

Optimization with nonlinear Perron eigenvectors

Francesco Tudisco
Gran Sasso Science Institute (Italy)



SIAM LA21 Minitutorial
Applied Nonlinear Perron-Frobenius Theory

May 18, 2021

I will present a Perron-Frobenius type result for nonlinear operators. This result is stated as a global optimization algorithm for a class of constrained opt problems.

1. Motivation: graph-based unsupervised learning aka graph partitioning
2. Perron-Frobenius theorem for (sub) multihomogeneous maps
3. Some example applications

Matrix singular vectors, once again

Consider the constrained optimization problem

$$\begin{cases} \text{optimize} & \mathbf{x}_1^\top M \mathbf{x}_2 \\ \text{subject to} & \|\mathbf{x}_1\| = \|\mathbf{x}_2\| = 1 \end{cases}$$

In general, $f(\mathbf{x}_1, \mathbf{x}_2) = \mathbf{x}_1^\top M \mathbf{x}_2$ is not convex. However, we know how to compute global max and global min:

singular vectors and singular values of M :

$$M_i \mathbf{x}_i = \lambda \mathbf{x}_i, \quad i = 1, 2$$

with $M_1 = MM^\top$, $M_2 = M^\top M$.

Nonlinear singular vectors

For *sufficiently smooth* homogeneous functions f, g , the problem

$$\begin{cases} \text{optimize} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}_1) = g(\mathbf{x}_2) = 1 \end{cases}$$

with $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2)$, can be brought down to

nonlinear singular vector problem

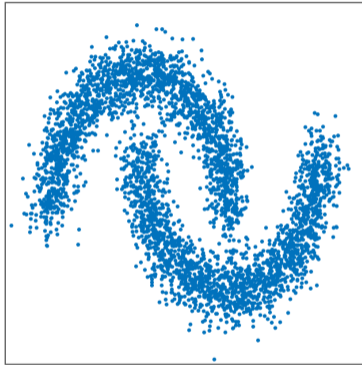
$$M_i(\mathbf{x})\mathbf{x}_i = \lambda\mathbf{x}_i, \quad i = 1, 2$$

where $M_i(\mathbf{x})$ are matrix-valued mappings, obtained differentiating f twice.

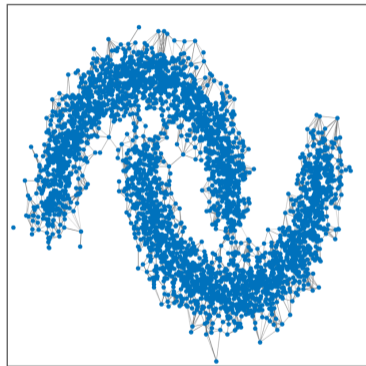
However, global max/min can be NP-hard...

Motivating example: graph clustering

Graph clustering

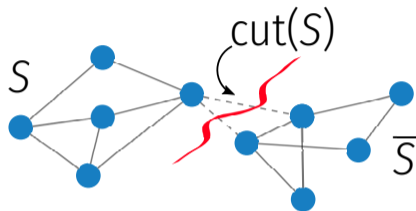


Graph clustering



$$\mathcal{G} = (V, E), V = \{1, \dots, n\}, E \subseteq V \times V$$

Balanced cut problem



$$\gamma(\mathcal{G}) = \min_{S \subseteq V} \frac{\text{cut}(S)}{\min\{|S|, |\bar{S}|\}}$$

$$\text{cut}(S) = \{ij \in E : i \in S, j \in \bar{S}\}, \quad \bar{S} = V \setminus S$$

General ratio-of-set-functions problem

$\gamma(\mathcal{G}) = \min_S \phi(S)/\psi(S)$ is the global minimum of the ratio of two set functions $\phi(S) = \text{cut}(S)$, $\psi(S) = \min\{|S|, |\bar{S}|\}$ such that:

1. ϕ, ψ are nonnegative
2. $\phi(V) = \psi(V) = 0$

General ratio-of-set-functions problem

$\gamma(\mathcal{G}) = \min_S \phi(S)/\psi(S)$ is the global minimum of the ratio of two set functions $\phi(S) = \text{cut}(S)$, $\psi(S) = \min\{|S|, |\bar{S}|\}$ such that:

1. ϕ, ψ are nonnegative
2. $\phi(V) = \psi(V) = 0$

In general, consider the problem

$$\min_{S \subseteq V} \vartheta(S), \quad \vartheta(S) = \frac{\phi(S)}{\psi(S)}$$

with $\phi, \psi : 2^V \rightarrow \mathbb{R}$ such that **1** and **2** hold

Homogeneous exact relaxation

Computing $\min \vartheta$ is in general NP-hard can we approximate it?

Homogeneous exact relaxation

Computing $\min \vartheta$ is in general NP-hard can we approximate it?

Theorem.

Given ϕ, ψ and $p \geq 1$, there exist homogeneous functions $f, g : \mathbb{R}^n \rightarrow \mathbb{R}$ of degree p such that, if λ is a solution to

$$\lambda = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) = 1, \end{cases}$$

then $\lambda \leq \min \vartheta \leq C^{p-1} \lambda^{1/p}$, in particular $\lambda \xrightarrow{p \rightarrow 1} \min \vartheta$.

Proof's sketch (for $p = 1$)

Consider the Lovasz extensions f, g of the functions ϕ and ψ .

1. $f(\mathbb{1}_S) = \phi(S), g(\mathbb{1}_S) = \psi(S)$

$$\min_{S \subseteq V} \frac{\phi(S)}{\psi(S)} \geq \min_{\mathbf{x} \in \mathbb{R}^n} \frac{f(\mathbf{x})}{g(\mathbf{x})} = \min_{\mathbf{x} \in \mathbb{R}^n} f\left(\frac{\mathbf{x}}{g(\mathbf{x})}\right)$$

2. $f(\mathbf{x}) = \sum_{i=0}^{n-1} \phi(S_{x_i}) |x_{i+1} - x_i| = \int_{-\infty}^{+\infty} \phi(S_t) dt$ where $S_t = \{k : x_k > t\}$.

$$\frac{f(\mathbf{x})}{g(\mathbf{x})} = \frac{\int_{-\infty}^{+\infty} \phi(S_t) dt}{\int_{-\infty}^{+\infty} \psi(S_t) dt} \geq \inf_t \frac{\phi(S_t)}{\psi(S_t)} \geq \min_{S \subseteq V} \frac{\phi(S)}{\psi(S)}$$

Based on

 [Hein, Setzer, NeurIPS 2012](#)

Back to the clustering problem

$\phi(S) = \text{cut}(S)$, $\psi(S) = \min\{|S|, |\bar{S}|\}$. The homogeneous functions f, g are:

$$f(\mathbf{x}) = \frac{1}{2} \sum_{ij=1}^n A_{ij} |x_i - x_j|^p \quad g(\mathbf{x}) = \|\mathbf{x} - \text{mean}(\mathbf{x})\mathbf{1}\|_p^p = 1$$

and we have that

$$\lambda = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) = 1 \end{cases} \iff \lambda = \text{smallest nonzero sol of:} \\ M_p(\mathbf{x})\mathbf{x} = \lambda\mathbf{x}$$

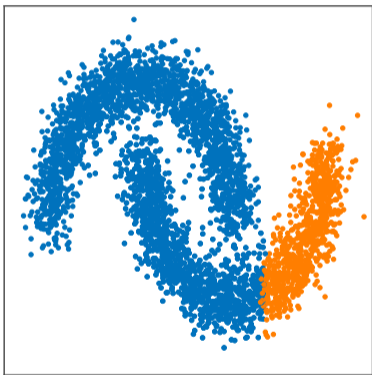
where M_p is a matrix-valued mapping, based on the graph p -Laplacian

Cheeger inequality

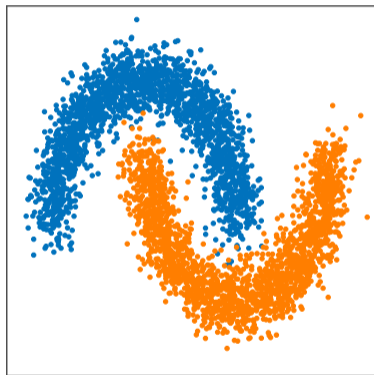
For $p = 2$ we obtain the famous Cheeger inequality

