

Distributed Collaborative Non-Convex Optimization
A Case Study on ℓ_0 Constrained Graph Estimation

Mengdi Wang

ORFE, Princeton University

SIAM CSE, March 2015



Collaborators



Ethan X. Fang



Han Liu

Distributed Collaborative Optimization

Consider the social welfare maximization problem

$$\max_{x_i \in \mathcal{X}_i, \forall i} u_0 \left(\frac{1}{N} \sum_{i=1}^N c_i(x_i) \right) + \frac{1}{N} \sum_{i=1}^N u_i(x_i)$$

- $u_i(x_i)$, $i = 1, \dots, N$, is the private utility of the i th user
- u_0 is the social utility associated with some common goods (e.g., total energy consumption, total air pollution)
- u_i, c_i, \mathcal{X}_i are **nonconvex**

Many decentralized problems have a natural **near-separable** structure.

Question:

Is there any intrinsic simplicity that we can leverage to design efficient distributed algorithms?

Application: Pricing of Public Resources

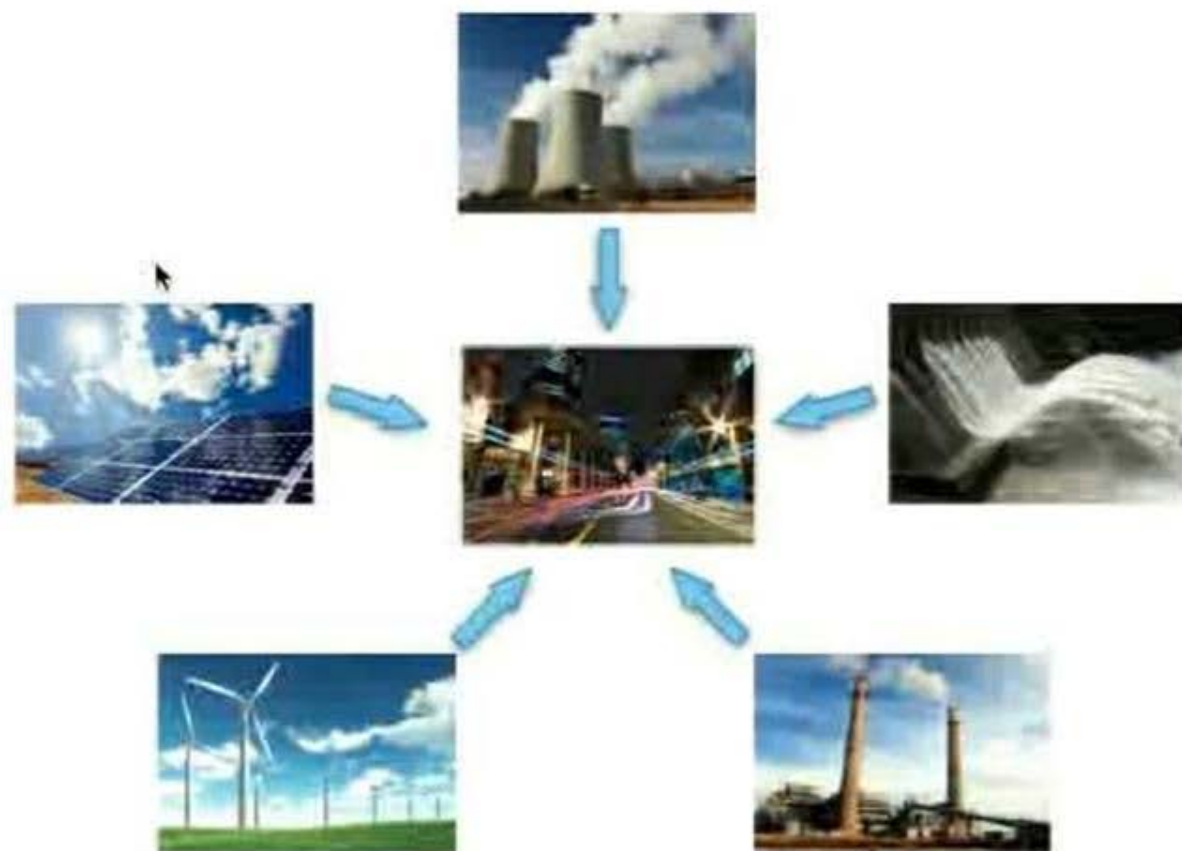


Figure : Individual users may have flexibility in selecting different energy usages, resulting in non-concave preference function.

