

Networked Federated Learning

Alexander Jung (Aalto University)

<https://www.linkedin.com/in/aljung/>



https://www.youtube.com/channel/UC_tW4Z_GfJ2WCnKDtwMuDUA

<https://twitter.com/alexjungaalto>



- GTVMin as NFL Principle
- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

• GTVMin as NFL Principle

- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

In a nutshell:

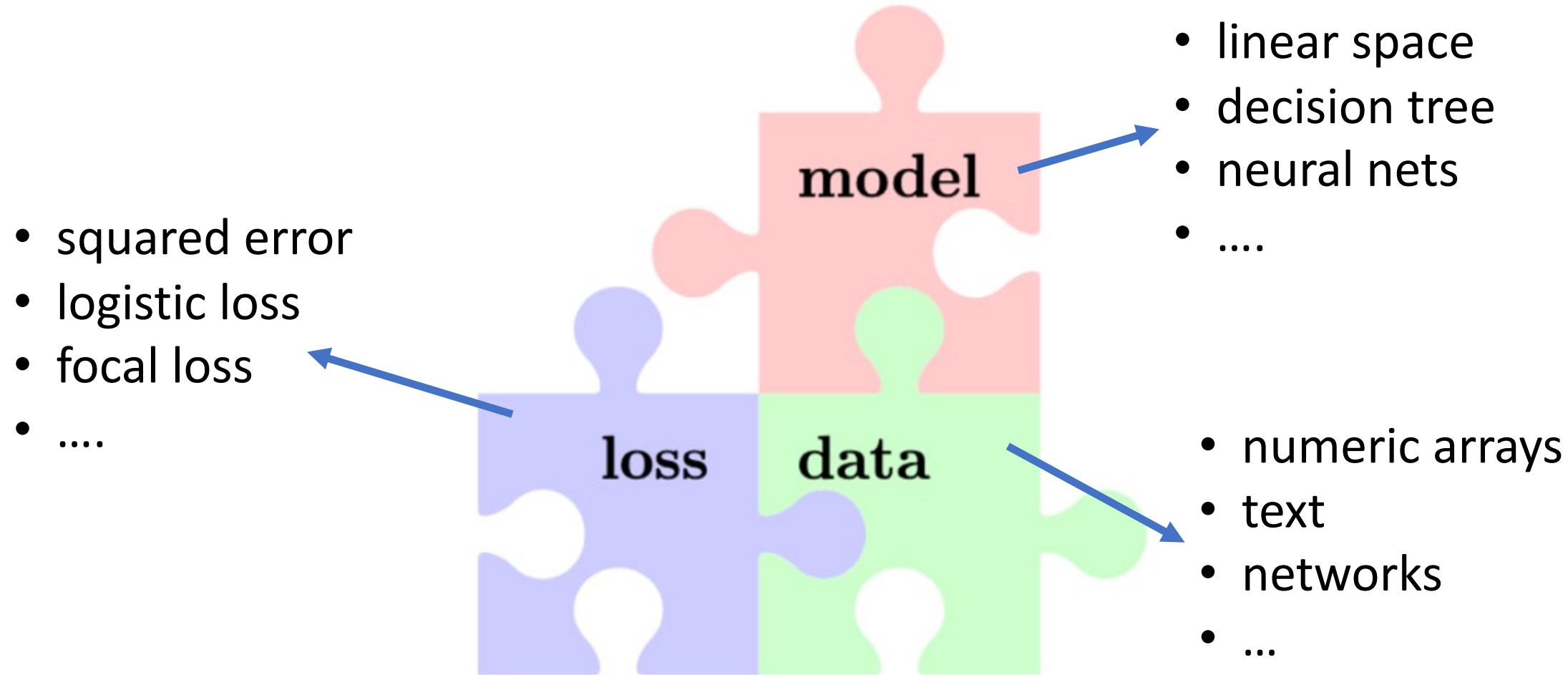
organize **data**, **models** and **computation** for
machine learning as **networks**.

Networked Federated Learning

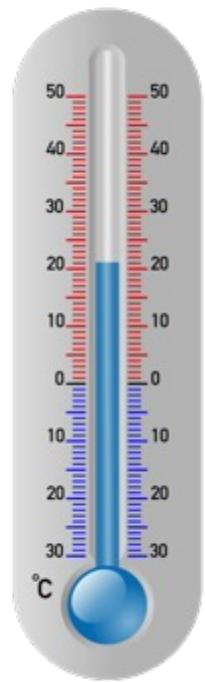
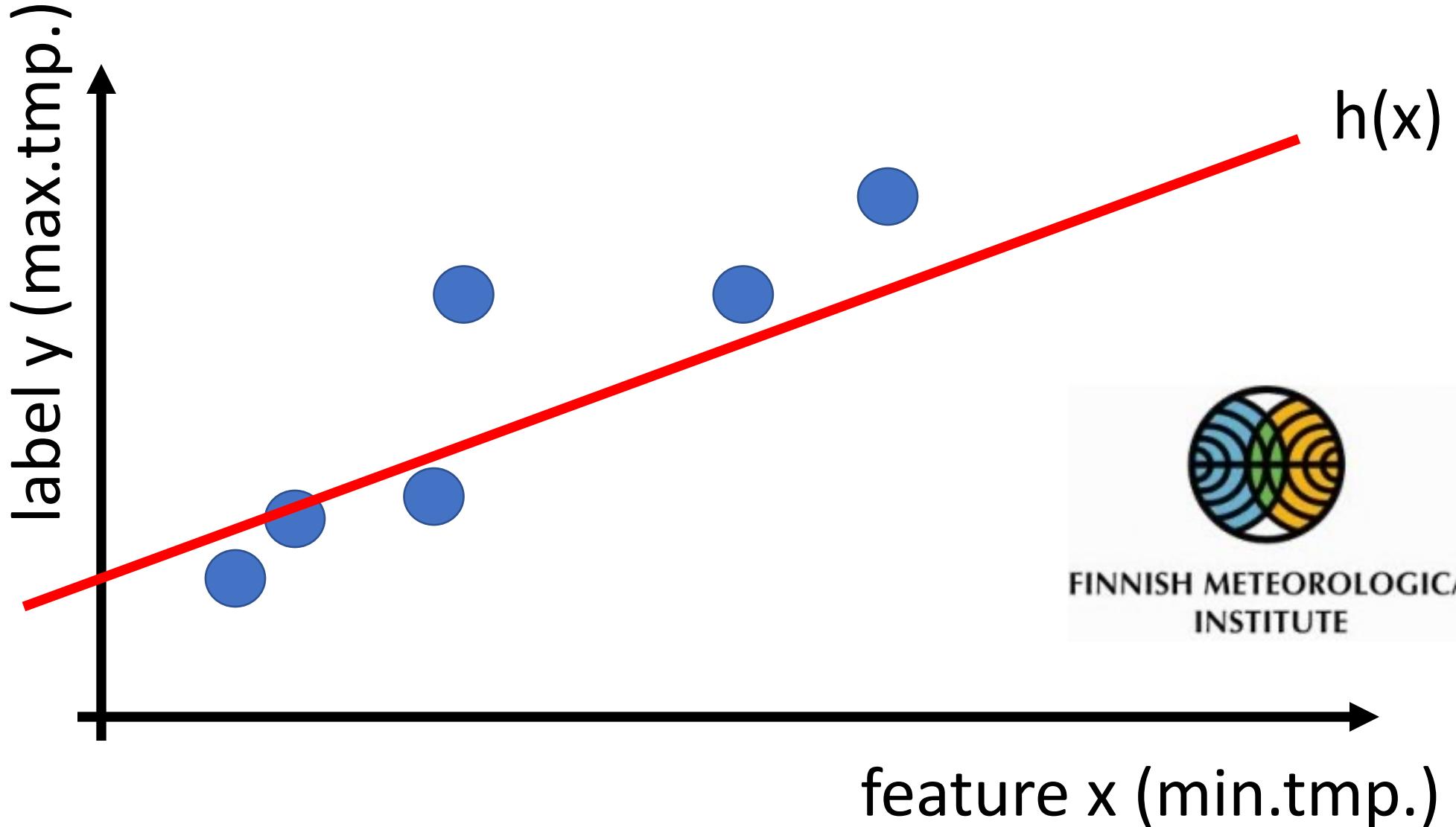
Federated Learning

Machine Learning

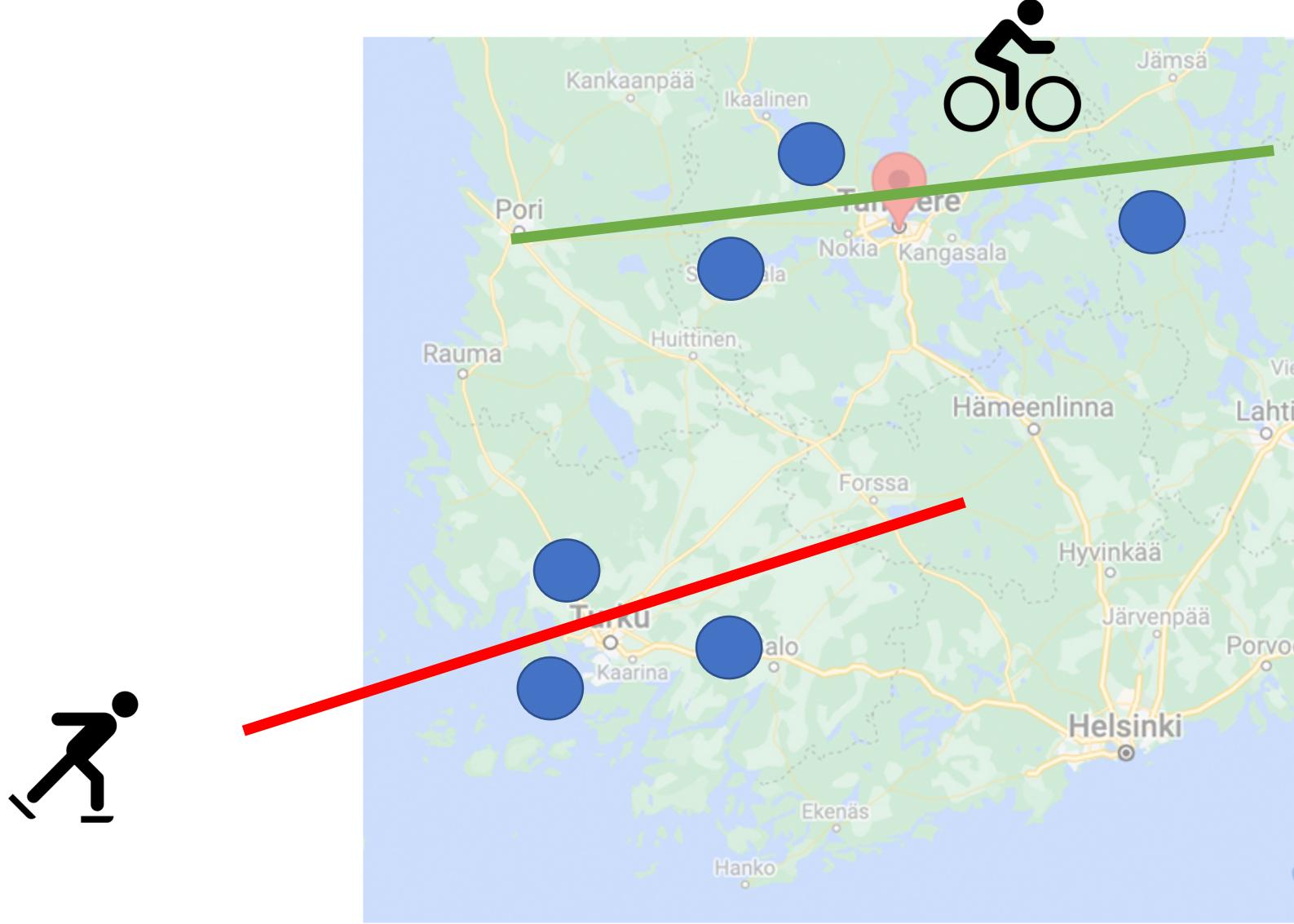
Three Components of ML



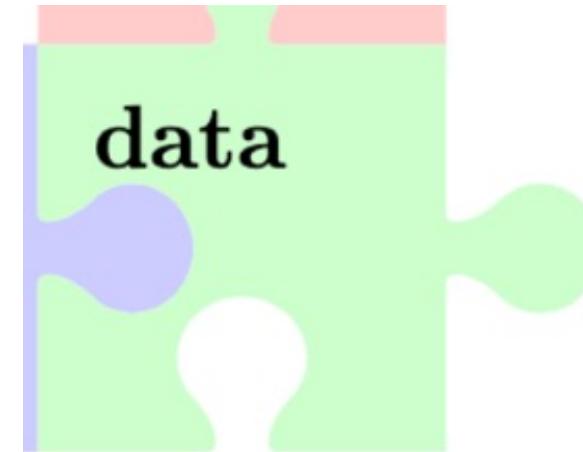
Plain Old Machine Learning.



Networked Federated Learning



Networked Data



Networked Data=Graph Database

The screenshot shows the neo4j Developer website. The top navigation bar includes the neo4j logo, "Developer", and links for "Docs", "Labs", and "Get Help". On the left, a sidebar under "DEVELOPER GUIDES" for "For Beginners" lists "Getting Started", "What is a Graph Database?", "Intro to Graph DBs Video Series", "Concepts: RDBMS to Graph", "Concepts: NoSQL to Graph", and "Getting Started Resources". The main content area shows the title "What is a Graph Database?", a "Beginner" skill level indicator, and a descriptive paragraph about graph databases.

Developer Guides / Getting Started / What is a Graph Database?

What is a Graph Database?

Beginner

A graph database stores nodes and relationships instead of tables, or documents. Data is stored just like you might sketch ideas on a whiteboard. Your data is stored without restricting it to a pre-defined model, allowing a very flexible way of thinking about and using it.

<https://neo4j.com/developer/graph-database/>

Weather Stations.



FINNISH METEOROLOGICAL
INSTITUTE

ImageNet.

“...ImageNet is an image database organized according to the WordNet hierarchy (currently only the nouns), in **which each node** of the hierarchy is depicted by **hundreds and thousands of images...**”

<https://image-net.org/>

WordNet.

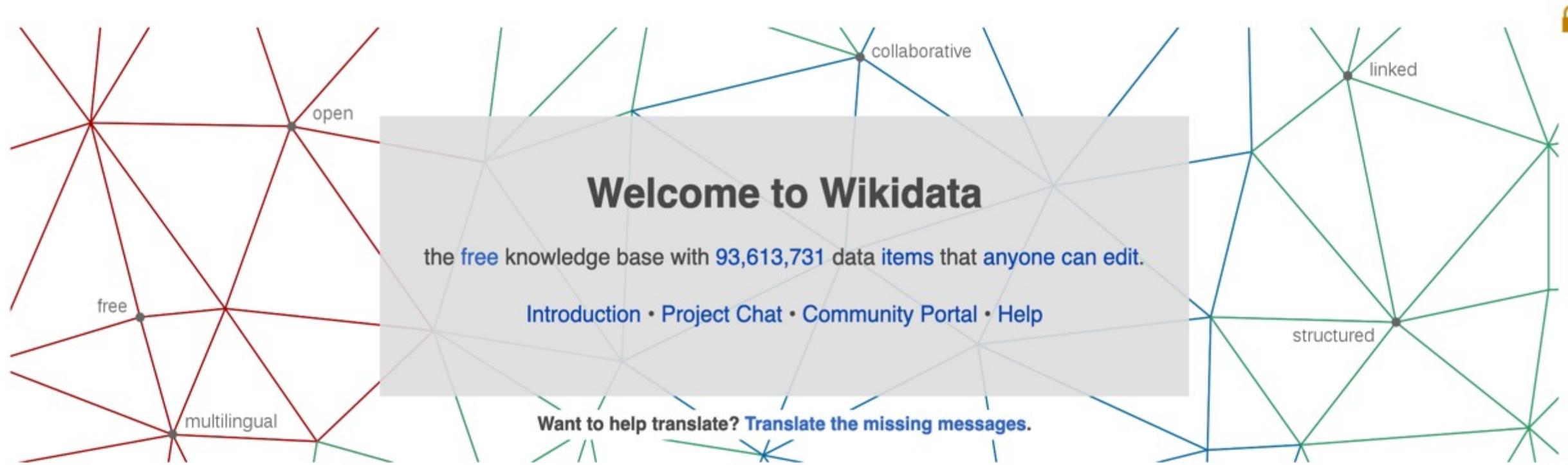
“...Nouns, verbs, adjectives and adverbs are grouped into sets of cognitive synonyms (synsets), each expressing a distinct concept... The resulting **network of meaningfully related words** and concepts can be navigated....”

<https://wordnet.princeton.edu/>

A. Jung, Networked Federated Learning

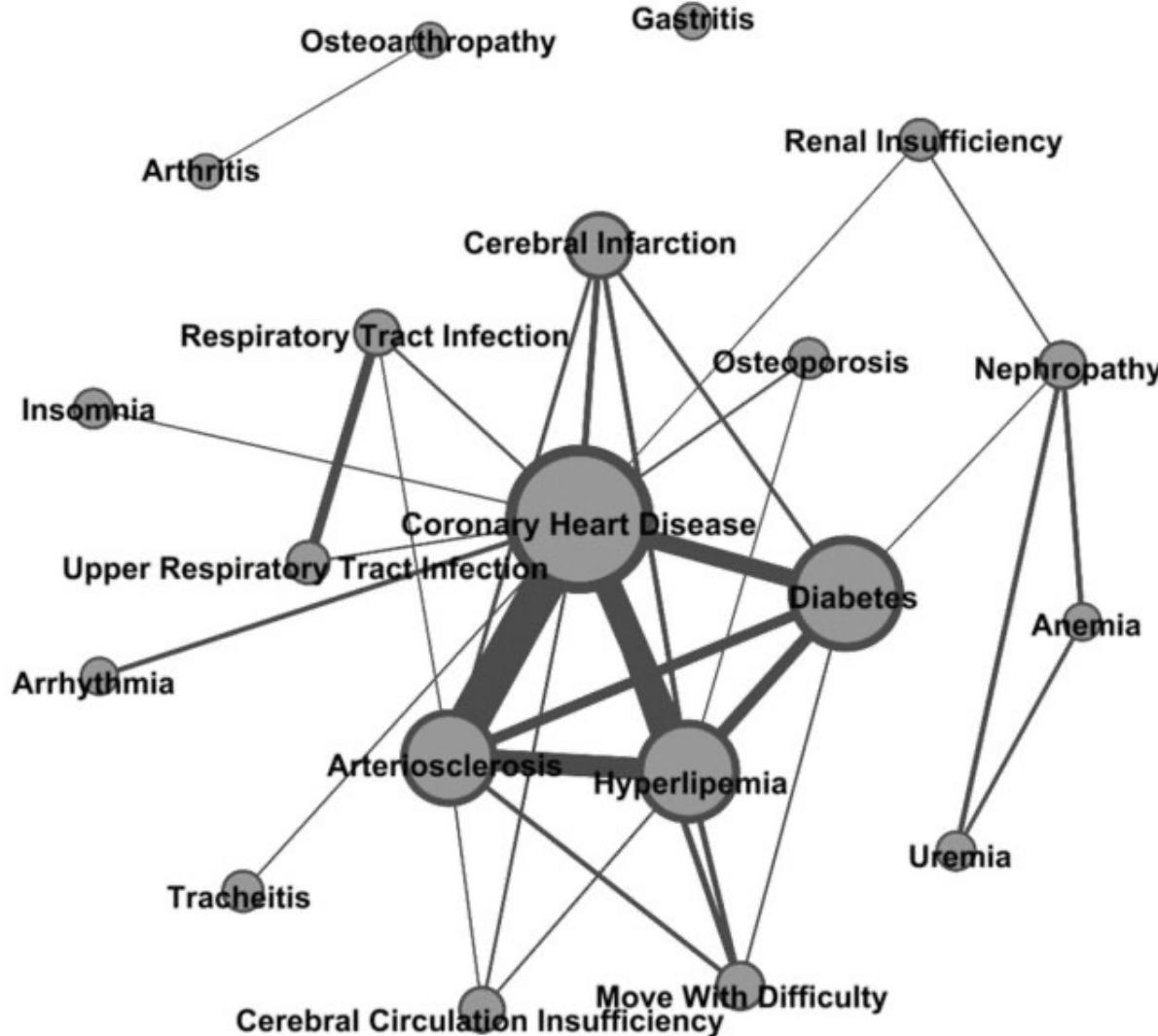
13

Wikidata.



