^{(e_-)}_{k+1}\right) - \left(\widehat{\mathbf{w}}^{(e_+)}_k - \widehat{\mathbf{w}}^{(e_-)}_k\right)\right)$

# Are GTVMin Solutions Any Good?

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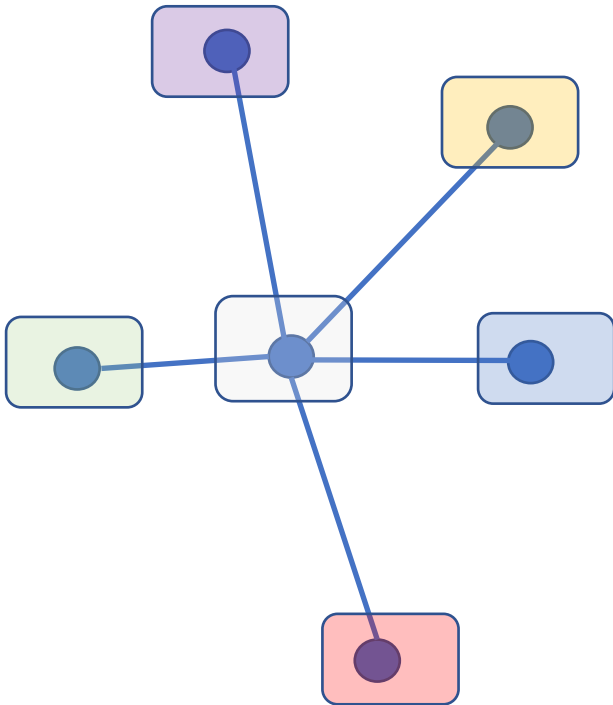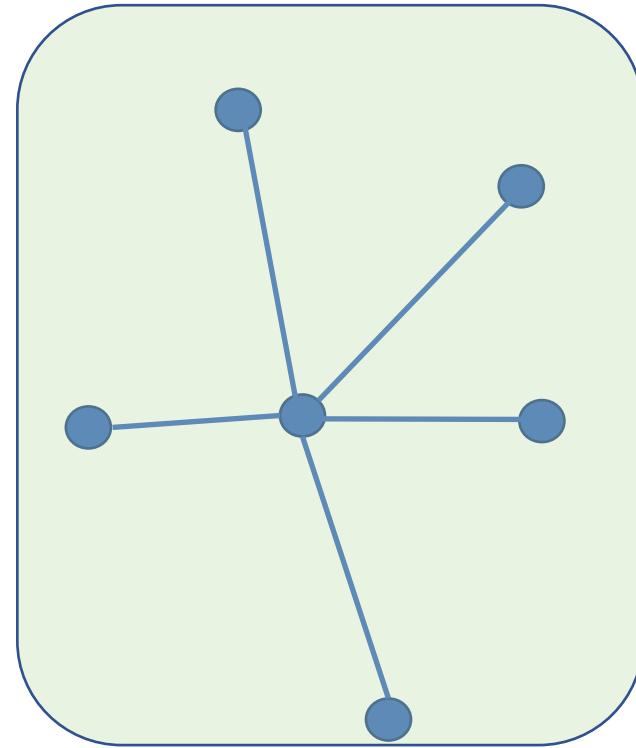
# Clustering



$\mathbf{w}^{(i)}$    $\lambda A_e$    $\mathbf{w}^{(j)}$

$\nabla L^{(i)}\big(\mathbf{w}^{(i)}\big)$

parameter vectors can only change over saturated links !

# Person. vs. Globalization

small $\lambda$, edges easily saturated

large $\lambda$, edges hard to saturate

# References

AJ, "On the Duality Between Network Flows and Network Lasso," in *IEEE Signal Processing Letters*, vol. 27, pp. 940-944, 2020.

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

Y. SarcheshmehPour, Y. Tian, L. Zhang, AJ, "Networked Federated Learning", <i>arXiv e-prints</i>, 2021.