Application: Graphical Models

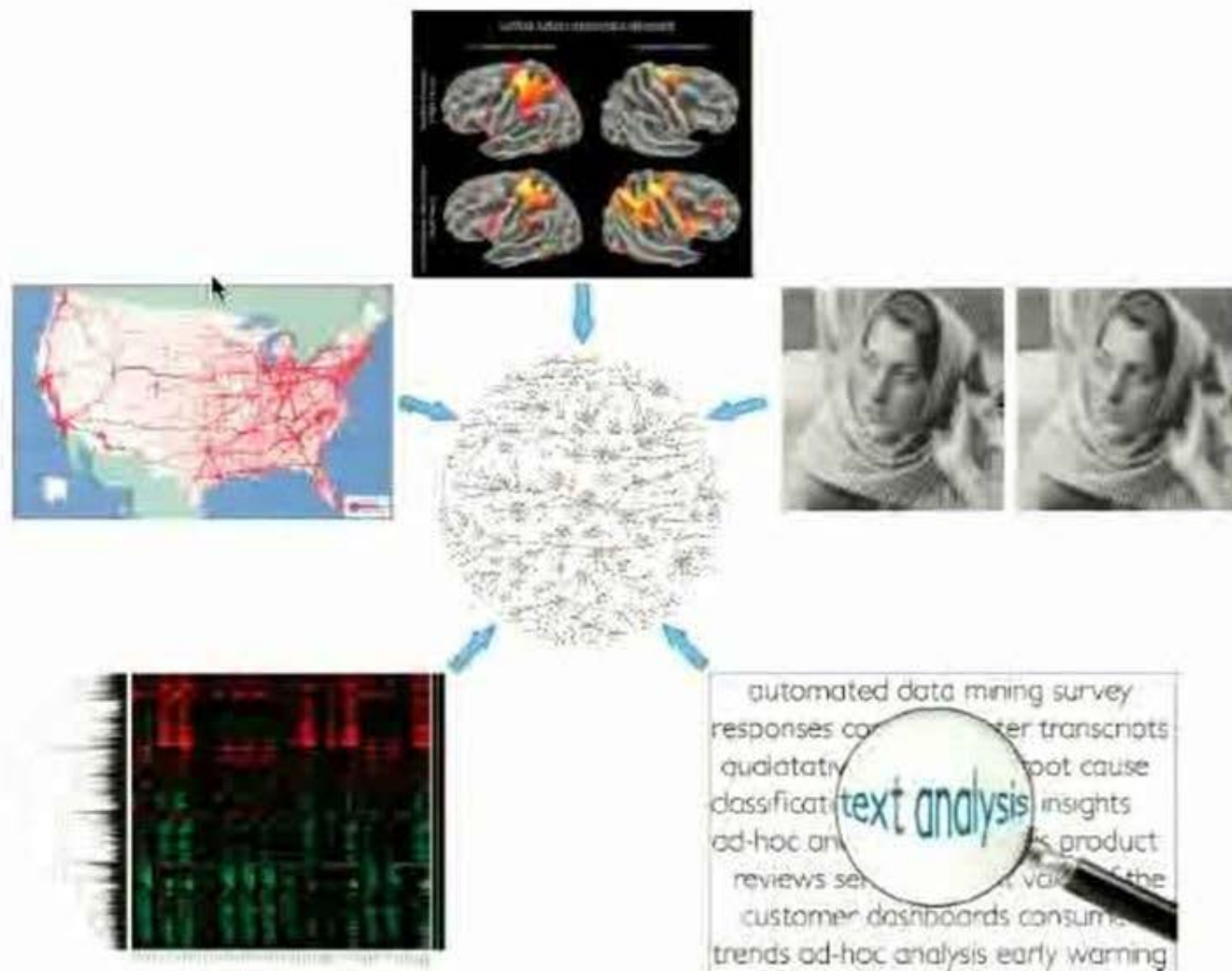


Figure 1.1 Graphical model is used in a wide range of applications.