https://www.wikidata.org/wiki/Wikidata:Main_Page

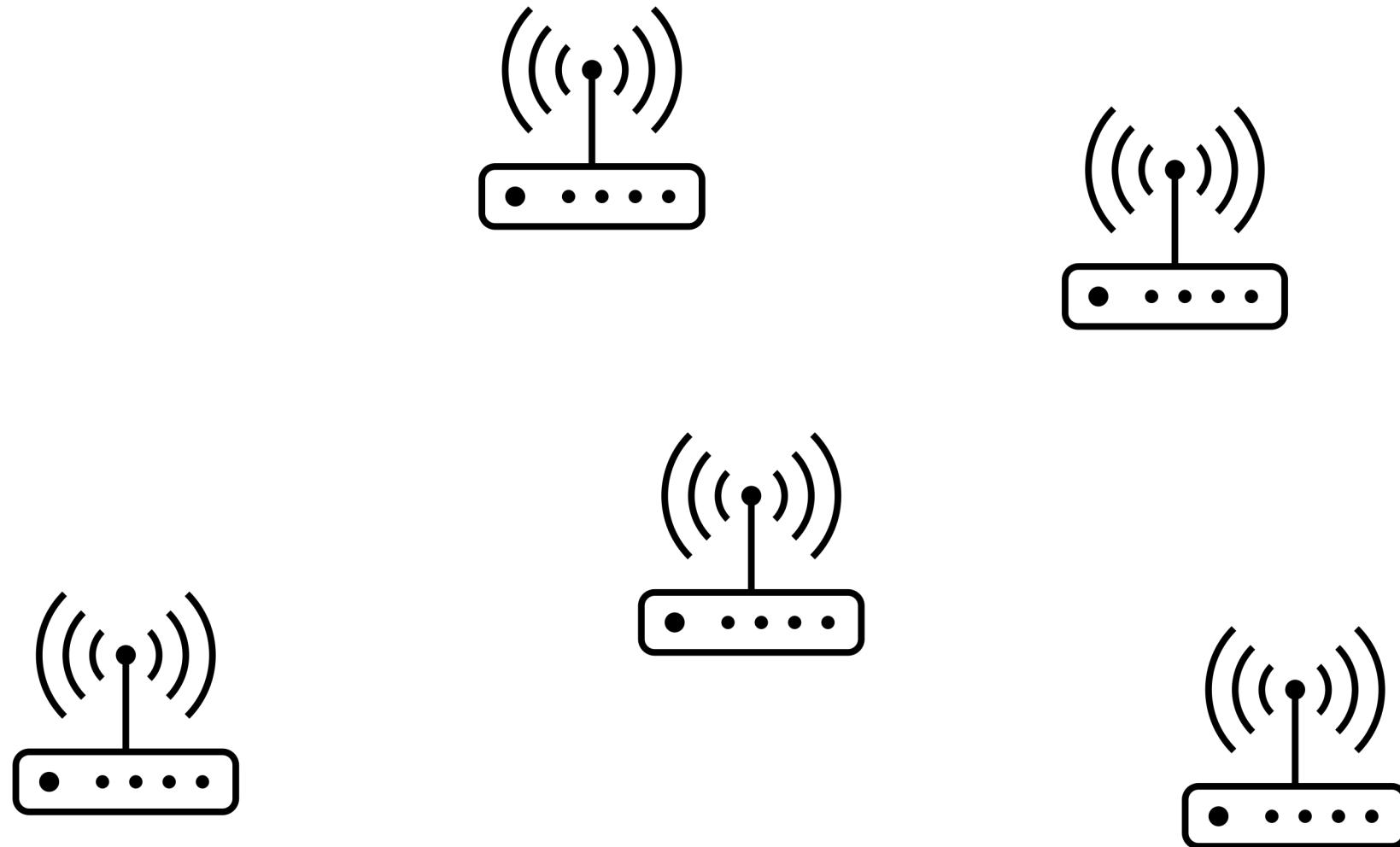
Diseases.



Liu, Jiaqi et.al..

Comorbidity Analysis According to Sex and Age in Hypertension Patients in China.
International Journal of Medical Sciences. 13. 99-107. 10.7150/ijms.13456.

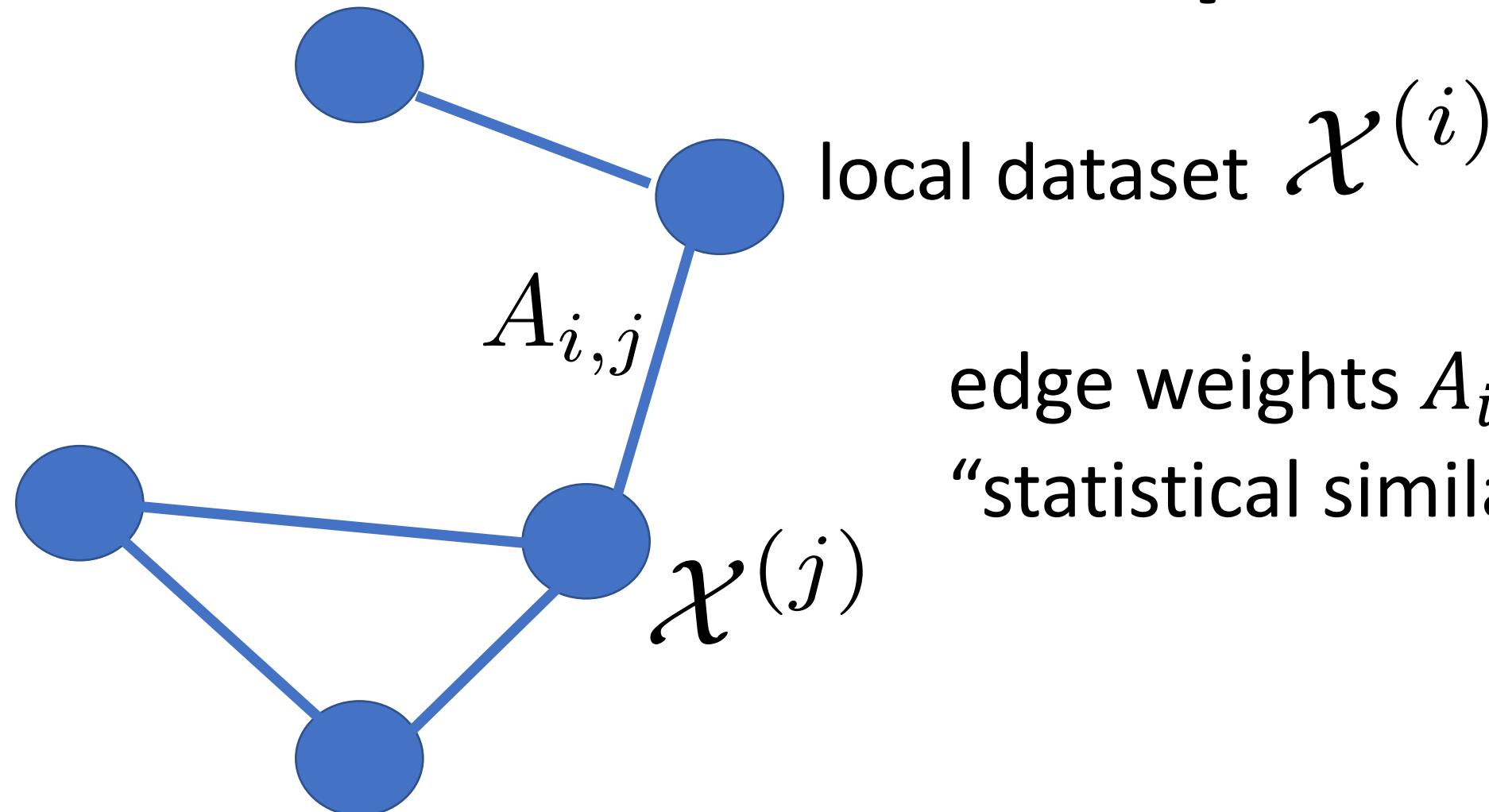
WSN.



Anchors.



Abstraction – The Empirical Graph.



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

How To Measure Statistical Sim.?

```
>>> from scipy.stats import ks_2samp
>>> import numpy as np
>>>
>>> np.random.seed(12345678)
>>> x = np.random.normal(0, 1, 1000)
>>> y = np.random.normal(0, 1, 1000)
>>> z = np.random.normal(1.1, 0.9, 1000)
>>>
>>> ks_2samp(x, y)
Ks_2sampResult(statistic=0.022999999999999909, pvalue=0.95189016804849647)
>>> ks_2samp(x, z)
Ks_2sampResult(statistic=0.4180000000000004, pvalue=3.7081494119242173e-77)
```

<https://stackoverflow.com/questions/10884668/two-sample-kolmogorov-smirnov-test-in-python-scipy>

https://en.wikipedia.org/wiki/Kolmogorov%E2%80%93Smirnov_test

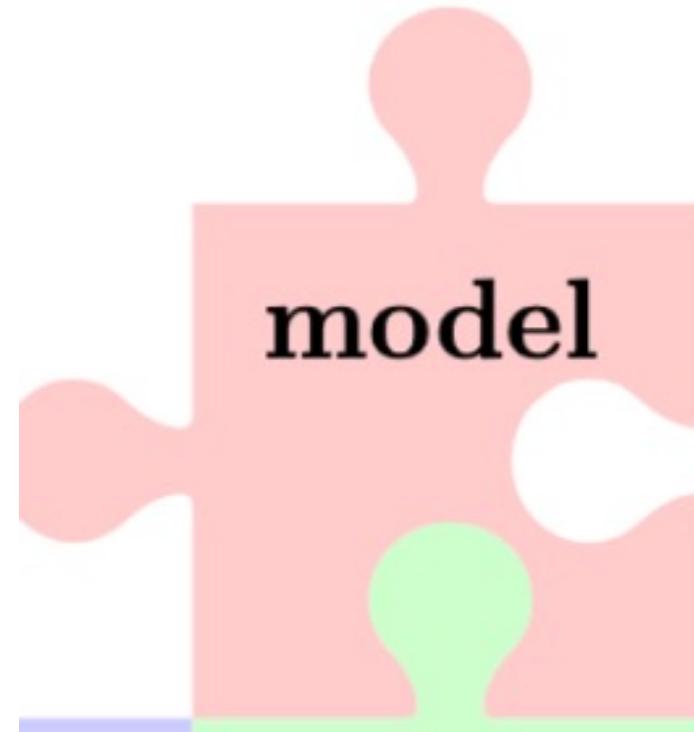
Geometric Dataset Distances via Optimal Transport

David Alvarez-Melis¹ Nicolò Fusi¹

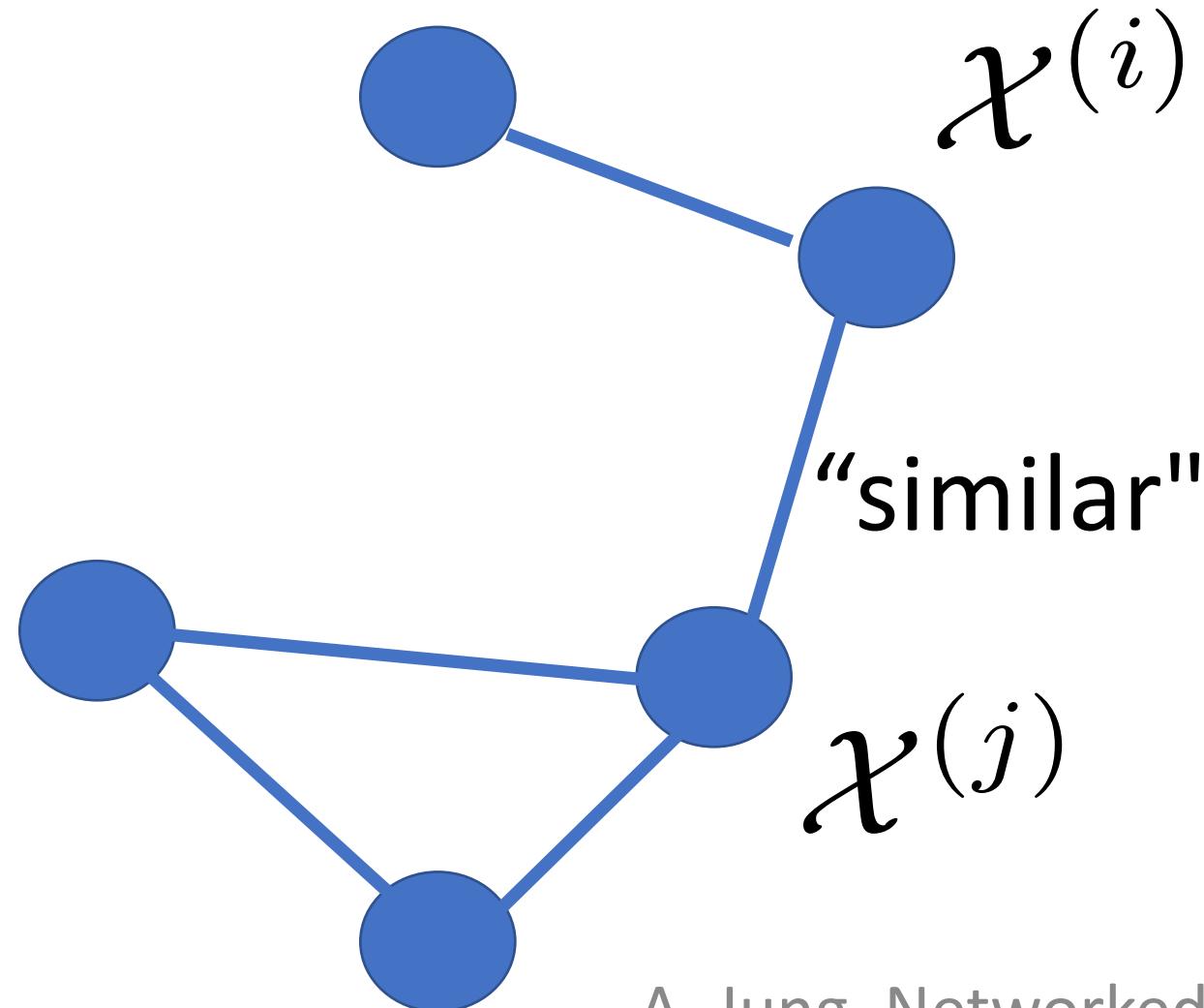
“In this work we propose an alternative notion of distance between datasets that (i) is model-agnostic, (ii) does not involve training,...

<https://arxiv.org/pdf/2002.02923.pdf>

Networked Models



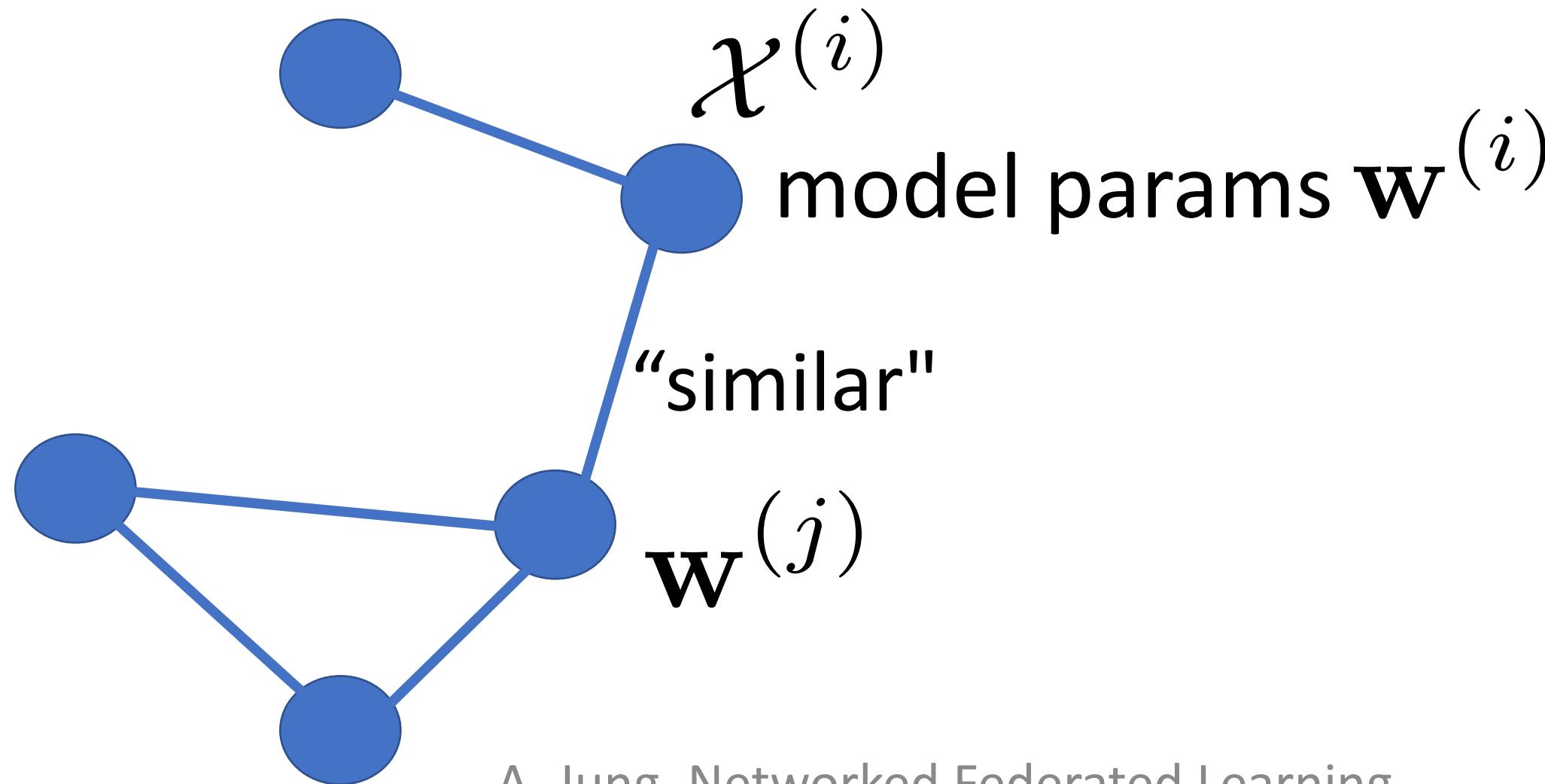
Networked Models.



local model for each node

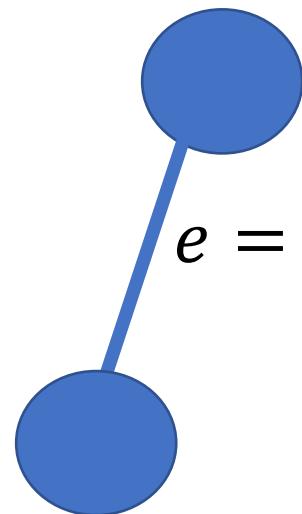
couple models at
connected nodes

Networked Parametric Models.