Graphical Model

In the graph, each node represents a random variable. An edge is presented between X_j and X_k if and only if they are **conditionally dependent** given other nodes, i.e.,

$$\mathbb{P}(X_j = x_j, X_k = x_k | \mathbf{X} \setminus \{j,k\}) = \mathbb{P}(X_j = x_j | \mathbf{X} \setminus \{j,k\}) \mathbb{P}(X_k = x_k | \mathbf{X} \setminus \{j,k\}).$$

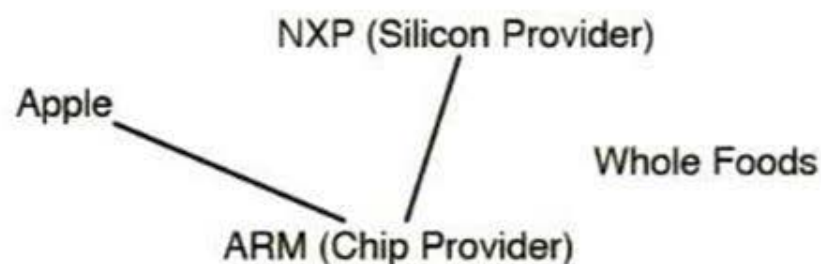
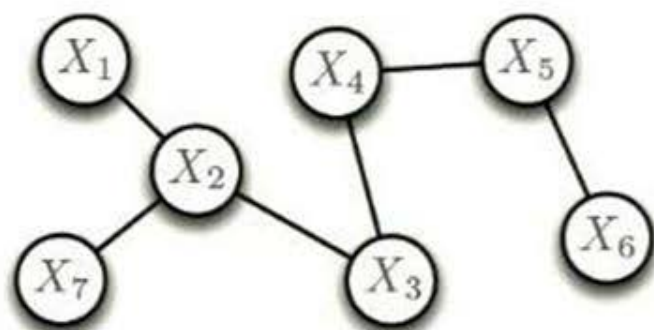
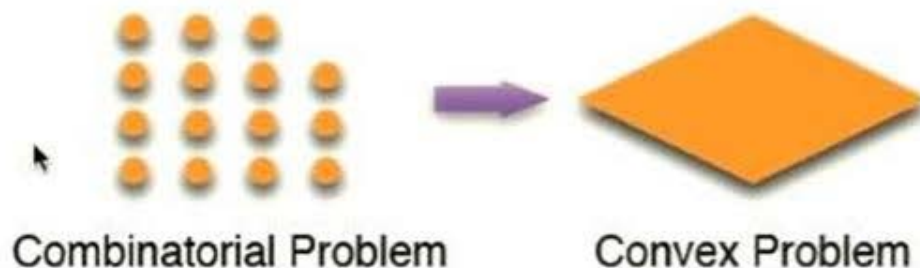


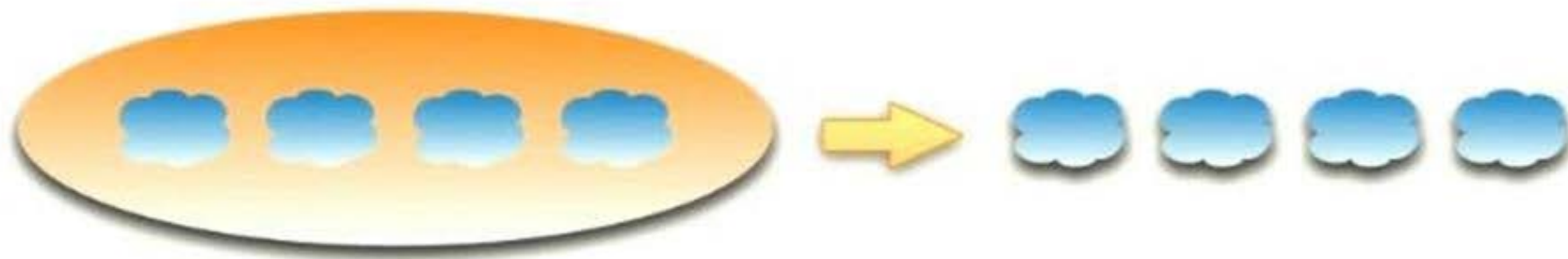
Figure : A class of graphical models explains the conditional independences.