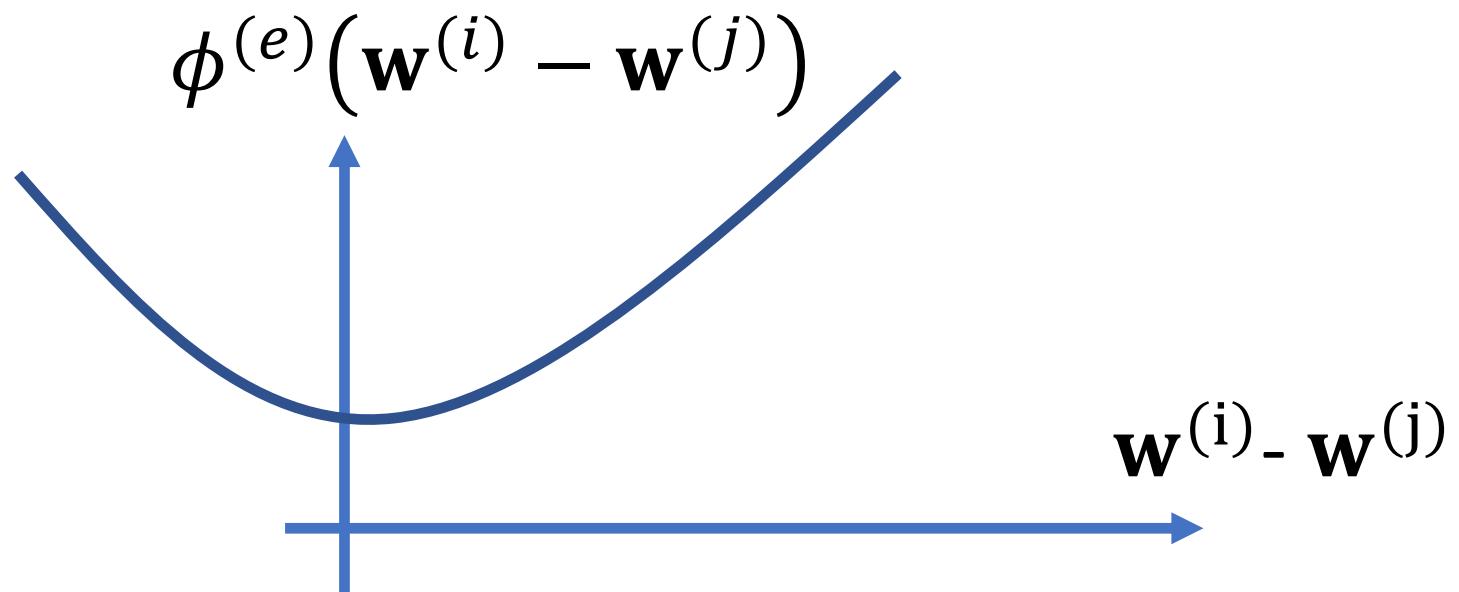
Smoothness/Clustering Assumption.

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

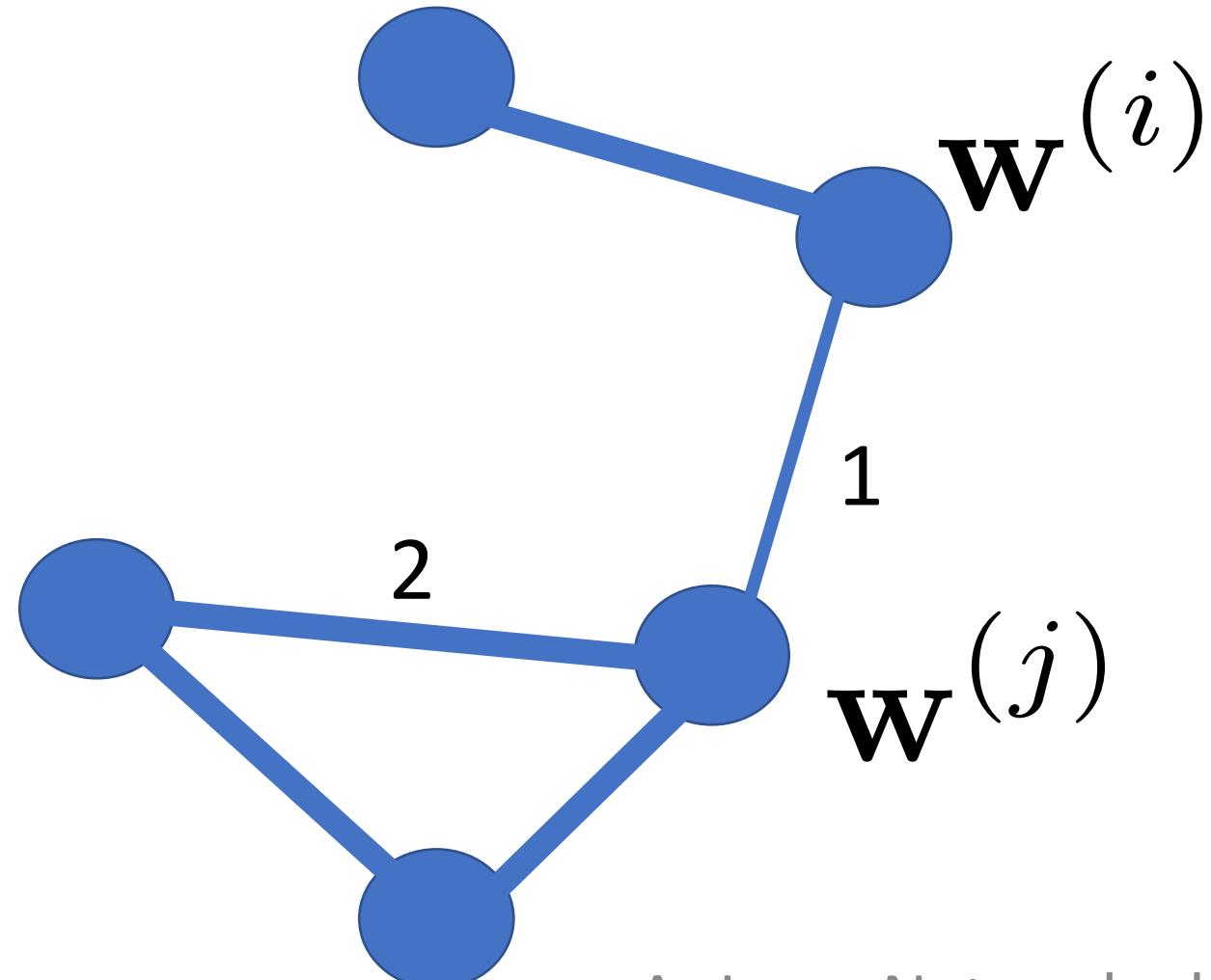


require similar params at ends of edge e

penalty function measures “**tension**”



Generalized Total Variation (GTV)



force params of well connected
nodes to be similar by requiring
a small GTV

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

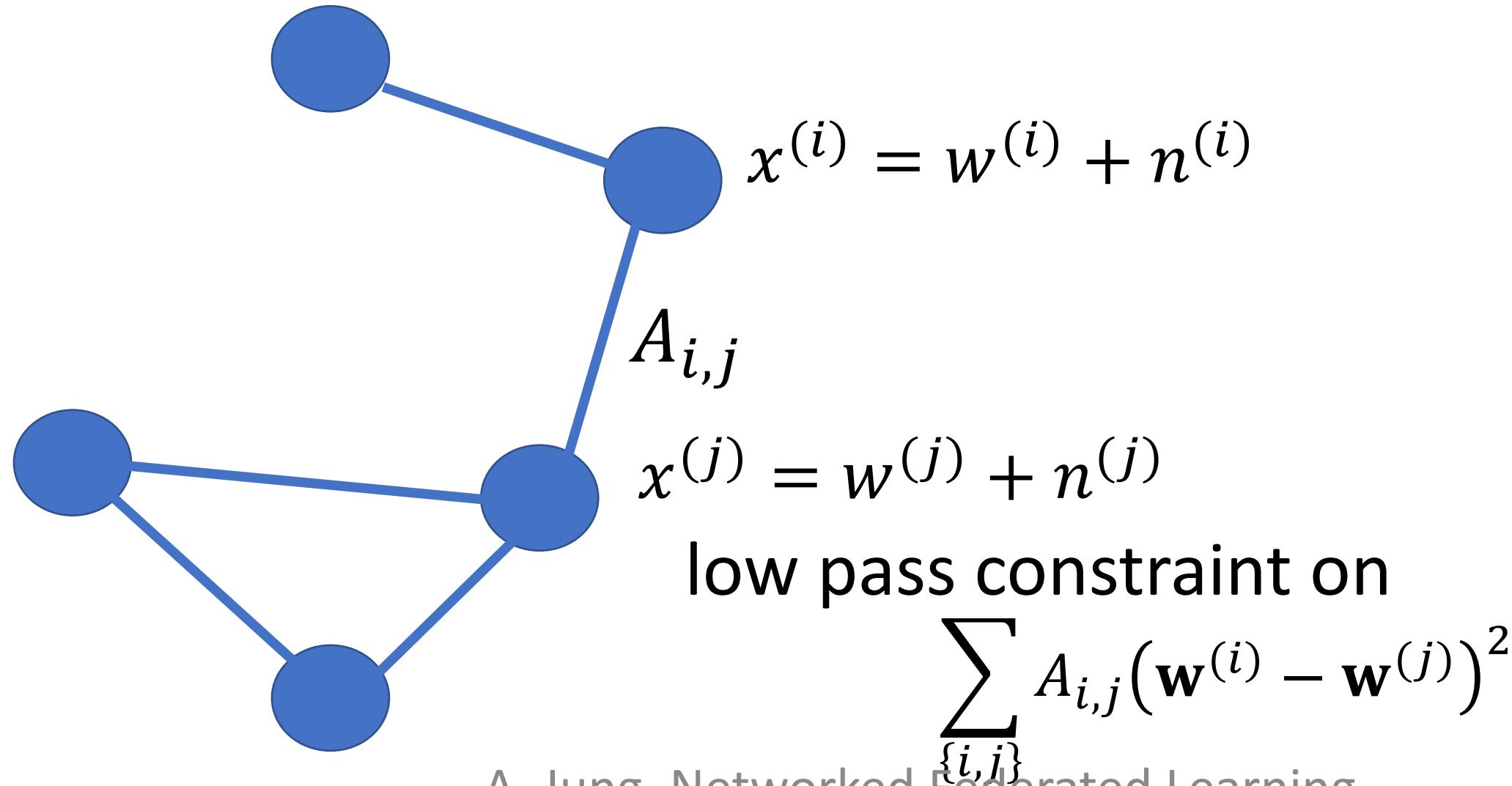
Two Special Cases of GTV.

total variation $\phi(\mathbf{u}) = \|\mathbf{u}\|_2$

graph Laplacian quadratic form is GTV with

$$\phi(\mathbf{u}) = \|\mathbf{u}\|_2^2$$

Smooth Graph Signals.



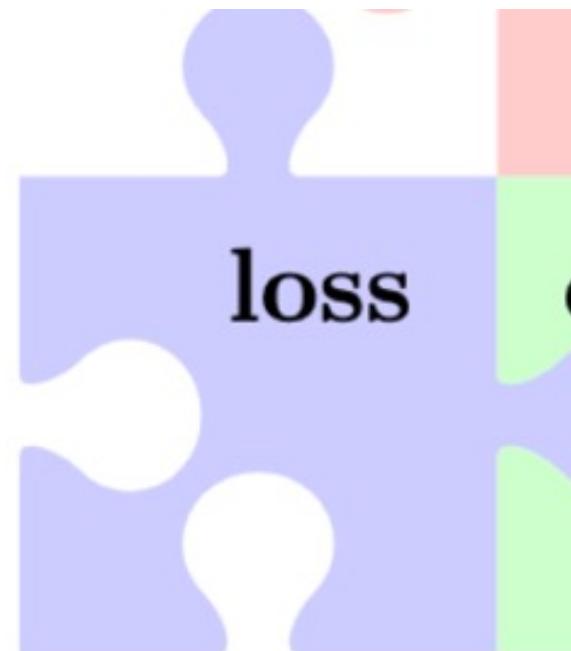
From now on,

GTVMin with penalty being a norm

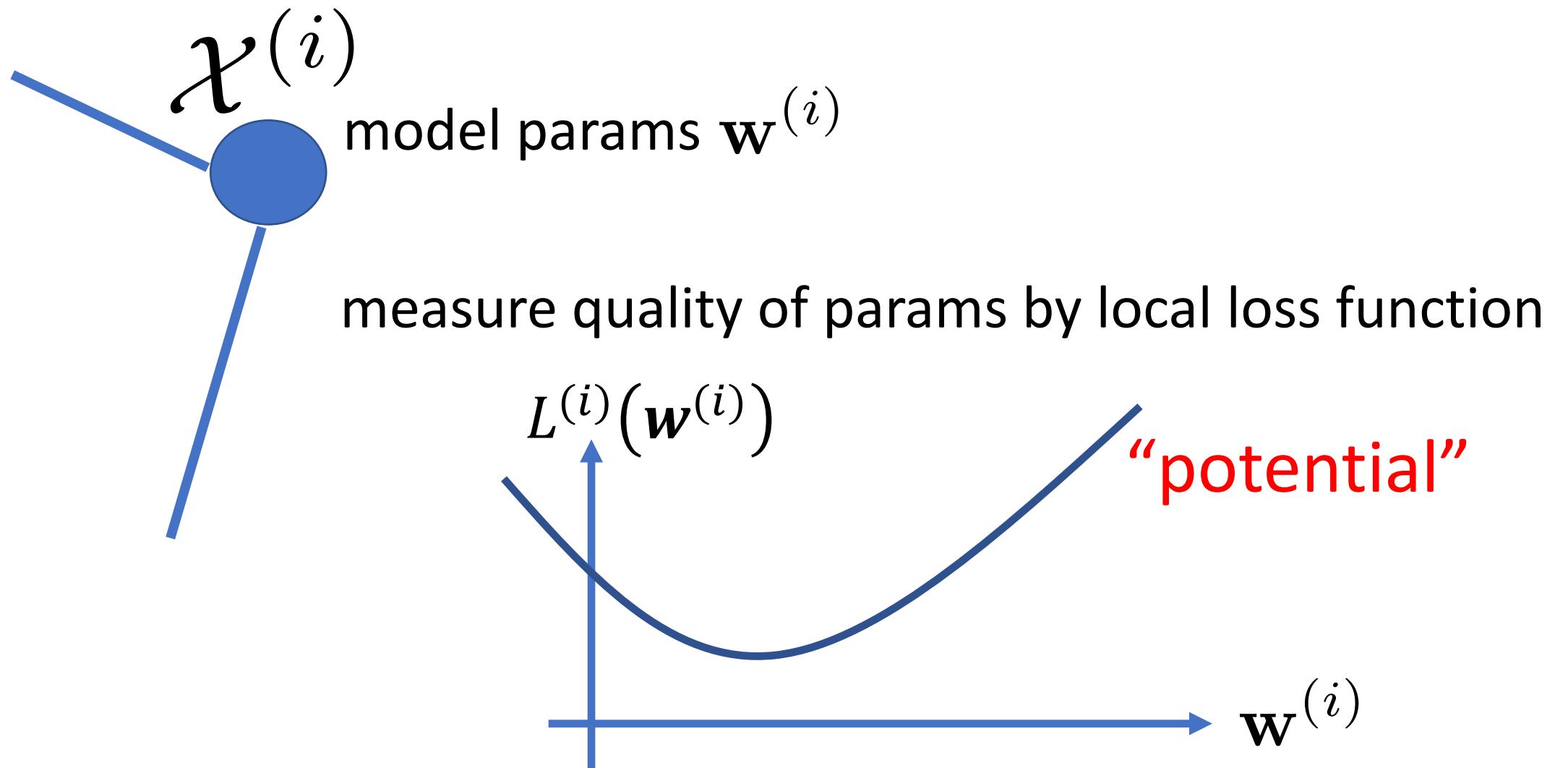
$$\phi(\mathbf{u}) = \|\mathbf{u}\|$$

(unless otherwise stated)

GTV Minimization.



Local Loss Functions.



GTV Minimization.

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

average local loss

increasing λ

“clusteredness”

training set M

Special Case: Network Lasso.

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

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

Special Case: “MOCHA”

$$\min_w \sum_{i \in M} L^{(i)}(w^{(i)}) + \lambda \sum_{\{i,j\}} A_{i,j} \|w^{(i)} - w^{(j)}\|^2$$

<https://papers.nips.cc/paper/7029-federated-m...> ▾ PDF

Federated Multi-Task Learning - NIPS Proceedings

by V Smith · 2017 · Cited by 501 — 3.2 MOCHA: A Framework for **Federated Multi-Task Learning**. In the **federated** setting, the aim is to train statistical models directly on the edge, and thus we solve (1) while assuming that the data $\{X_1, \dots, X_m\}$ is distributed across m nodes or devices.

- GTVMin as NFL Principle
- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

“Massaging” GTV Minimization.

$$\hat{\mathbf{w}} \in \arg \min_{\mathbf{w} \in \mathcal{W}} f(\mathbf{w}) + g(\mathbf{D}\mathbf{w})$$

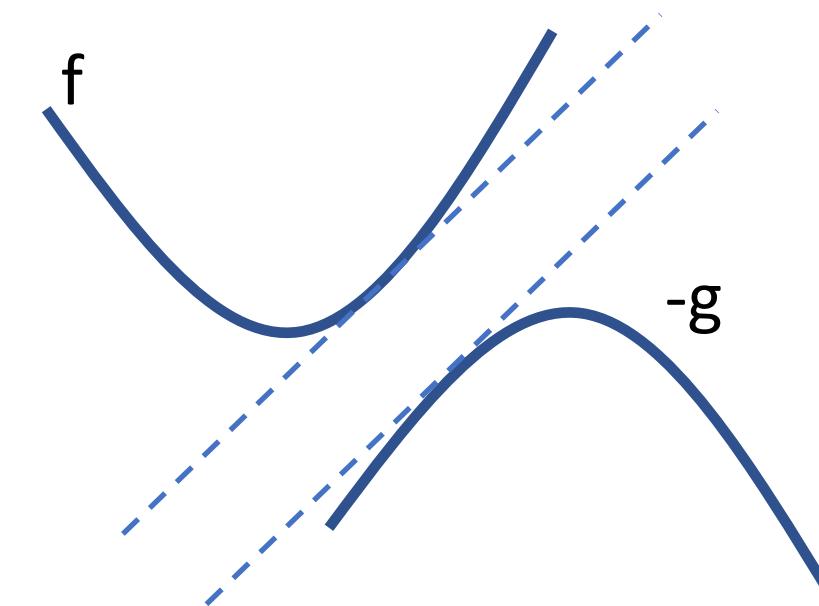
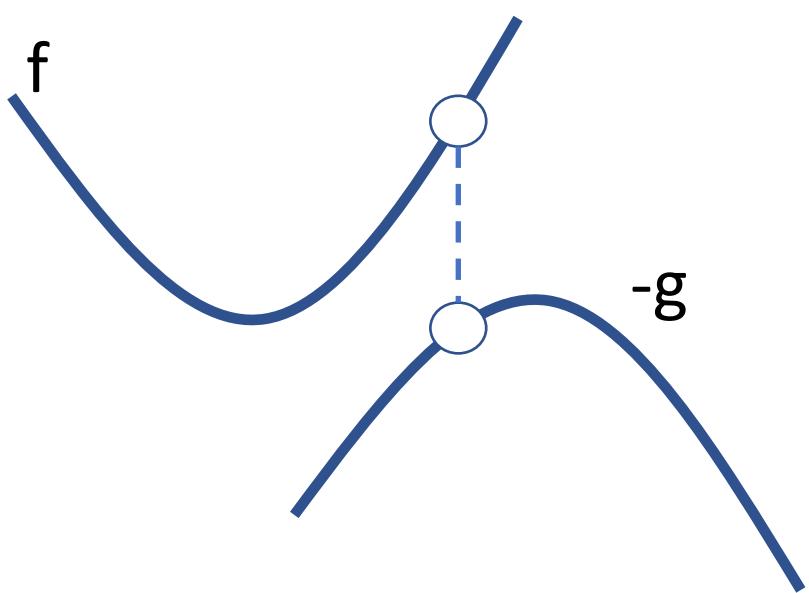
with $f(\mathbf{w}) := \sum_{i \in \mathcal{V}} L_i(\mathbf{w}^{(i)})$, and $g(\mathbf{u}) := \lambda \sum_{e \in \mathcal{E}} A_e \phi(\mathbf{u}^{(e)})$.

with incidence matrix/operator

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

Fenchel's Duality

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



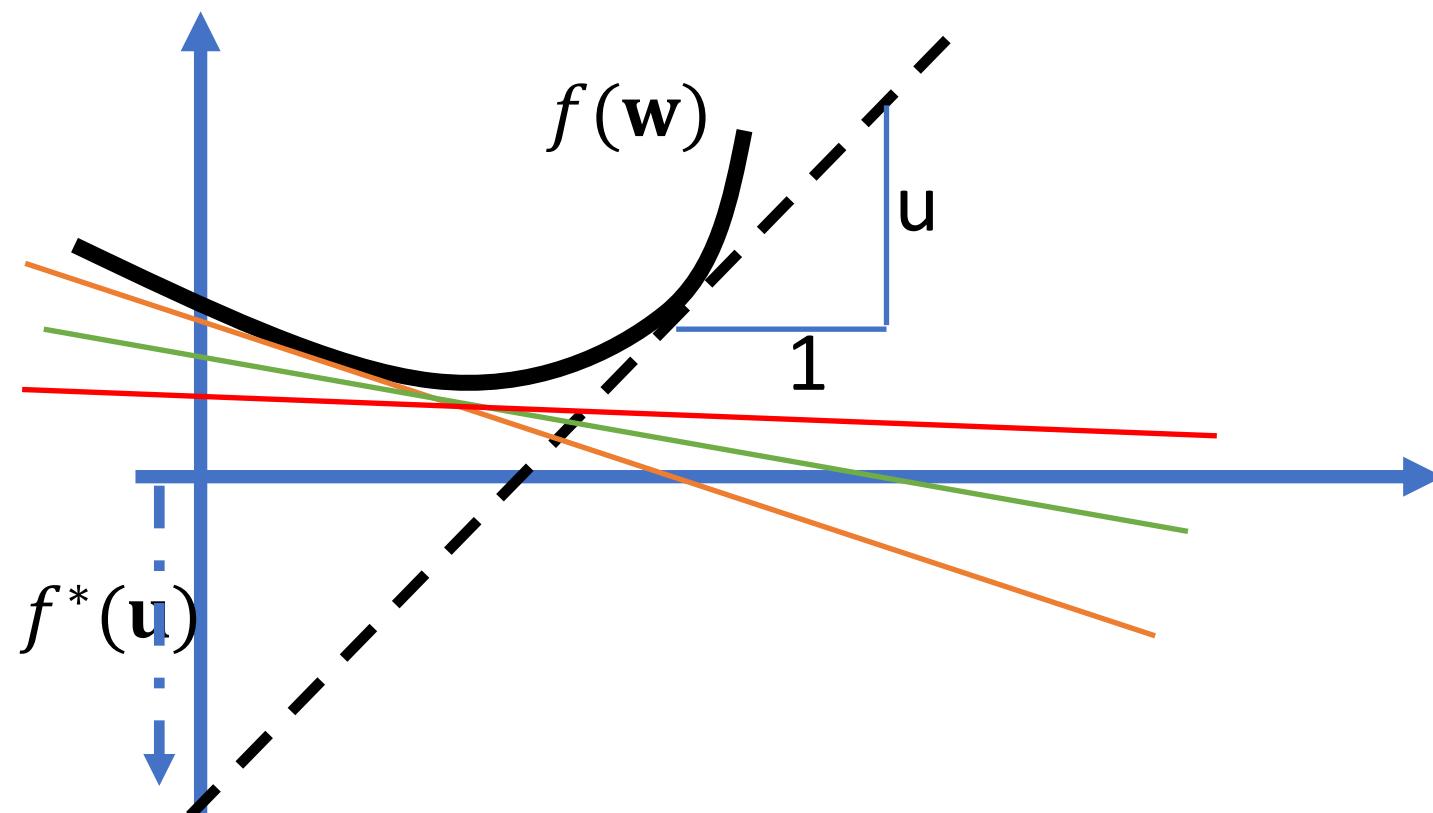
R. T. Rockafellar, *Convex Analysis*. Princeton, NJ: Princeton Univ. Press, 1970.

https://en.wikipedia.org/wiki/Fenchel%27s_duality_theorem

Convex Conjugate.

$$f^*(\mathbf{w}) := \sup_{\mathbf{z} \in \mathbb{R}^{n|\mathcal{V}|}} \mathbf{w}^T \mathbf{z} - f(\mathbf{z})$$

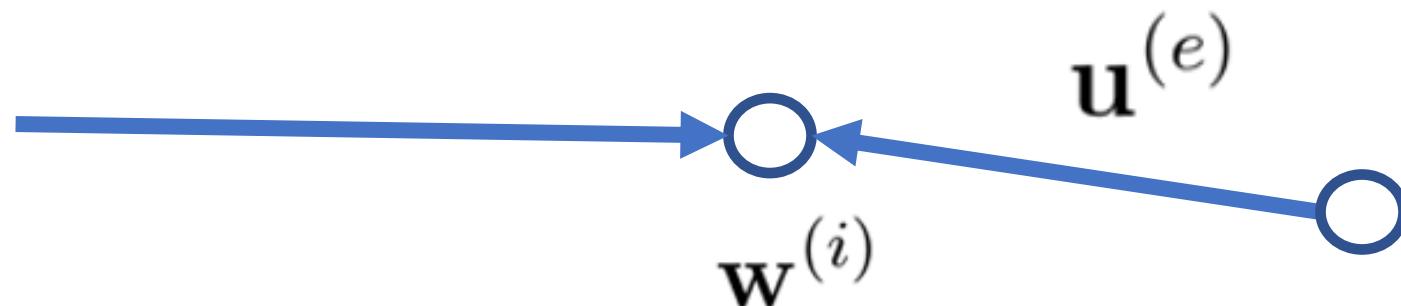
$$g^*(\mathbf{u}) := \sup_{\mathbf{z} \in \mathbb{R}^{n|\mathcal{E}|}} \mathbf{u}^T \mathbf{z} - g(\mathbf{z})$$



The Dual of GTVMin.

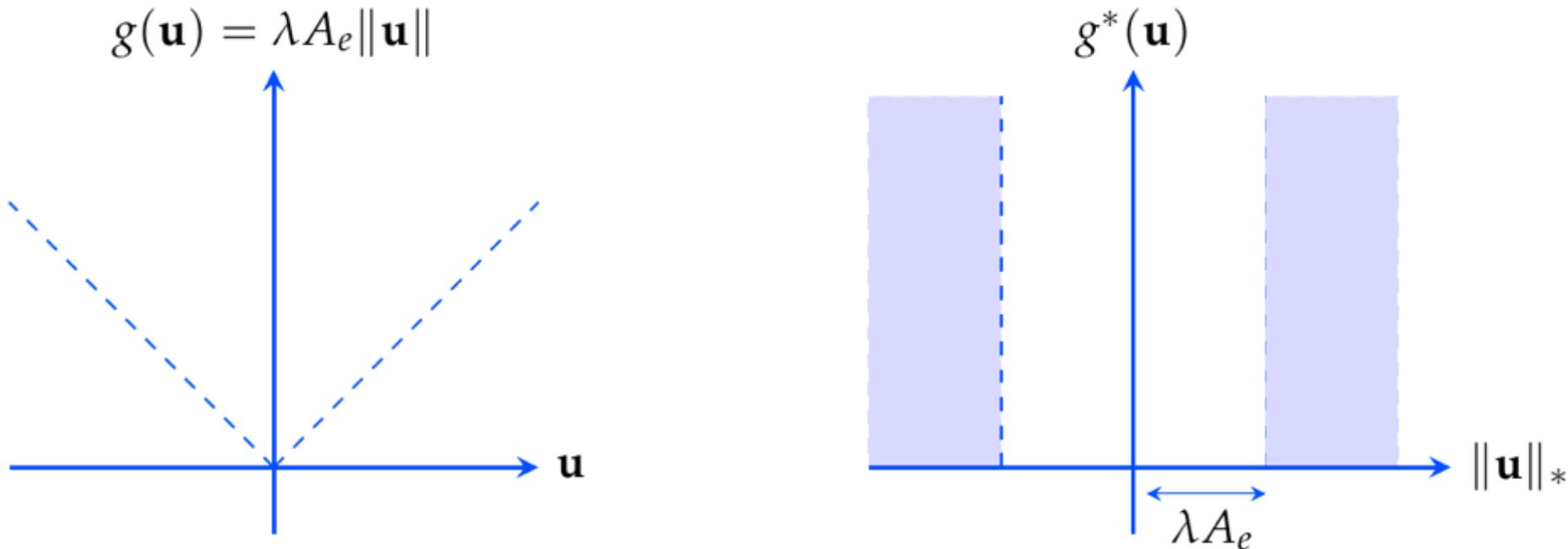
$$\max_{\mathbf{u} \in \mathcal{U}} - \sum_{i \in \mathcal{V}} L_i^* (\mathbf{w}^{(i)}) - \lambda \sum_{e \in \mathcal{E}} A_e \phi^* \left(\mathbf{u}^{(e)} / (\lambda A_e) \right)$$

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 $i \in \mathcal{V}$.



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

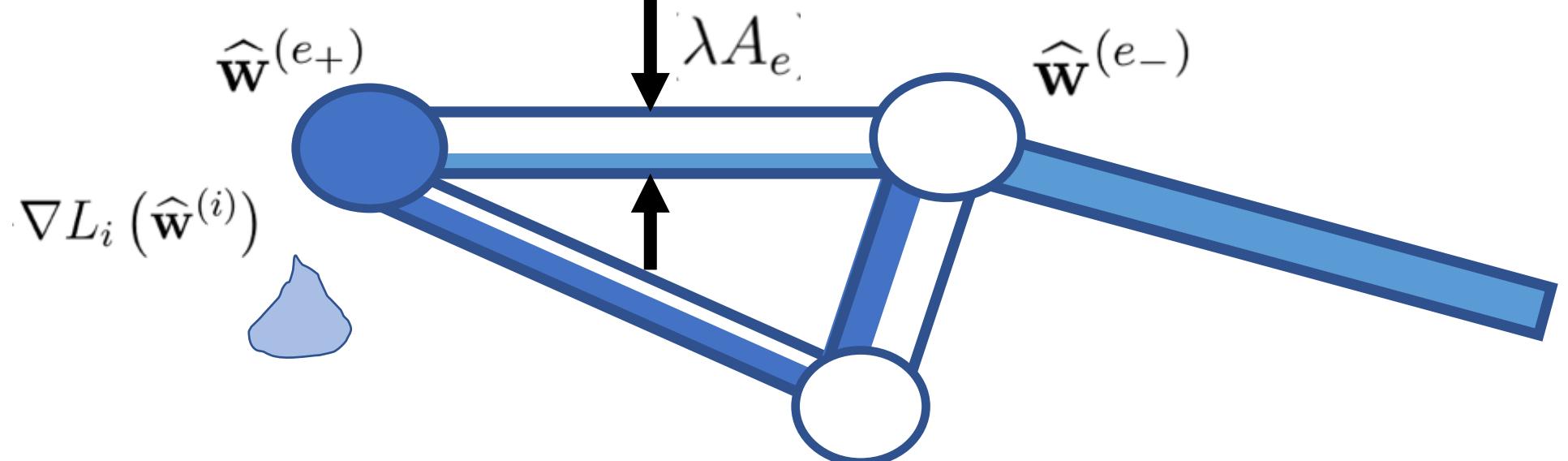
Convex Conjugate of Norm.



Primal and Dual Optimality.

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

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



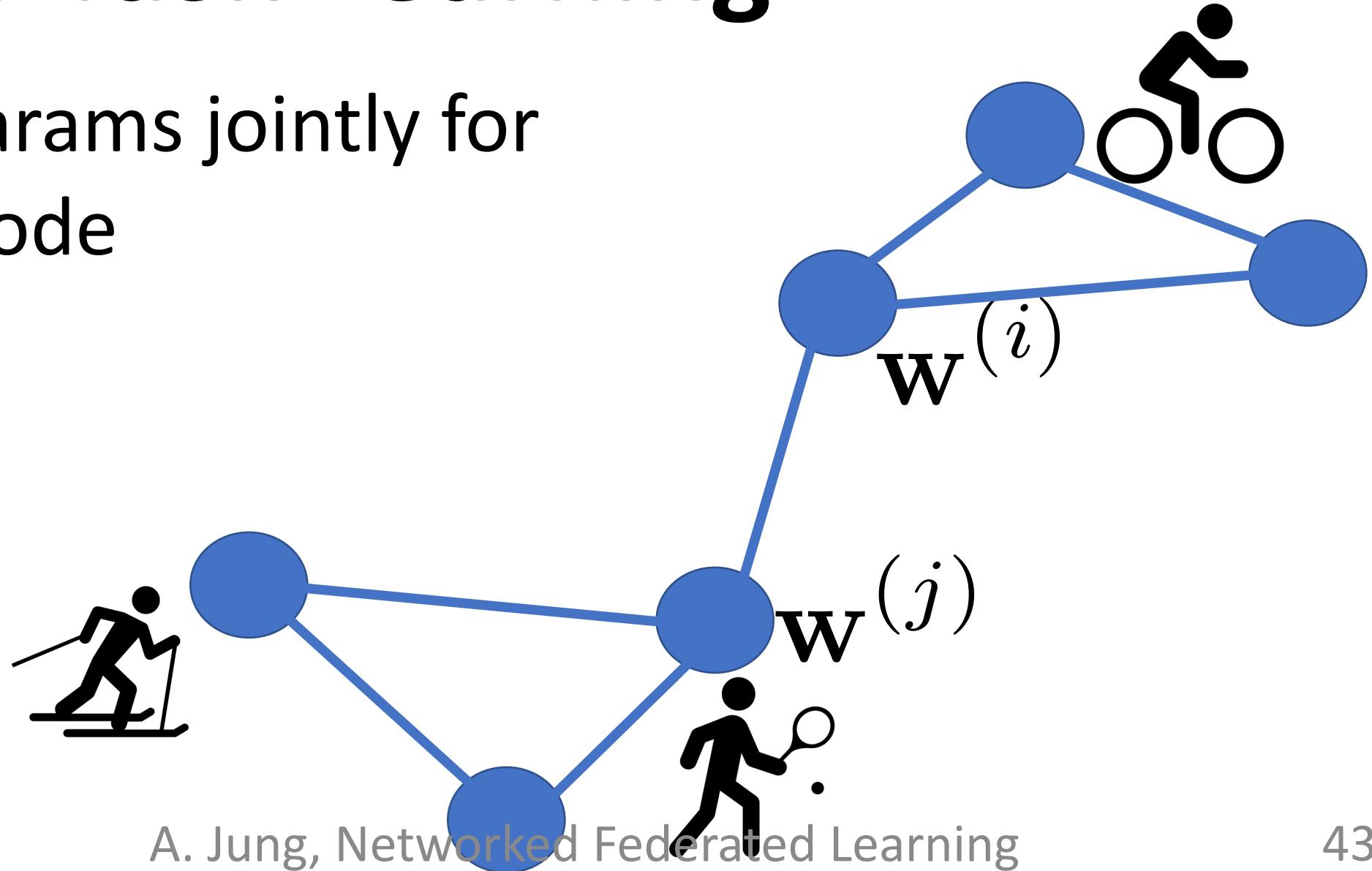
- GTVMin as NFL Principle
- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

Smooth Graph Sig. Recovery

$$\min_w \sum_{i \in M} (y^{(i)} - w^{(i)})^2 + \lambda \sum_{\{i,j\}} A_{i,j} (w^{(i)} - w^{(j)})^2$$

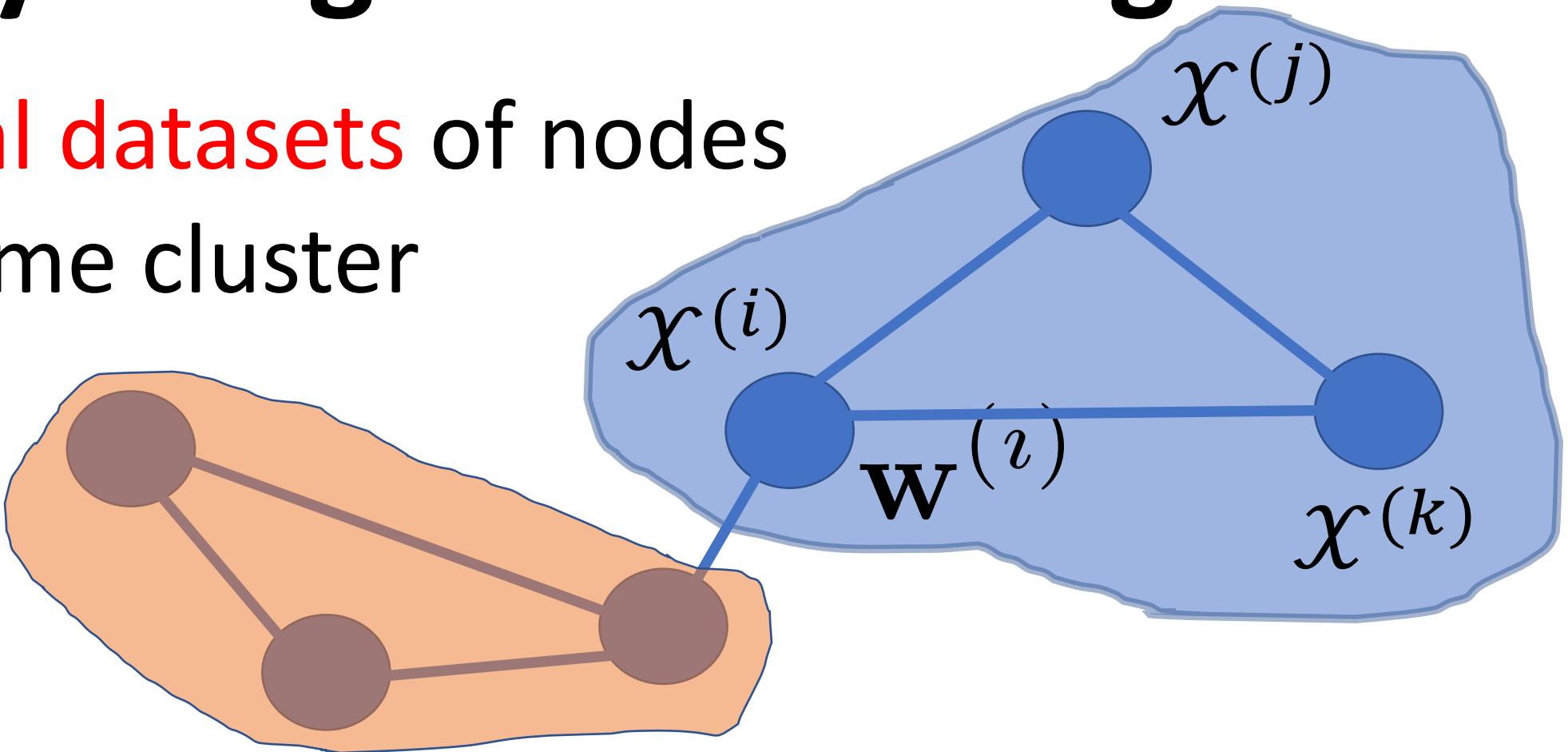
Multi-Task Learning

learn params jointly for
every node



Locally Weighted Learning

pool local datasets 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

$$\min_{\mathbf{w}} \sum_{i \in M} \left\| \mathbf{w}^{(i)} - \mathbf{a}^{(i)} \right\|^2 + \lambda \sum_{\{i,j\}} A_{i,j} \left\| \mathbf{w}^{(i)} - \mathbf{w}^{(j)} \right\|_p$$

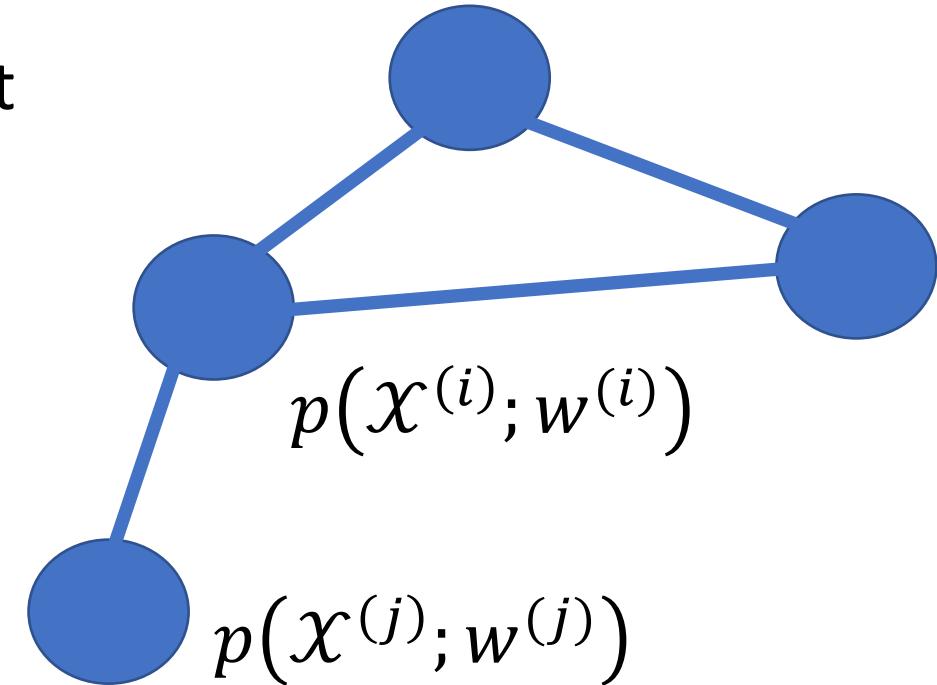
D. Sun, K.-C. Toh, Y. Yuan;