Computational Challenge

This problem is combinatorial and NP-hard.



Convex relaxation would incur some **unavoidable** statistical errors. In practice, each subproblem can be easy to solve. However, the whole problem is difficult due to the global sparsity constraint.



Sparse Graph Estimation

$$\min_{\beta_1, \dots, \beta_d} \frac{1}{d} \sum_{j=1}^d \ell_j(\beta_j), \text{ subject to } \sum_{j=1}^d \|\beta_j\|_0 \leq K.$$

Why ℓ_0 ?

We care about estimation error as well as computation efficiency.

- The ℓ_0 constraint exactly describes the statistic modeling assumptions:

“The total number of arcs is bounded by $K \ll d^2$.”

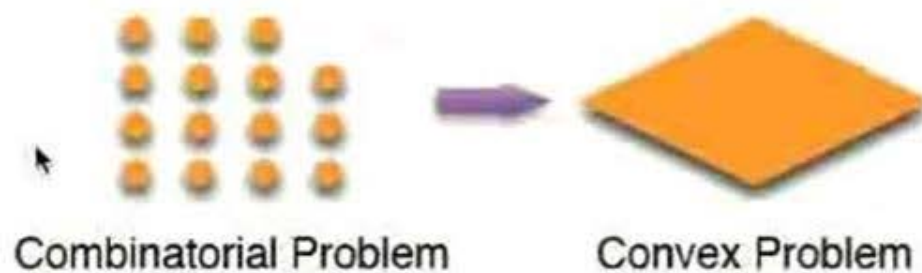
- Avoids the statistical estimation gap between ℓ_0 and ℓ_1 norm

Objective: Estimate the True Graph

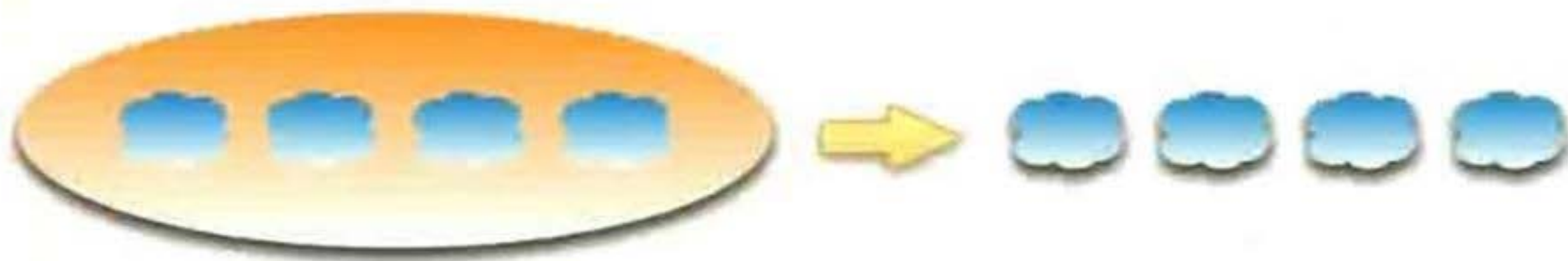
- Let β^* be the *true* coefficients (not the optimal solution to the optimization problem)
- We want to find an (approximate) solution $\hat{\beta}$ to the optimization problem and use it as an estimator of β^*

Computational Challenge

This problem is combinatorial and NP-hard.



Convex relaxation would incur some **unavoidable** statistical errors.
In practice, each subproblem can be easy to solve. However, the whole problem is difficult due to the global sparsity constraint.



Sparse Graph Estimation

$$\min_{\beta_1, \dots, \beta_d} \frac{1}{d} \sum_{j=1}^d \ell_j(\beta_j), \text{ subject to } \sum_{j=1}^d \|\beta_j\|_0 \leq K.$$

Why ℓ_0 ?

We care about estimation error as well as computation efficiency.

- The ℓ_0 constraint exactly describes the statistic modeling assumptions:

“The total number of arcs is bounded by $K \ll d^2$.”

- Avoids the statistical estimation gap between ℓ_0 and ℓ_1 norm

Objective: Estimate the True Graph

- Let β^* be the *true* coefficients (not the optimal solution to the optimization problem)
- We want to find an (approximate) solution $\hat{\beta}$ to the optimization problem and use it as an estimator of β^*

Intuition from (non-)Convex Geometry

The intuition traces way back to Shapley-Folkman Theorem that

“The sum of many nonconvex sets is close to being convex.”