Convex Clustering: Model, Theoretical Guarantee and Efficient Algorithm, JMLR, 22(9):1–32, 2021

(Probabilistic) Graphical Model

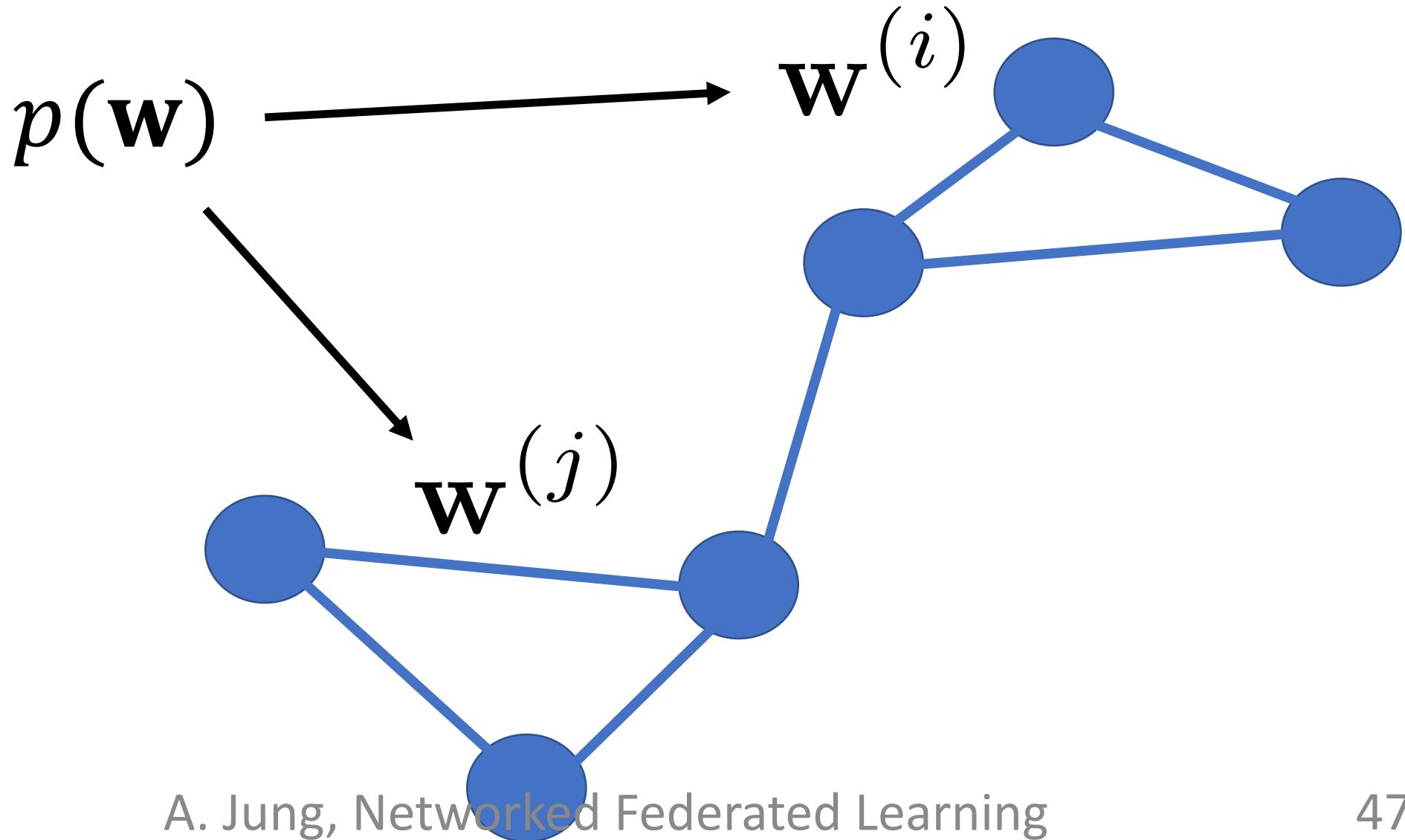
separate prob. space for each local dataset

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



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

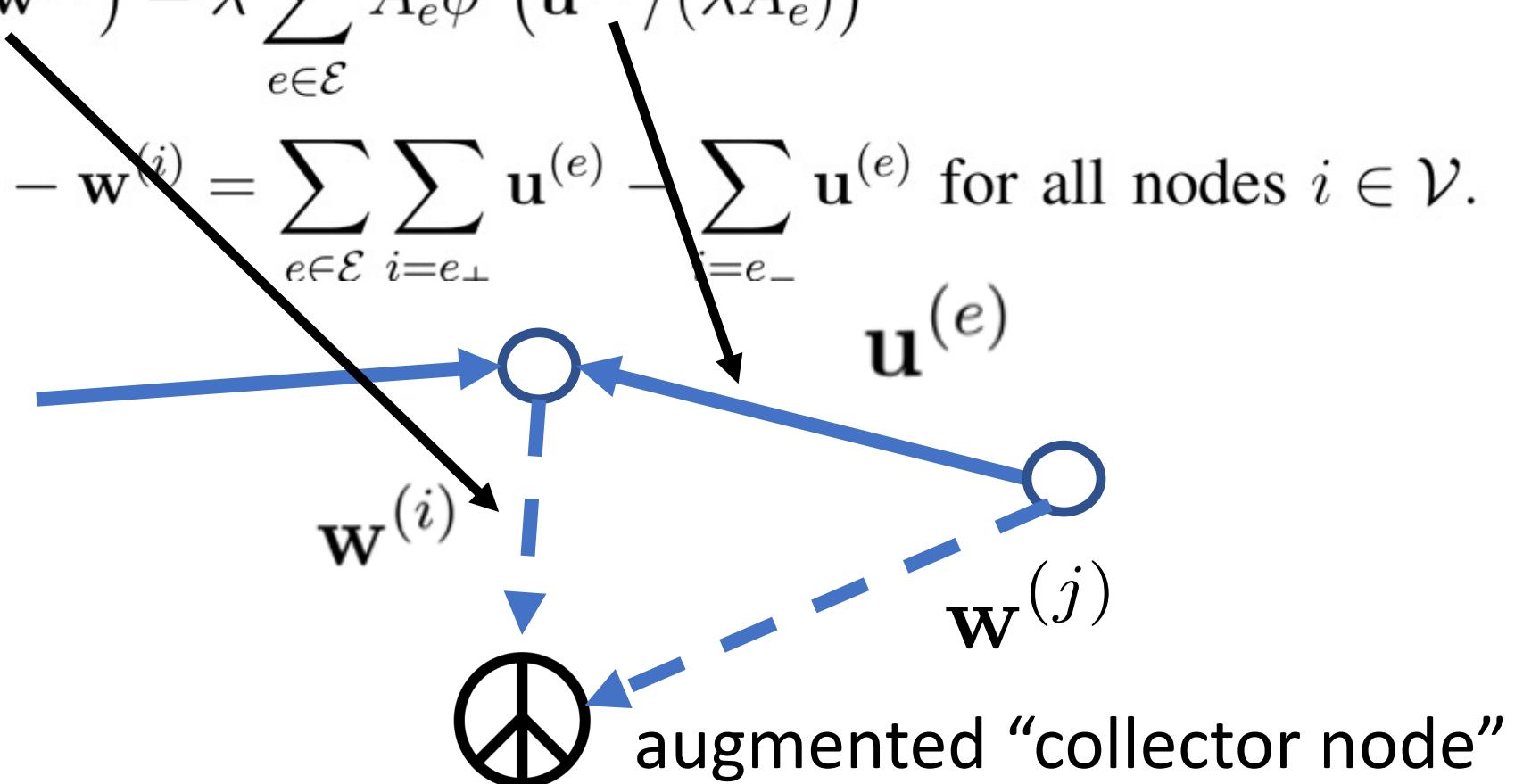
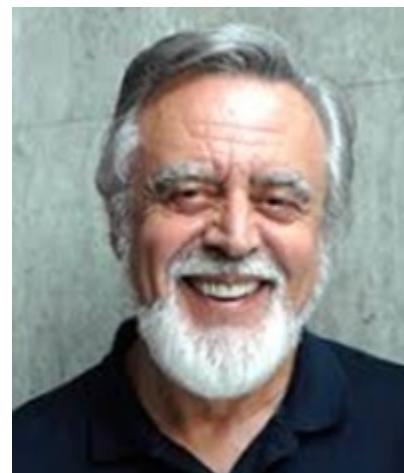
Approx. Hierarch. Bayes' Model



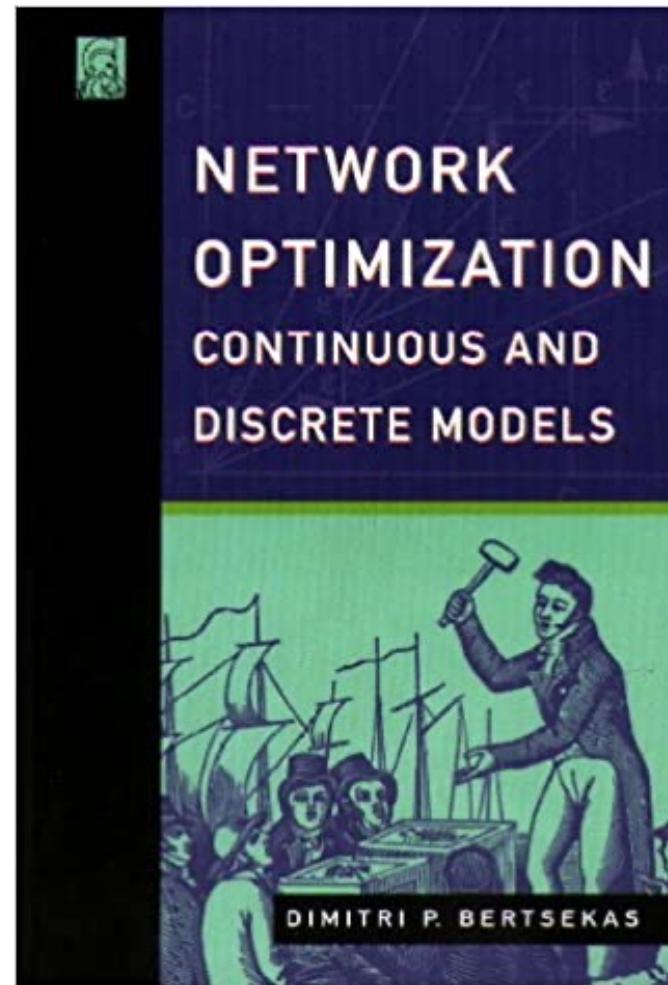
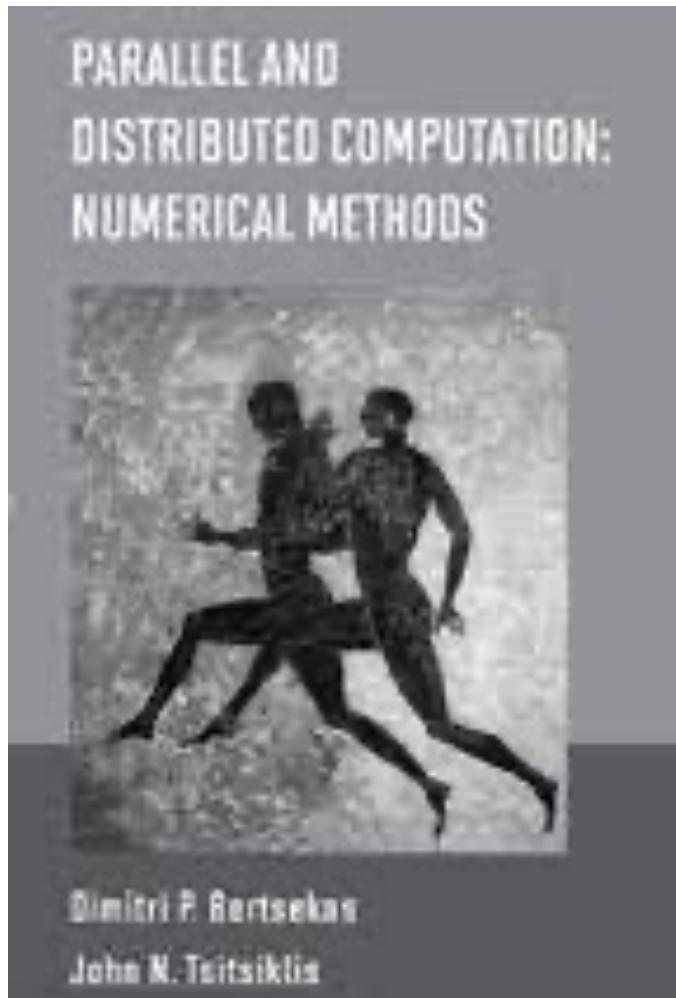
Non-Linear Min-Cost-Flow

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

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 $i \in \mathcal{V}$.



Non-Linear Min-Cost-Flow



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

Kirchhoff's Current Law

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

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

Generalized Ohm Law

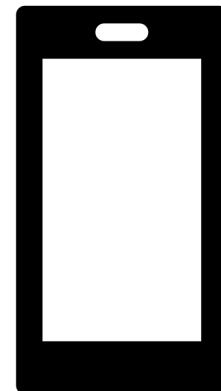
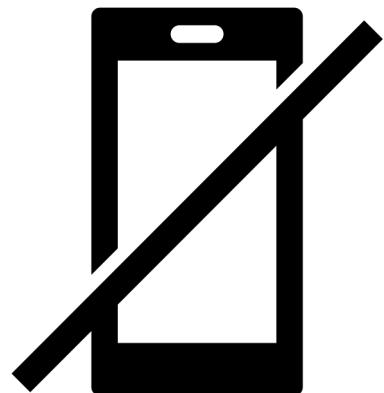
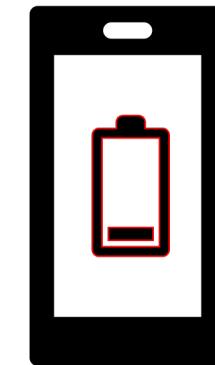
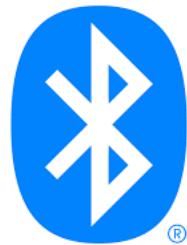
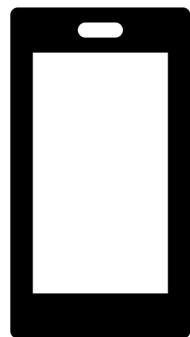
- GTVMin as NFL Principle
- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

Computational Aspects.

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

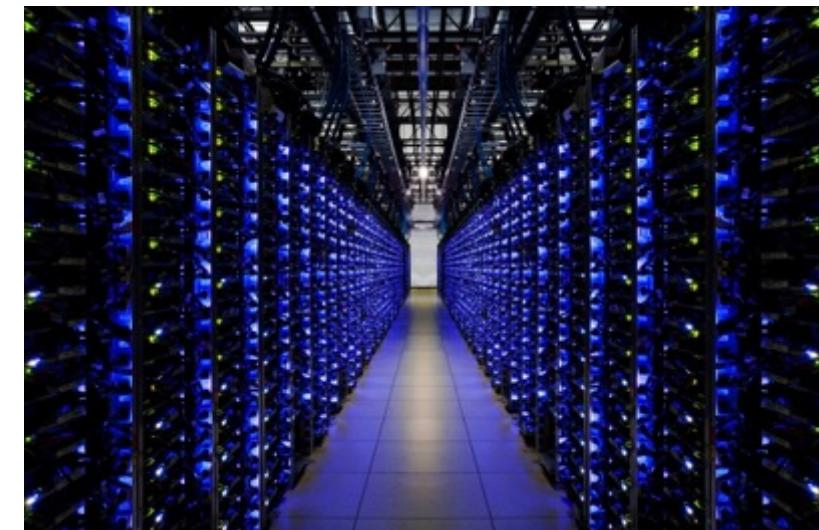
- solve in ad-hoc nets of low-cost devices
- robustness against node/link failures
- robustness against “stragglers”

Our Toy NFL Setting



Another NFL Setting...

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



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

Two Main Flavours

- Primal (Gradient) Methods
- Primal-Dual Methods

Primal (Gradient) Methods



Gradient Descent

$$\min_w \sum_{i \in M} L^{(i)}(w^{(i)}) + \lambda \sum_{\{i,j\}} A_{i,j} \phi(w^{(i)} - w^{(j)})$$

$f(w)$

optimality condition $\nabla f(w) = 0$

$$w^{(k+1)} = w^{(k)} - \alpha^{(k)} \nabla f(w^{(k)})$$

Subgradient Descent (SGD)

$$\min_w \sum_{i \in M} L^{(i)}(w^{(i)}) + \lambda \sum_{\{i,j\}} A_{i,j} \phi(w^{(i)} - w^{(j)})$$

$f(w)$

optimality condition $0 \in \partial f(w)$

$$w^{(k+1)} = w^{(k)} - \alpha^{(k)} g^{(k)} \quad g^{(k)} \in \partial f(w^{(k)})$$

Distributed SGD

A. Nedić and A. Olshevsky, "Distributed Optimization Over Time-Varying Directed Graphs," in *IEEE Transactions on Automatic Control*, 2015,

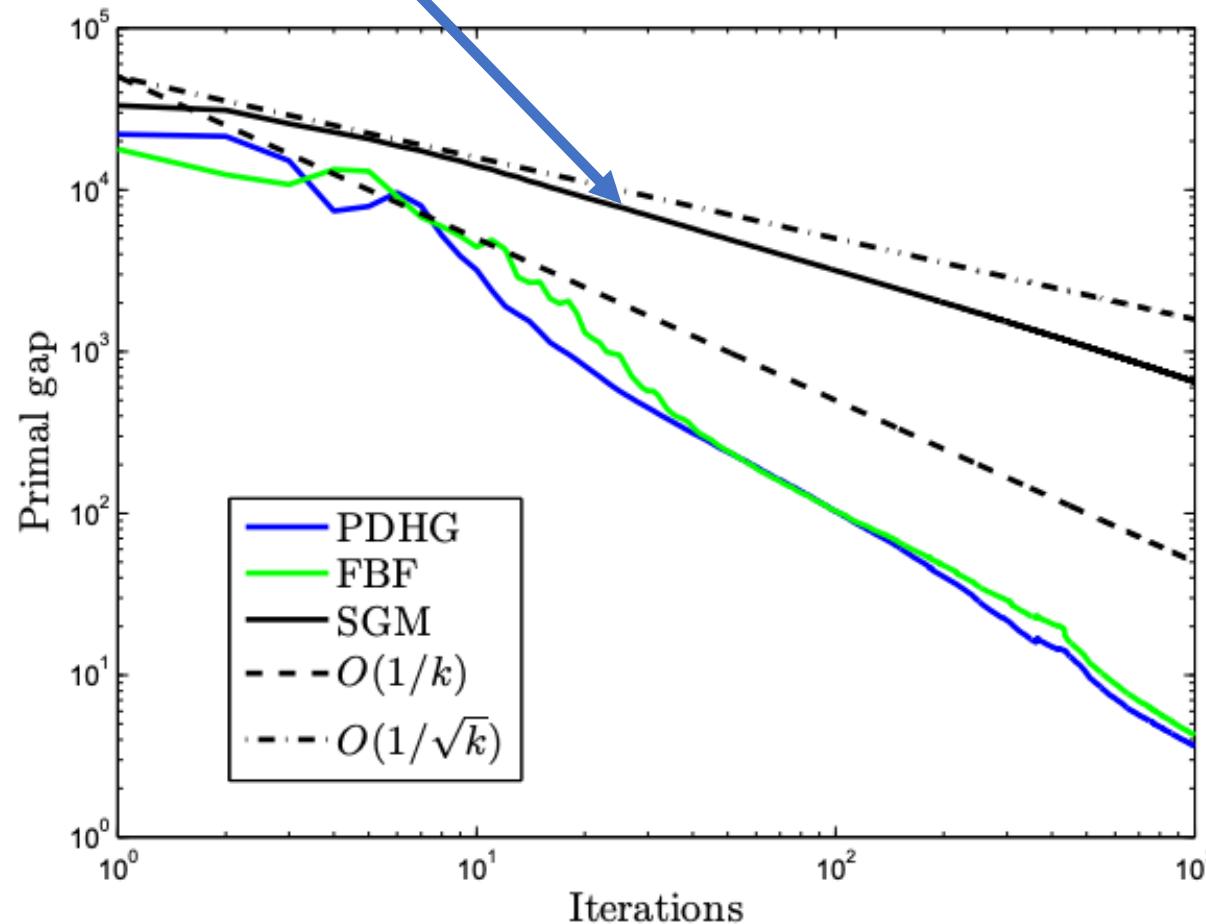


A. Nedic (M.S., University of Belgrade, 1991)

A. Nedić and A. Olshevsky, "Stochastic Gradient-Push for Strongly Convex Functions on Time-Varying Directed Graphs," in *IEEE Transactions on Automatic Control*, 2016,

A. Nedic and A. Ozdaglar, "Distributed Subgradient Methods for Multi-Agent Optimization," in *IEEE Transactions on Automatic Control*, Jan. 2009.