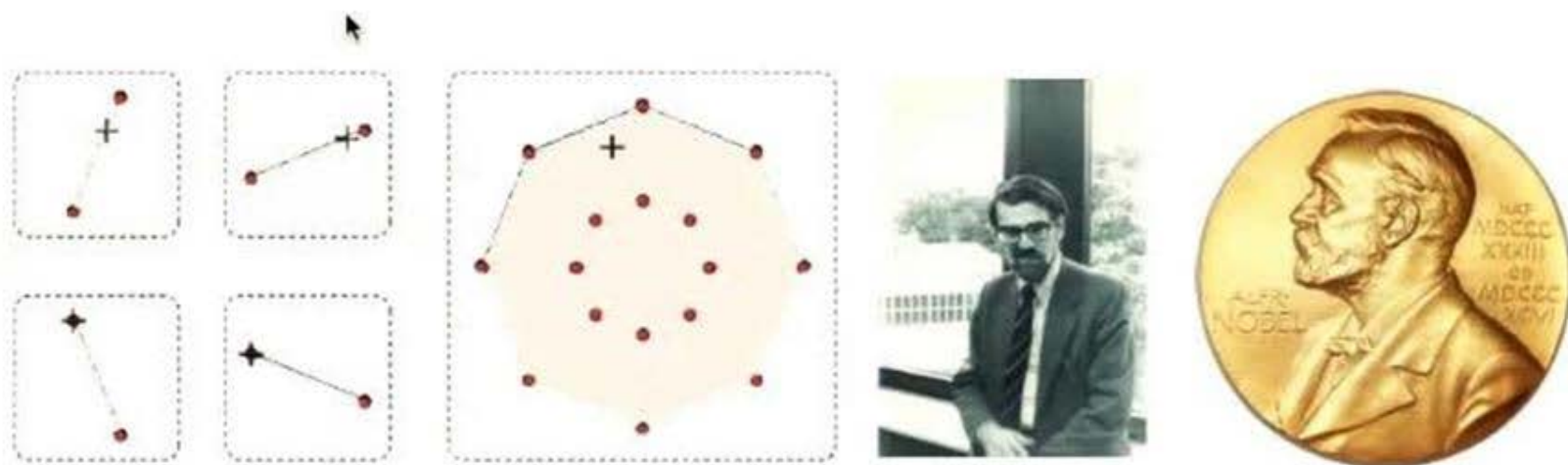


Figure : Shapley-Folkman Theorem reveals the insight that the sum of many nonconvex sets is close to being convex. Lloyd Shapley, 2012 Nobel Economics Laureate.

Blessing of Large Scales

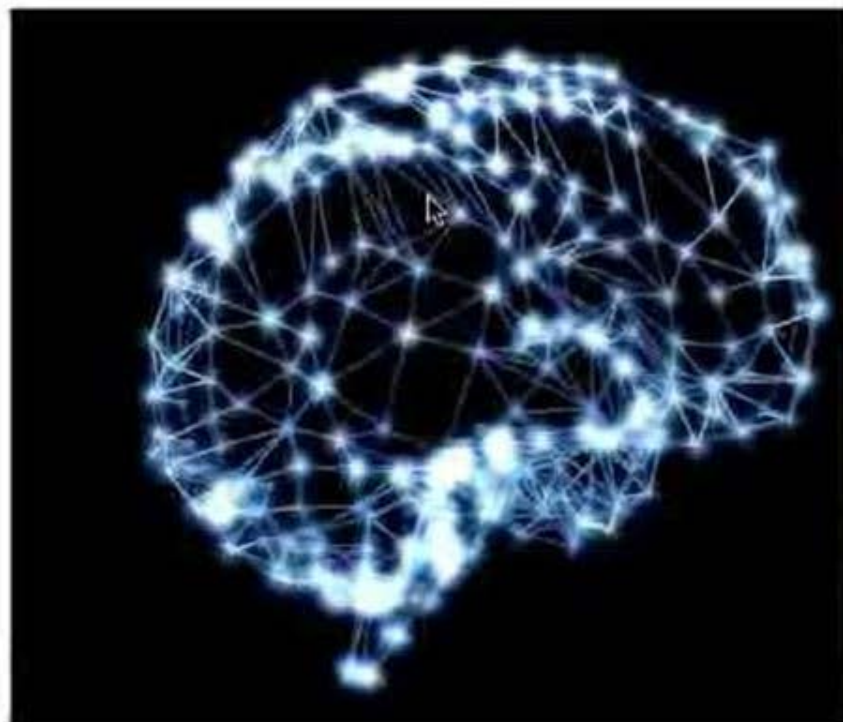
We use the dual solution $\hat{\beta}_j$'s as the estimators. Let β_j^* 's denote the truces and n denote sample size. We get

$$\frac{1}{d} \sum_{j=1}^d \|\hat{\beta}_j - \beta_j^*\|_2^2 = \underbrace{\mathcal{O}\left(\frac{(K/d) \log d}{n}\right)}_{\text{unavoidable}} + \underbrace{\mathcal{O}\left(\frac{\delta_n}{d}\right)}_{\text{duality gap}},$$

- $\mathcal{O}\left(\frac{(K/d) \log d}{n}\right)$ matches the information-theoretic lower bound.
- $\mathcal{O}\left(\frac{\delta_n}{d}\right)$ is the duality gap (difference between ℓ_0 problem and its dual) vanishes as graph size increases.
- As $n \rightarrow \infty$, $\delta_n \rightarrow \max_{j=1, \dots, d} D_{KL}\left(\mathbb{P}(\mathbf{X}_j) \parallel \mathbb{P}(\mathbf{X}_j | \mathbf{X}_{-j})\right)$. The duality gap constant converges to the maximum **Kullback-Leibler divergence** between the marginal distribution of a node j and its conditional distribution.
- For large graphs:

duality gap \ll unavoidable statistical uncertainty

Brain Voxel Network Estimation



The NEW ENGLAND
JOURNAL of MEDICINE

nature **methods**

Techniques for life scientists and chemists

Home | Current issue | Contents | Research | Archive | Authors & referees | About the journal

Home | archive | issue | commentary | full text

NATURE METHODS | COMMENTARY



Making sense of brain network data

Olaf Sporns

Nature Methods 10, 491–493 (2013) | doi:10.1038/nmeth.2485

Published online 30 May 2013

Figure: Understanding brain network is an initial step to understand our brain.

fMRI Data

The fMRI (Functional Magnetic Resonance Imaging) technology measures brain activity by detecting associated changes in blood flow. Number of voxels is $d \geq 30,000$.