SGD Requires Many Iter.



Complexity of SGD

Theorem 3. Let $L_\ell, R > 0$ and $\gamma \in (0, 1]$. There exists a matrix W of eigengap $\gamma(W) = \gamma$, and n functions f_i satisfying (A2), where n is the size of W , such that for all $t < \frac{d-2}{2} \min(\tau/\sqrt{\gamma}, 1)$ and all $i \in \{1, \dots, n\}$,

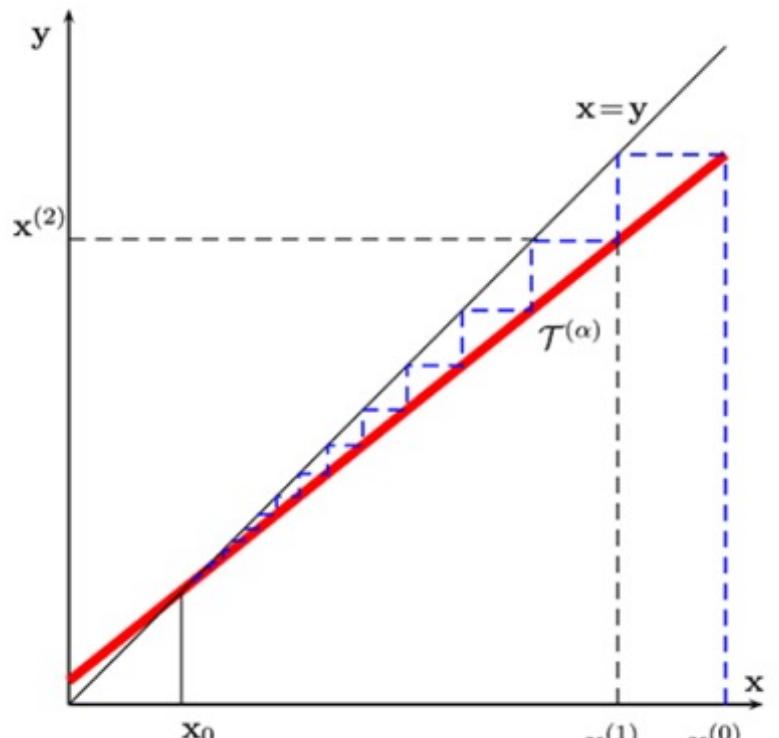
$$\bar{f}(\theta_{i,t}) - \min_{\theta \in B_2(R)} \bar{f}(\theta) \geq \frac{RL_\ell}{108} \sqrt{\frac{1}{(1 + \frac{2t\sqrt{\gamma}}{\tau})^2} + \frac{1}{1+t}}. \quad (19)$$

K. Scaman, F. Bach, S. Bubeck, L. Massoulié, Y Lee, Optimal Algorithms for Non-Smooth Distributed Optimization in Networks, NeurIPS 2018.

SGD as Fixed Point Iteration

$$w^{(k+1)} = \mathcal{T}^{(k)}(w^{(k)})$$

with $\mathcal{T}^{(k)}(w^{(k)}) = w^{(k)} - \alpha^{(k)} \partial f(w^{(k)})$



AJ, “A Fixed-Point of View on Gradient Methods for Big Data”, Front. Appl. Math. Stat., 2017.

Plenary on Fixed-Point Tools

$$w^{(k+1)} = \mathcal{T}^{(k)}(w^{(k)})$$



Jean-Christophe
Pesquet

Jean-Christophe Pesquet (Il in 1987, the Ph.D. and HDE in 1999, he was a Maître de C University Paris-Est, and fro the university. He is curre Director of the CVN (Inria te 2021. In 2005, J.-C. Pesque was a member of the SPTM IEEE SPL (2004-2006). He journal (2010-2015), and a r now an associate editor of methods in data science.

Fixed Point Strategies in Signal and Image Processing

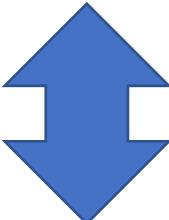
Primal-Dual Methods



Primal-Dual Optimality Conditions.

(assuming convexity of loss functions and GTV penalty)

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



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

this is again a fixed-point problem !

Proximal Point Algorithm.

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

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

solve iteratively by proximal point algorithm

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

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

After Some Manipulations.

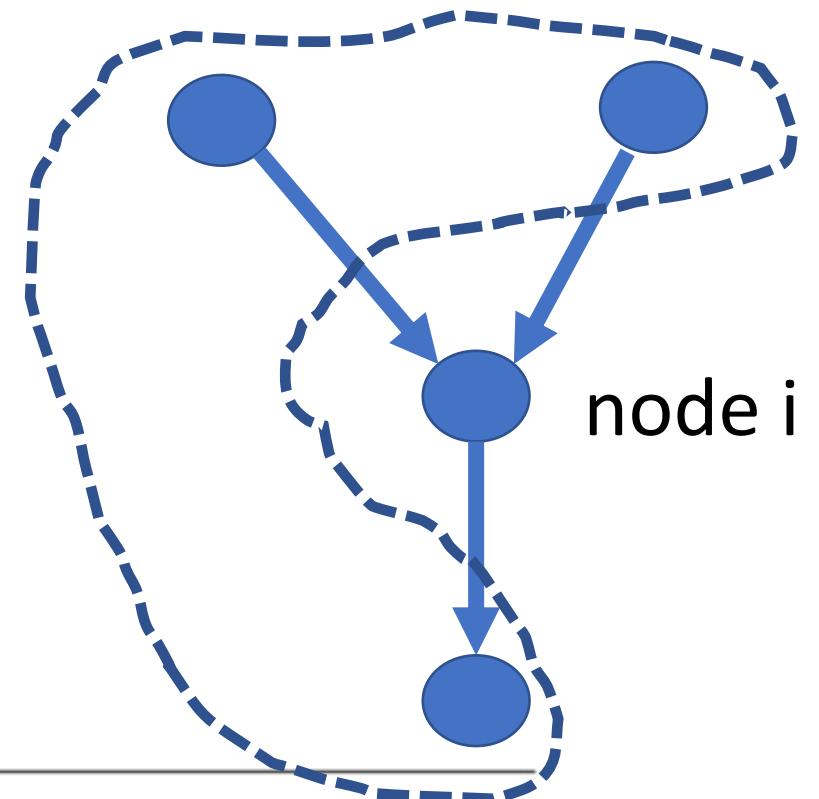
Algorithm 1 Primal-Dual Method for Networked FL

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

Initialize: $k := 0$; $\hat{\mathbf{w}}_0 := \mathbf{0}$; $\hat{\mathbf{u}}_0 := \mathbf{0}$; $\sigma_e = 1/2$ and $\tau_i = 1/|\mathcal{N}_i|$

```

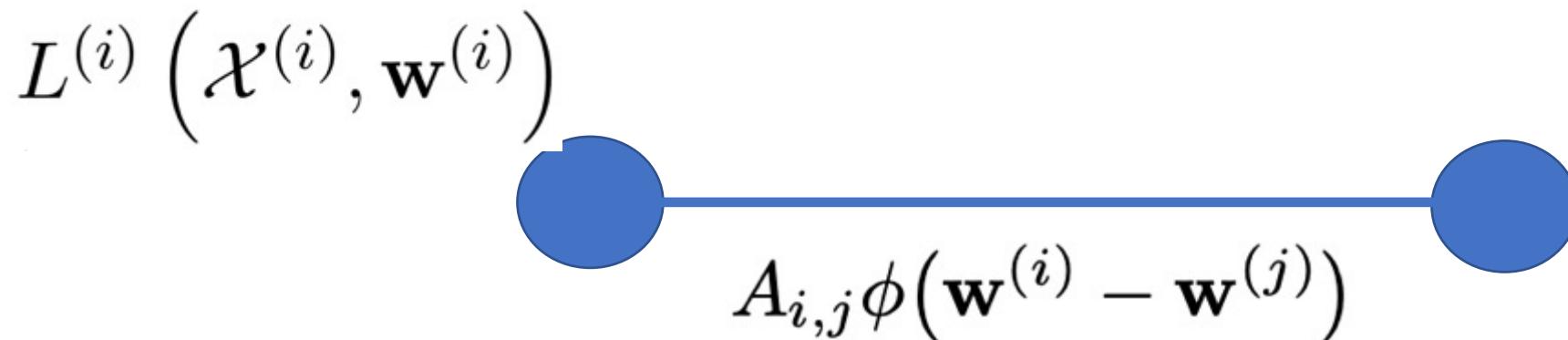
1: while stopping criterion is not satisfied do
2:   for all nodes  $i \in \mathcal{V}$  do
3:      $\hat{\mathbf{w}}_{k+1}^{(i)} := \hat{\mathbf{w}}_k^{(i)} - \tau_i \sum_{e \in \mathcal{E}} D_{e,i} \hat{\mathbf{u}}_k^{(e)}$ 
4:   end for
5:   for nodes in the training set  $i \in \mathcal{M}$  do
6: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



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.



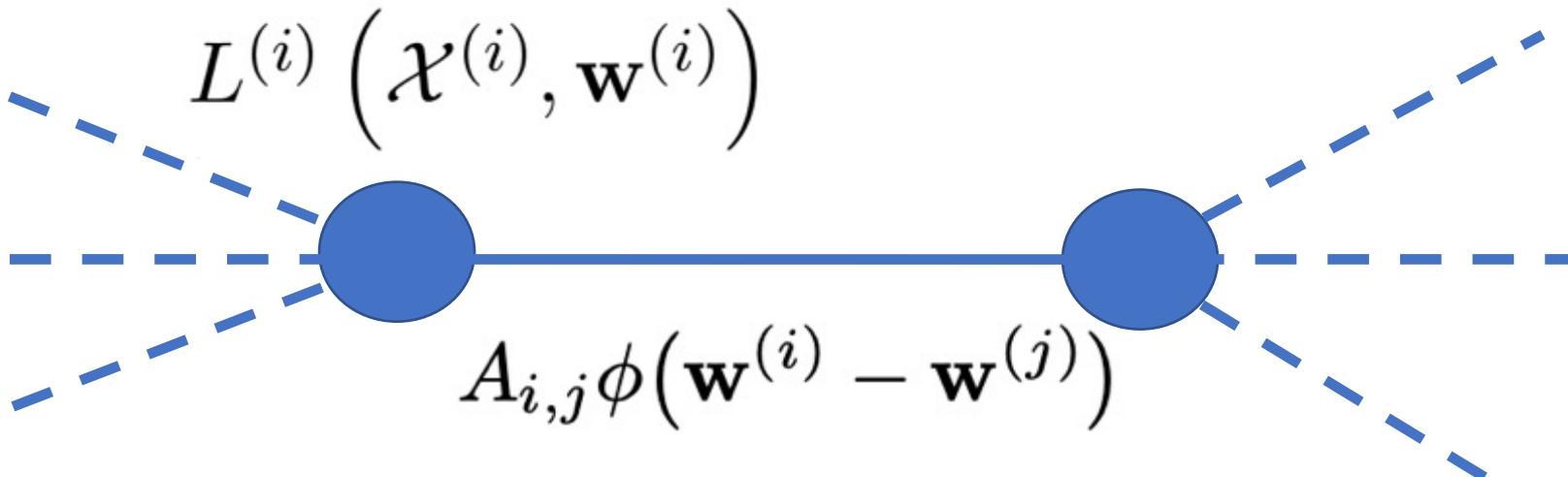
node-wise

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

edge-wise

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

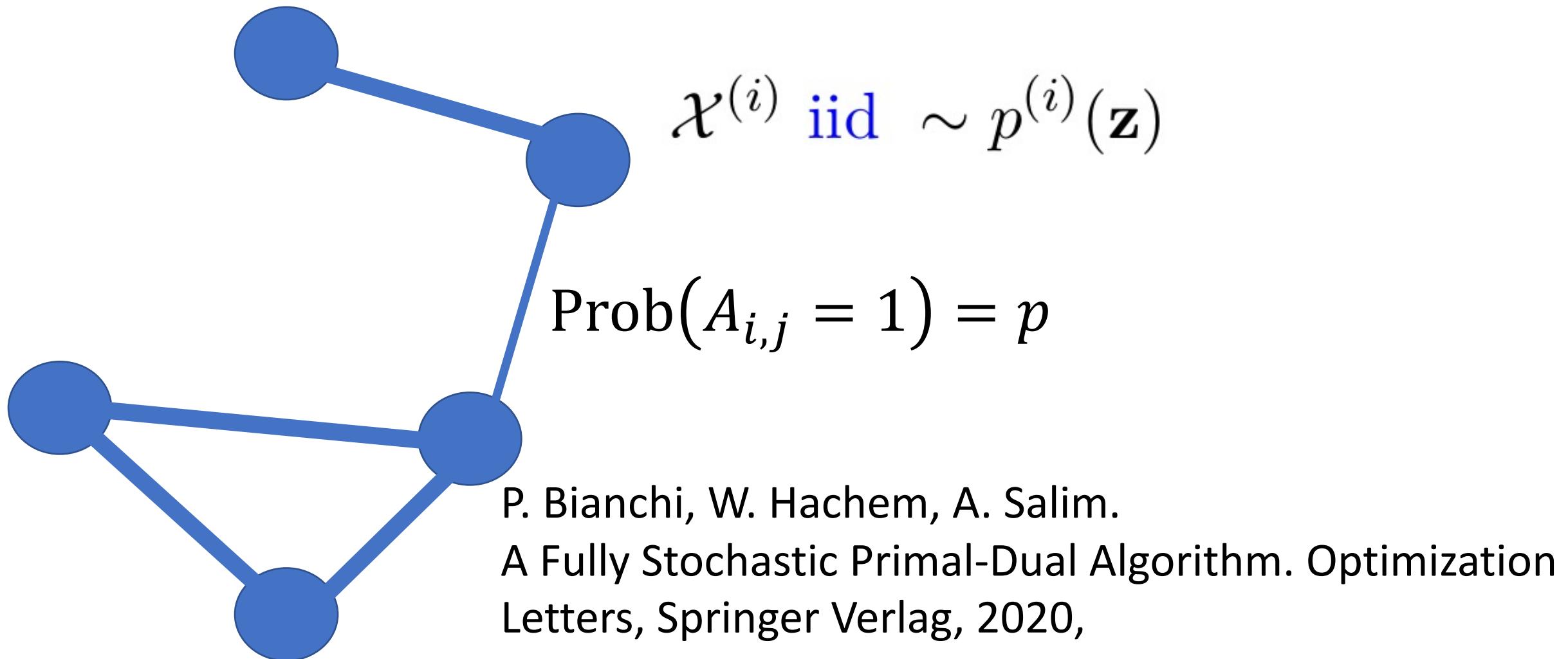
Spreading Local Results.



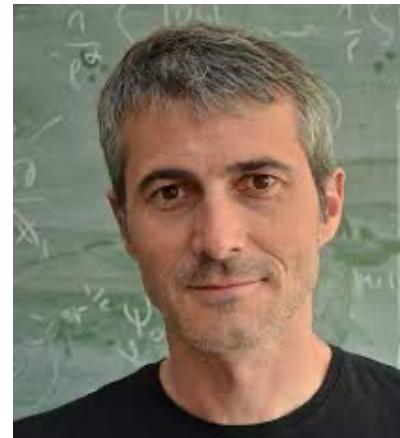
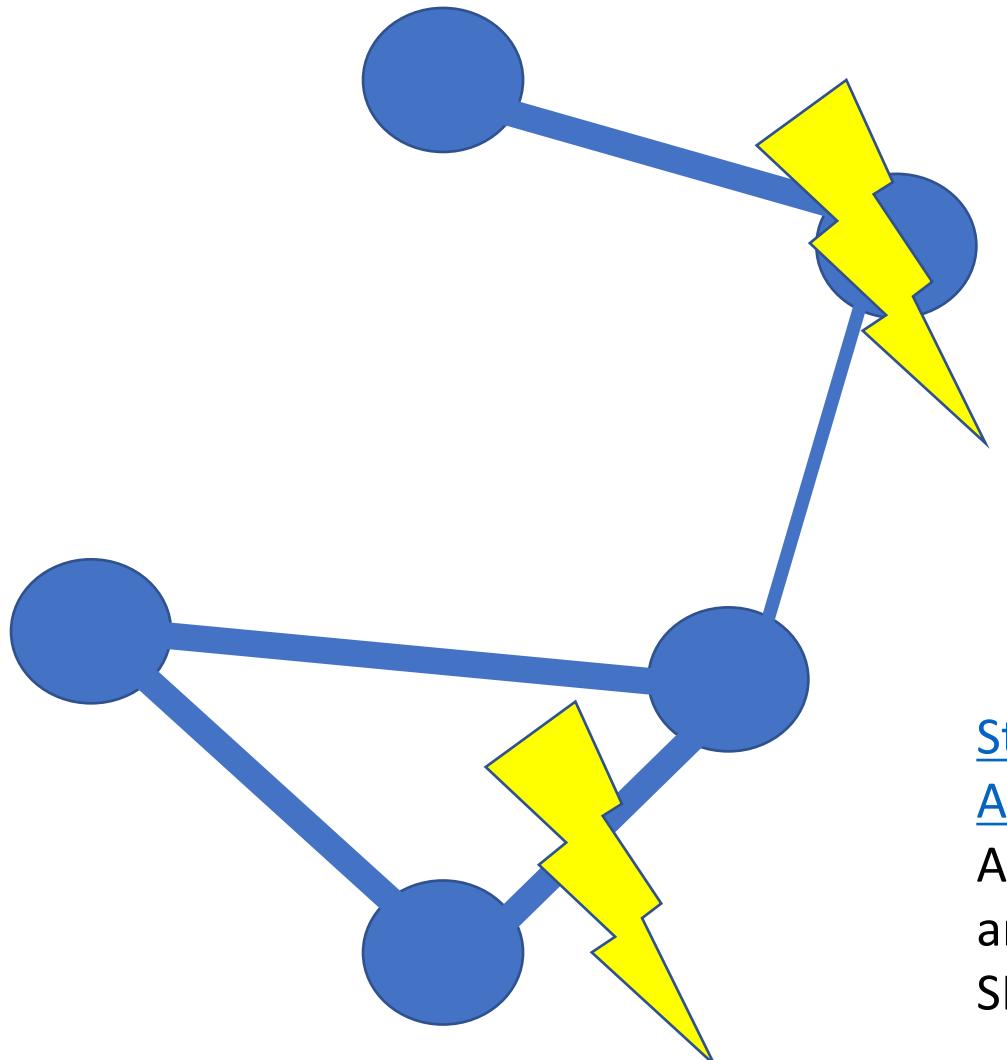
```
2:   for all nodes  $i \in \mathcal{V}$  do
3:      $\hat{\mathbf{w}}_{k+1}^{(i)} := \hat{\mathbf{w}}_k^{(i)} - \tau_i \sum_{e \in \mathcal{E}} D_{e,i} \hat{\mathbf{u}}_k^{(e)}$ 
4:   end for
```