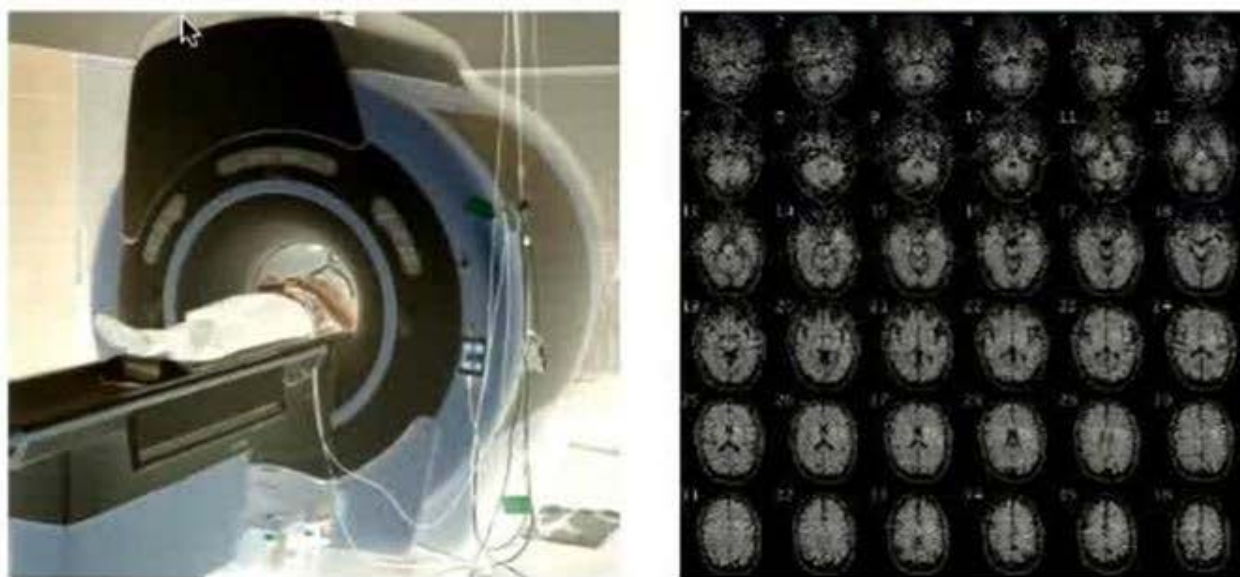


Figure : fMRI data.

Real Data Result

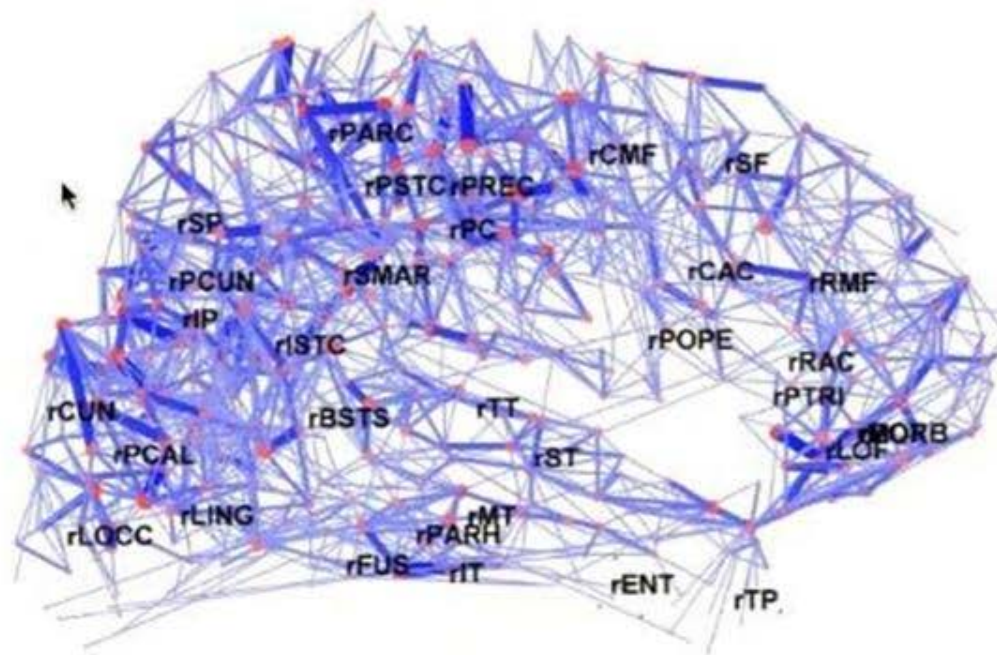


Figure : The estimated graph reveals that voxels related to a same function are highly connected.

More Generally: Distributed Collaborative Optimization

$$\max_{x_i \in \mathcal{X}_i, \forall i} u_0 \left(\frac{1}{N} \sum_{i=1}^N c_i(x_i) \right) + \frac{1}{N} \sum_{i=1}^N u_i(x_i)$$

- u_0 is the convex social utility function
- u_i is the individual utility of the i th player; $u_0, u_i, c_i, \mathcal{X}_i$ are nonconvex

Fenchel Dual Problem

$$\inf_{\lambda \in \mathbb{R}^m} \left\{ -u_0^*(\mu) + \frac{1}{N} \sum_{i=1}^N \sup_{x_i \in \mathcal{X}_i} \{ \lambda' c_i(x_i) + u_i(x_i) \} \right\}.$$

Duality Gap Results

- The duality gap vanishes to 0 as $N \rightarrow \infty$
- There exists a dual solution (out of many bad dual solutions) that achieves the duality gap.

Decentralized Coordination Algorithm: A Sketch

- At time t , the current price vector is λ_t .
- Under an incentive mechanism, each user reports multiple solutions to the penalized problem

$$\max_{x_i \in \mathcal{X}_i} u_i(x_i) - \lambda_t' c_i(u_i)$$

and receives a reward for providing flexibility.

- The central decision maker chooses an admissible solution for each user, and updates the price vector from λ_t to λ_{t+1} .

Dynamic Convergence to Social Optimum

- The iteration converges to an ϵ -social optimum with $\epsilon =$ duality gap.
- If users are willing to report their utility values, there exists an algorithm achieving regret

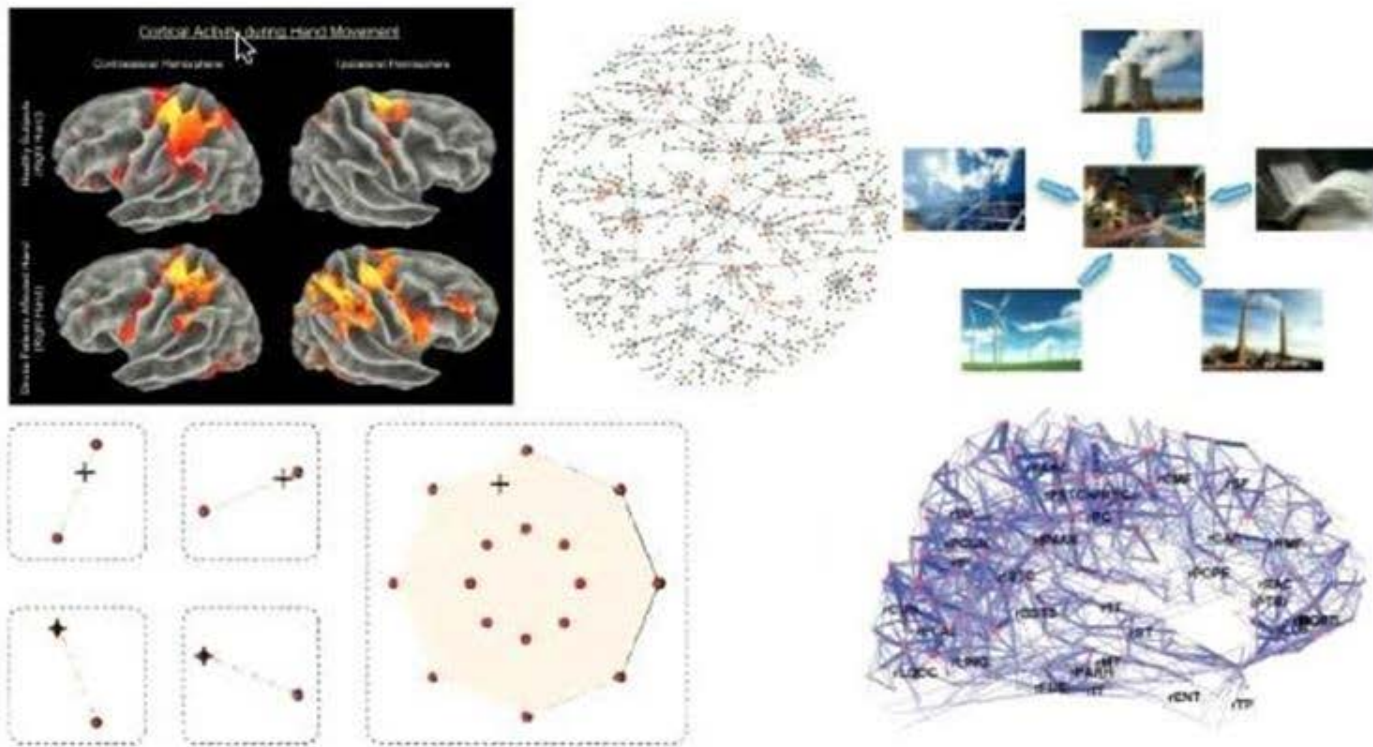
$$\mathcal{O}(\log T + \epsilon T)$$

- If users don't report their utility values, there is an algorithm achieving regret

$$\mathcal{O}(\sqrt{T} + \epsilon T)$$

Summary

$$\min_{\beta_1, \dots, \beta_d} \frac{1}{d} \sum_{j=1}^d \ell_j(\beta_j), \text{ subject to } \sum_{j=1}^d g_j(\beta_j) \leq b.$$



Vanishing duality gap in large distributed systems.