```
8:   for all edges  $e \in \mathcal{E}$  do
9:      $\hat{\mathbf{u}}_{k+1}^{(e)} := \hat{\mathbf{u}}_k^{(e)} + \sigma_e (2(\hat{\mathbf{w}}_{k+1}^{(e_+)} - \hat{\mathbf{w}}_{k+1}^{(e_-)}) - (\hat{\mathbf{w}}_k^{(e_+)} - \hat{\mathbf{w}}_k^{(e_-)}))$ 
```

Probabilistic Networked Data



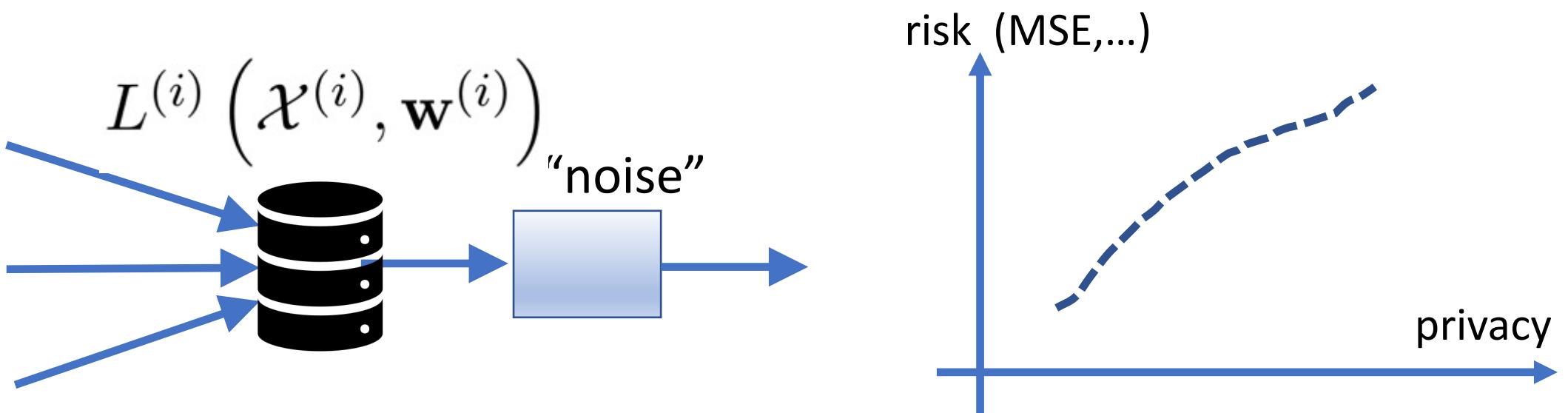
Random Node/Link Failures.



[Stochastic Primal-Dual Hybrid Gradient Algorithm with
Arbitrary Sampling and Imaging Applications](#)

Antonin Chambolle, Matthias J. Ehrhardt, Peter Richtárik,
and Carola-Bibiane Schönlieb
SIAM Journal on Optimization 2018 28:4, 2783-2808

Privacy-Preservation.



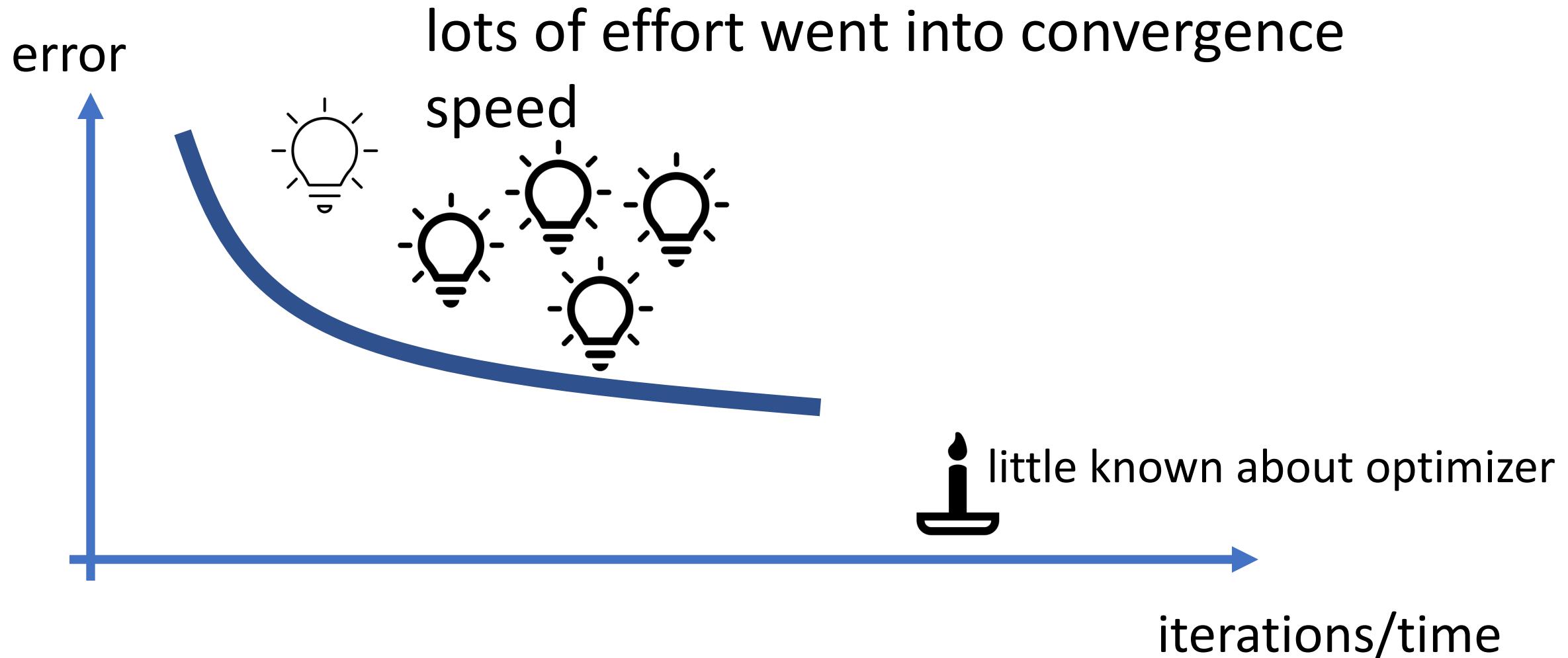
- F. Shang, T. Xu, Y. Liu, H. Liu, L. Shen and M. Gong, "Differentially Private ADMM Algorithms for Machine Learning," in *IEEE Transactions on Information Forensics and Security*, vol. 16, pp. 4733-4745, 2021, doi: 10.1109/TIFS.2021.3113768.
- J. C. Duchi, M. I. Jordan, and M. J. Wainwright, "Local privacy and statistical minimax rates," in Proc. IEEE Annu. Symp. Found. Comput. Sci., pp. 429–438, 2013.

Bottom Line.

PD method solves GTVMin in distributed,
robust and **privacy-friendly way**

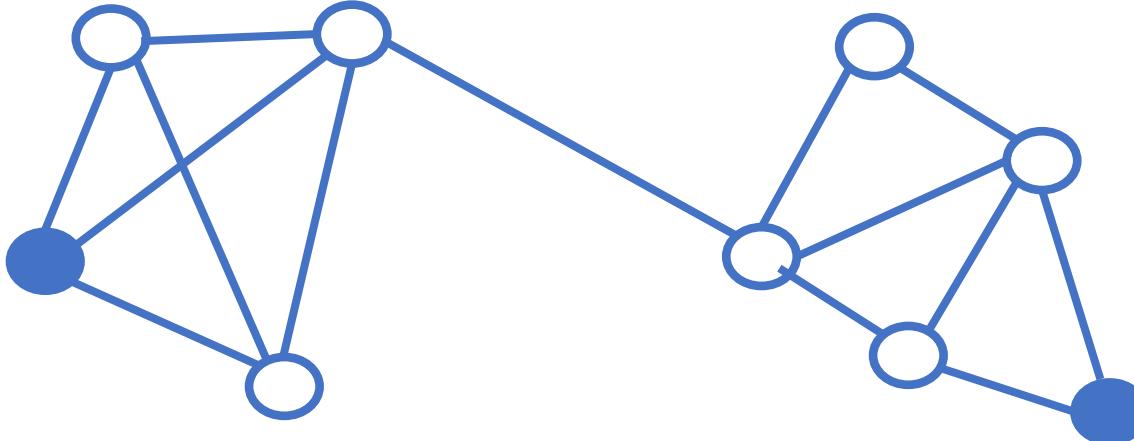
....., however

Compute vs. Accuracy



Are GTVMin Solutions Any Good?

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



training/sampling
set \mathcal{M}

which combination of signal model (choice of ϕ) and sampling set \mathcal{M} ensure solutions of GTVMin are “sensible” ?

- GTVMin as NFL Principle
- The Dual of GTVMin
- Interpretations
- Computational Aspects
- Statistical Aspects

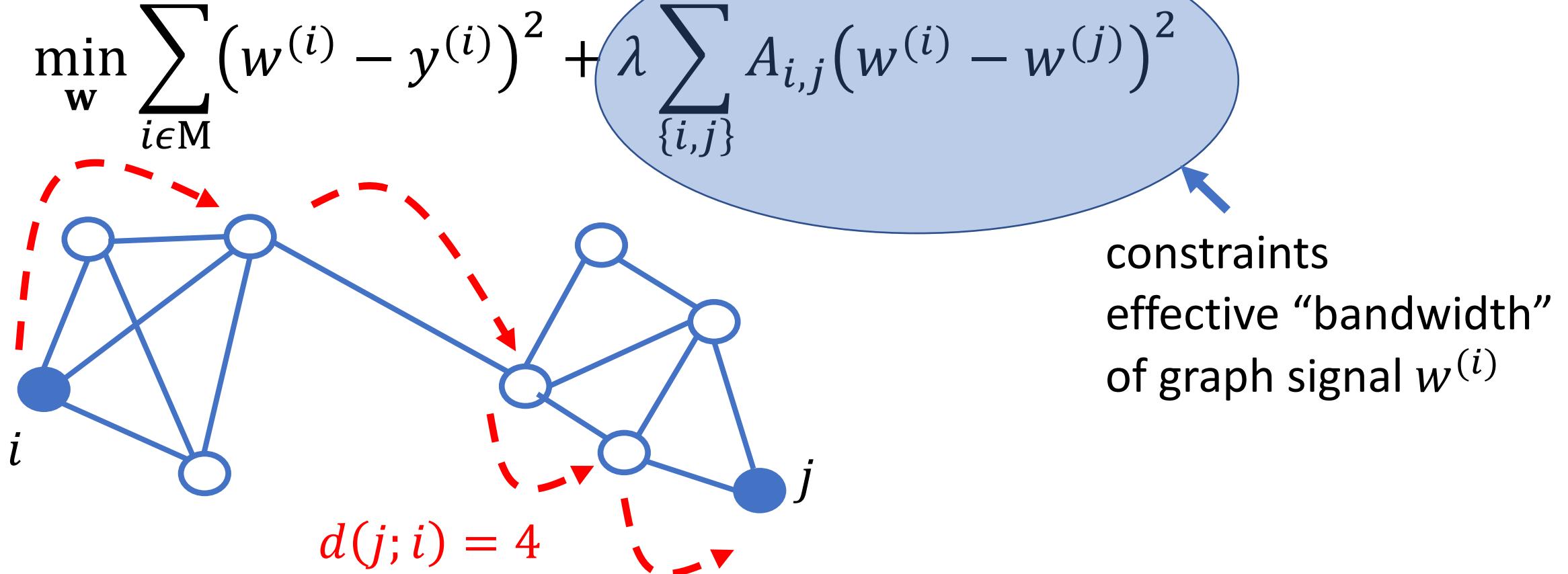
Statistical Aspects of GTVMin

$$\min_w \sum_{i \in M} L^{(i)}(w^{(i)}) + \lambda \sum_{\{i,j\}} A_{i,j} \phi(w^{(i)} - w^{(j)})$$

statistical properties of GTVMin solutions?

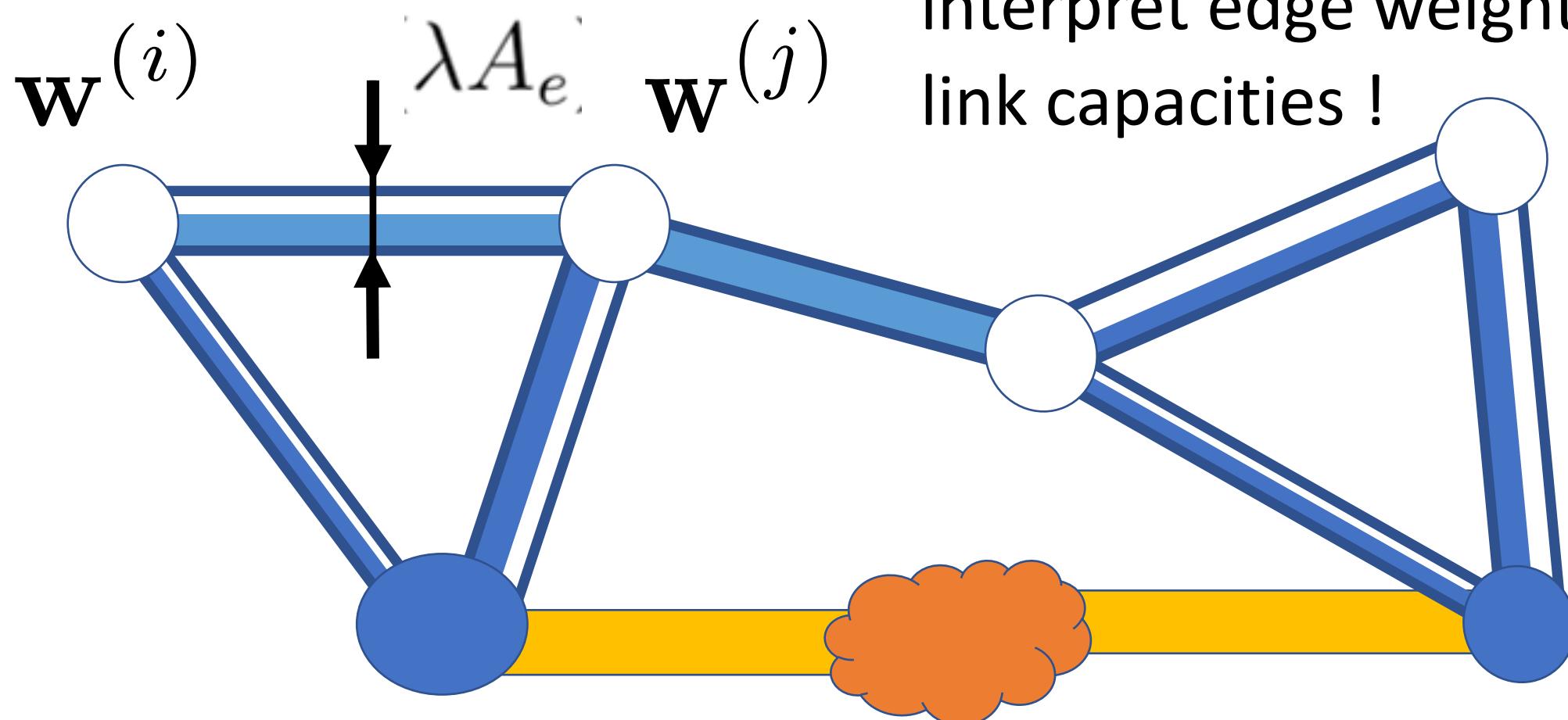
- sampling theorems (**signal processing**)
- generalization bounds (**ML**)

Signal Processing Perspective.



M. Tsitsvero, S. Barbarossa and P. Di Lorenzo, "Signals on Graphs: Uncertainty Principle and Sampling," in *IEEE Transactions on Signal Processing*, 2016,

Our Perspective: Flows.



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

Why Flows?

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

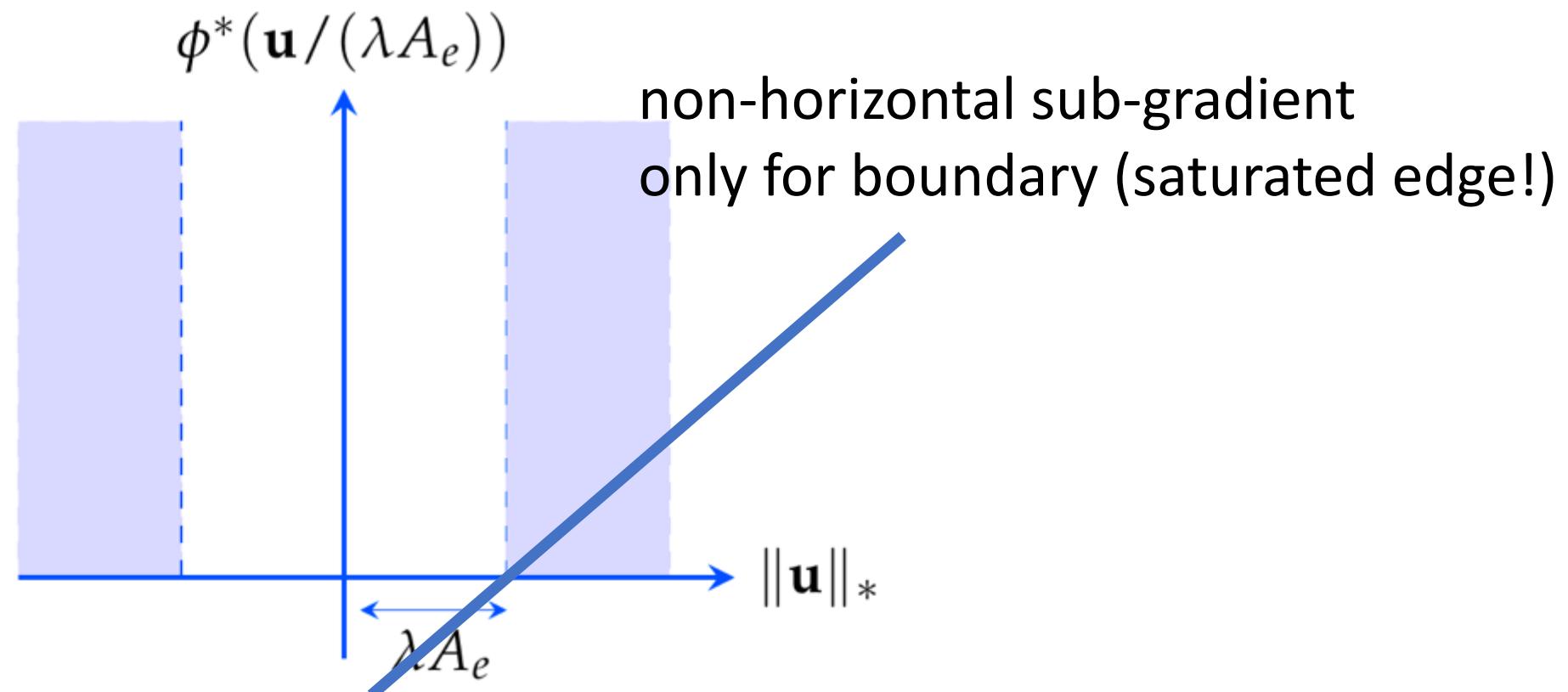
$$\widehat{\mathbf{w}}^{(e_+)} - \widehat{\mathbf{w}}^{(e_-)} \in (\lambda A_e) \partial \phi^*(\widehat{\mathbf{u}}^{(e)} / (\lambda A_e)) \text{ for every edge } e \in \mathcal{E}.$$

the solutions of GTVMin is a flow $\widehat{\mathbf{u}}^{(e)}$,

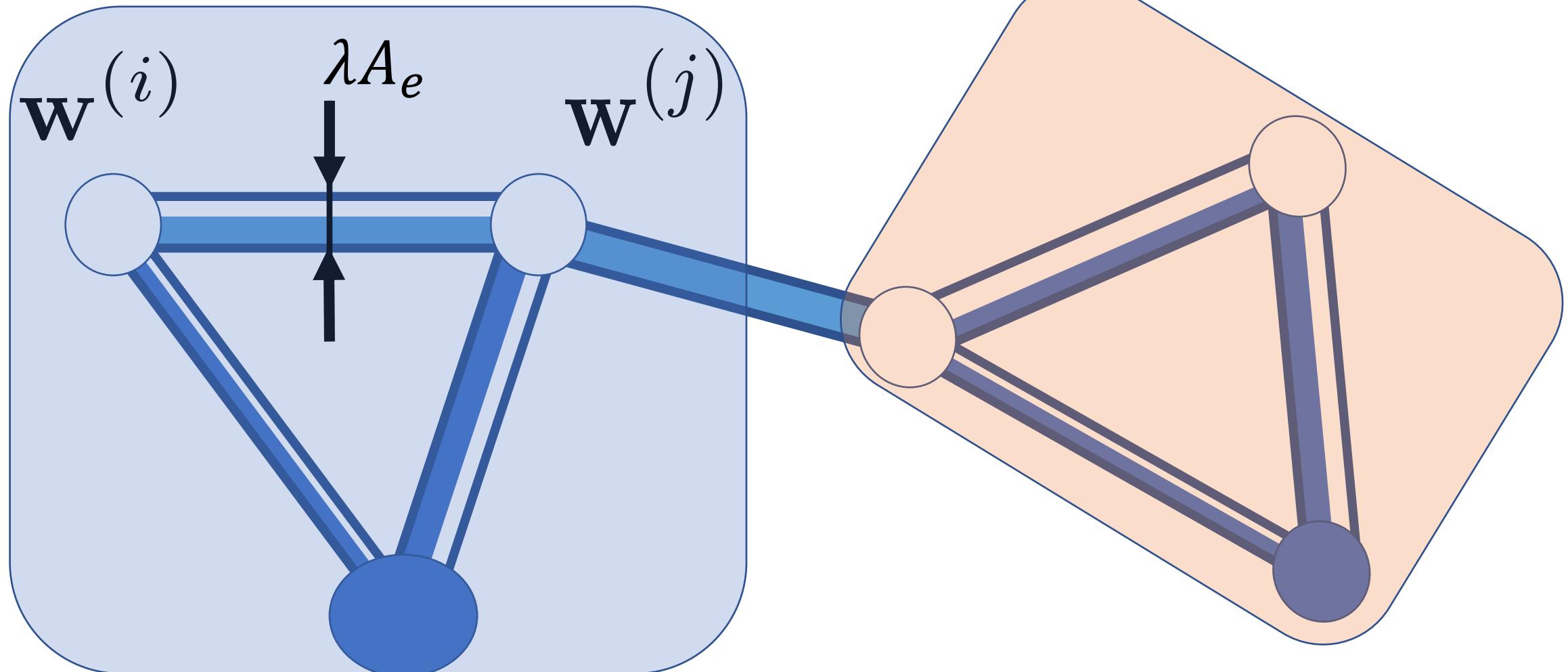
properties of the flow coupled with properties of GTVMin solution!

Primal-Dual Witness

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

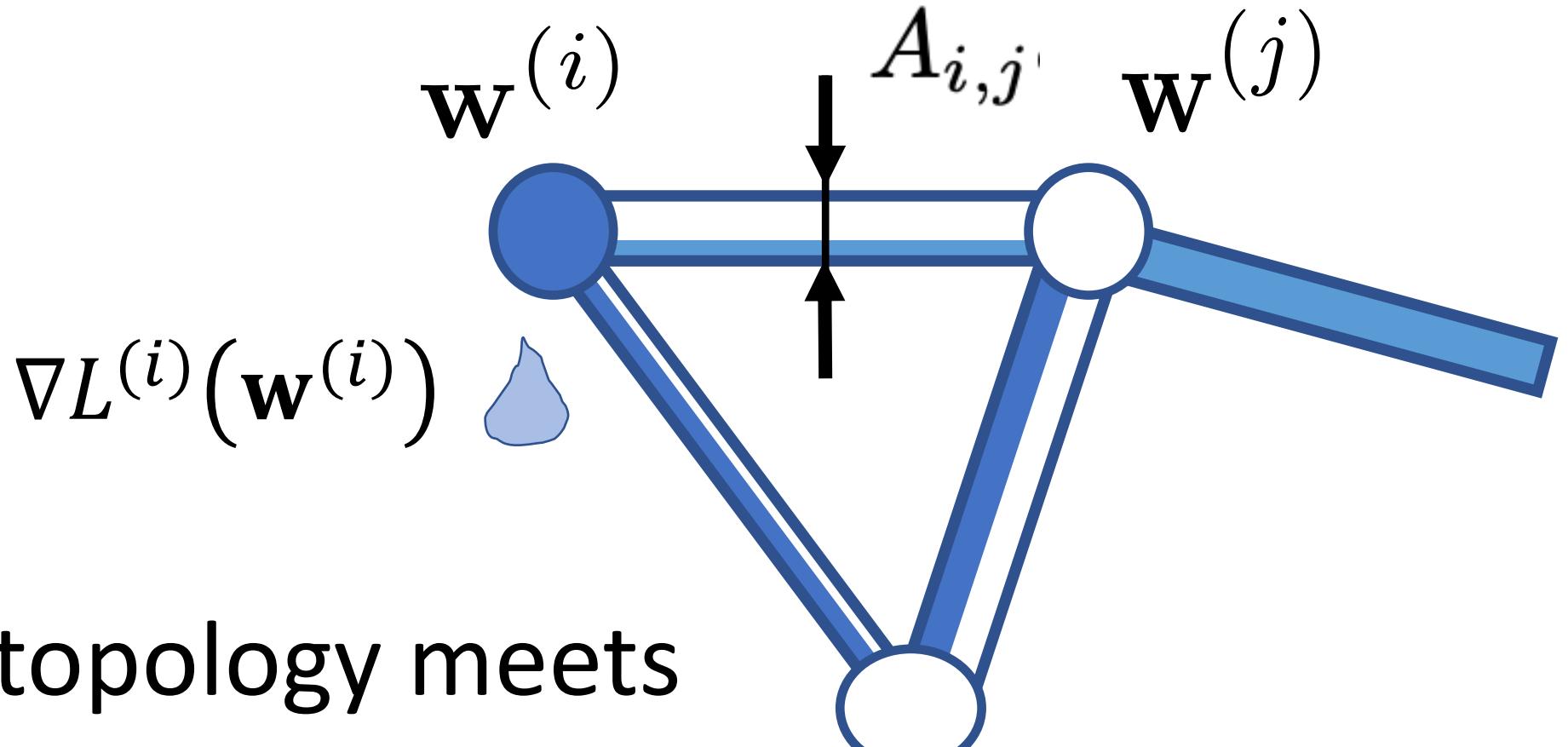


Cluster-wise Pooling.



parameter vectors can only change over saturated links !

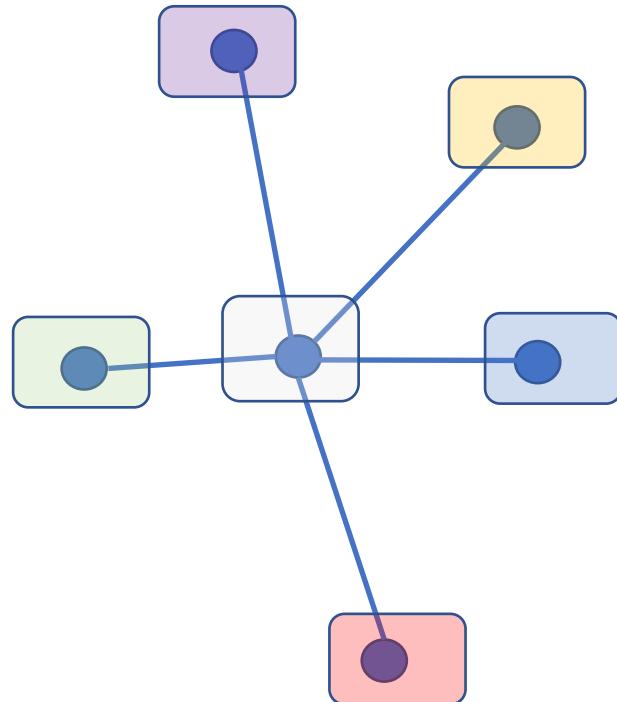
Leaky Training Set.



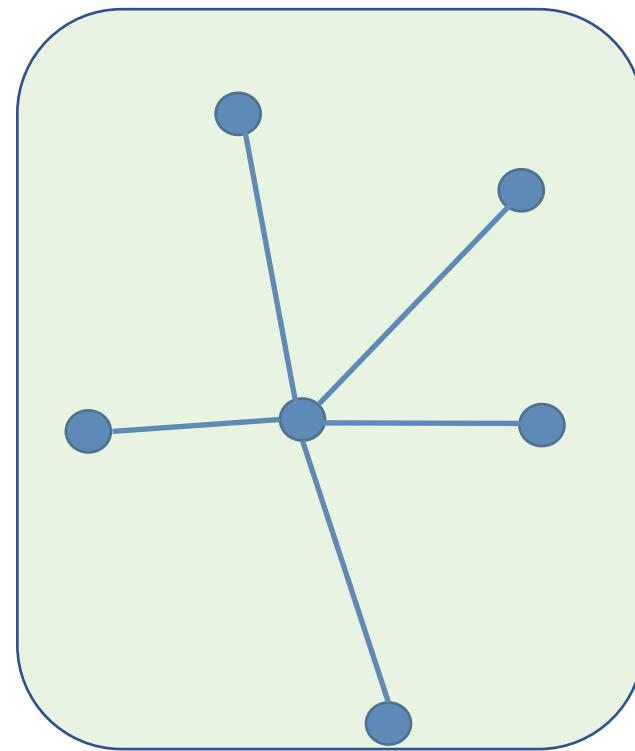
network topology meets
geometry of loss functions !

Personalization vs. Globalization

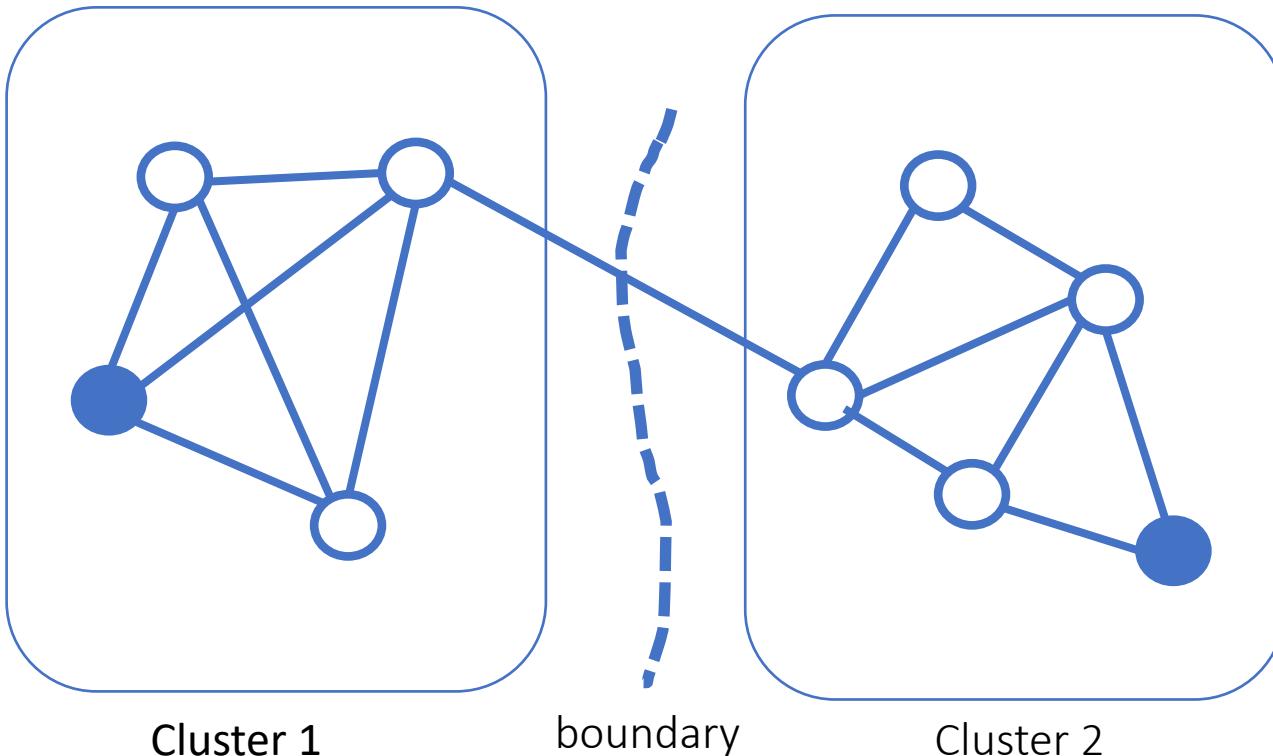
small λ , edges easily saturated



large λ , edges hard to saturate

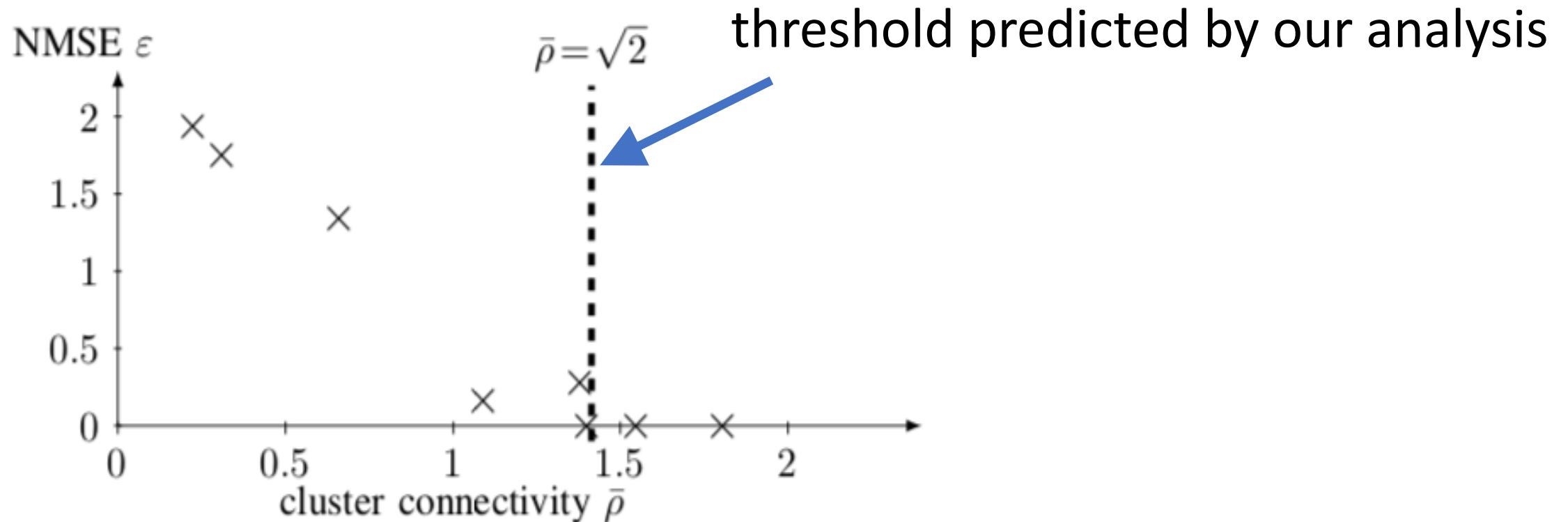


Define Cluster by Boundary Flow.



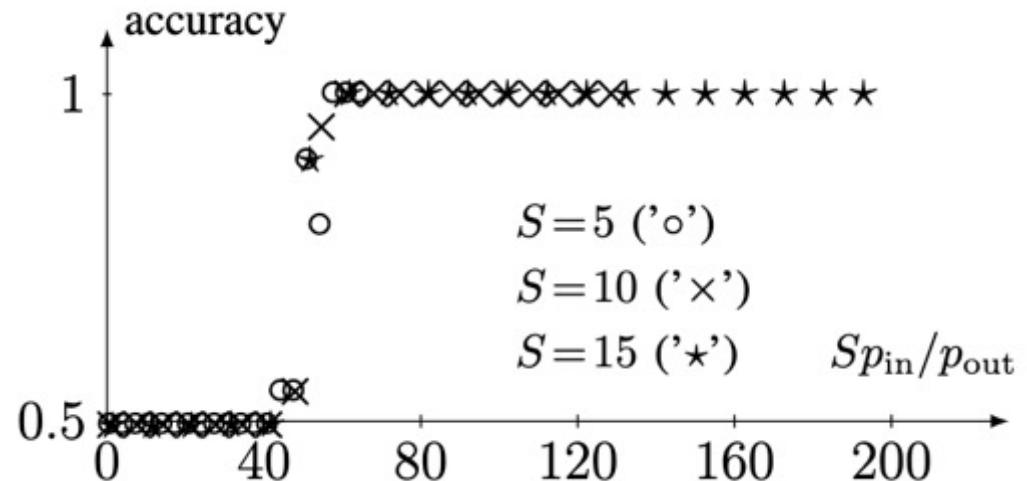
connectivity measured by
flow ρ that can be routed
over boundary edge

Statistical Error vs. Connectivity.



A. Jung and N. Tran, "Localized Linear Regression in Networked Data," in *IEEE Signal Processing Letters*, vol. 26, no. 7, pp. 1090-1094, July 2019.

Clustering Assumption in SBM.



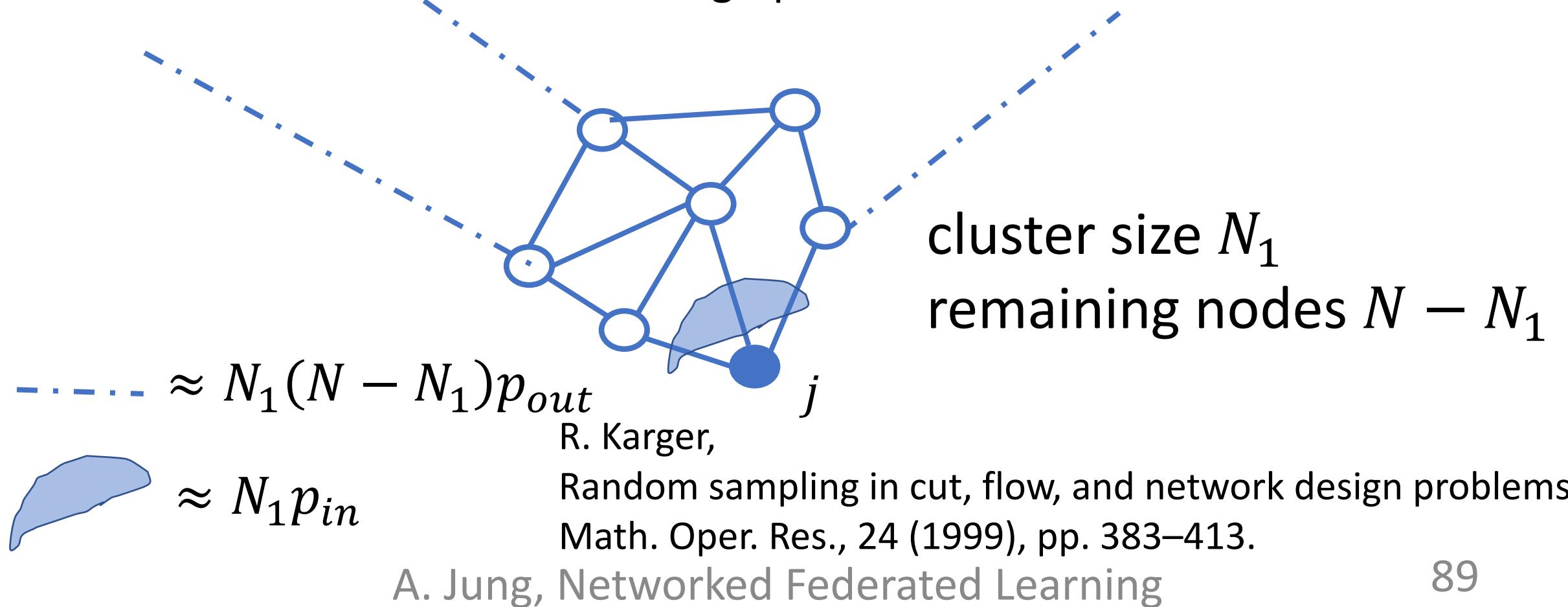
- intra-cluster edge prob p_{in}
- inter-cluster edge prob p_{out}
- S training nodes in each cluster
- critical value for S^*p_{in}/p_{out}

A. Jung,

"Clustering in Partially Labeled Stochastic Block Models via Total Variation Minimization,"
54th Asilomar Conference on Signals, Systems, and Computers, 2020,

Mathematical Devices.

- flow conservation/Hoffman's circulation theorem
- concentration of cuts in random graphs



Wrap Up.

- GTVMin paradigm for NFL
- dual of GTVMin = non-lin. minimum-cost flow
- solve GTVMin. with **primal-dual method**
- **scalable and robust** message passing
- GTV min. adaptively pools similar datasets

Thank you for
your attention!

Interested in a Phd in my group?

<https://www.visitfinland.com/en/>

happiest country in the world

All Images News Maps Videos More

About 26.100.000 results (1,15 seconds)

Finland

1. **Finland.** For the fifth year in a row, Finland is number one when it comes to happiness. 31 Mar 2022