- $M_2 = L = \text{diag}(A\mathbb{1}) - A = \text{Graph Laplacian Matrix}$
- $\lambda = \text{Fiedler eigenvalue (or algebraic connectivity)}$
- $\min \vartheta = \gamma(\mathcal{G}) = \text{graph Cheeger constant}$
- $\lambda \leq \gamma(\mathcal{G}) \leq C\sqrt{\lambda}$

Linear (vs) nonlinear spectral clustering



$p = 2$



$p = 1$

Drawback

f and g are nonlinear and nonconvex in general, solving

$$\lambda = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) = 1, \end{cases}$$

can be very challenging.

E.g. think at the result $\lambda \xrightarrow{p \rightarrow 1} \text{Cheeger constant}$

Drawback

f and g are nonlinear and nonconvex in general, solving

$$\lambda = \begin{cases} \min_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}) \\ \text{subject to} & g(\mathbf{x}) = 1, \end{cases}$$

can be very challenging.

E.g. think at the result $\lambda \xrightarrow{p \rightarrow 1}$ Cheeger constant

However, we can compute λ to an arbitrary accuracy when f and g are **nonnegative** and **sub-multihomogeneous**.

Perron–Frobenius theorem for sub-multihomogeneous mappings

We are going to consider the

optimization of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ such that
the gradient of f is **sub-multihomogeneous**

To this end, we first introduce this concept.

Multihomogeneous mappings (two variables)

Suppose $\mathbf{x} = (\mathbf{x}_1, \mathbf{x}_2) \in \mathbb{R}^n$ with $\mathbf{x}_1 \in \mathbb{R}^{m_1}$, $\mathbf{x}_2 \in \mathbb{R}^{m_2}$.

Partition the gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ as:

$$\partial f = \begin{bmatrix} \partial_1 f \\ \partial_2 f \end{bmatrix}, \quad \partial_i f = \partial_{\mathbf{x}_i} f = \text{partial derivative w.r.t. variables in } \mathbf{x}_i$$

If there exists a 2×2 matrix Θ such that

$$\begin{cases} \partial_1 f(\lambda \mathbf{x}_1, \mathbf{x}_2) = \Theta_{11} \partial_1 f(\mathbf{x}) & \partial_1 f(\mathbf{x}_1, \lambda \mathbf{x}_2) = \Theta_{12} \partial_1 f(\mathbf{x}) \\ \partial_2 f(\lambda \mathbf{x}_1, \mathbf{x}_2) = \Theta_{21} \partial_2 f(\mathbf{x}) & \partial_2 f(\mathbf{x}_1, \lambda \mathbf{x}_2) = \Theta_{22} \partial_2 f(\mathbf{x}) \end{cases}$$

then ∂f is multihomogeneous. We write this compactly as $f \in \text{hom}'(\Theta)$.

Multihomogeneous functions have multihomogeneous gradient

An example of $f \in \text{hom}'(\Theta)$ are multihomogeneous functions.
In fact, if f is such that

$$f(\mathbf{x}_1, \dots, \lambda \mathbf{x}_j, \dots, \mathbf{x}_m) = \lambda^{\delta_j} f(\mathbf{x})$$

then it is easy to verify that

$$\partial_j f(\mathbf{x}_1, \dots, \lambda \mathbf{x}_j, \dots, \mathbf{x}_m) = \lambda^{\Theta_{ij}} \partial_j f(\mathbf{x})$$

with

$$\Theta = \begin{bmatrix} \delta_1 - 1 & \delta_2 & \dots & \delta_s \\ \delta_1 & \delta_2 - 1 & \dots & \delta_s \\ \vdots & & \ddots & \vdots \\ \delta_1 & \dots & \delta_{s-1} & \delta_s - 1 \end{bmatrix} = \mathbb{1}\delta^\top - I.$$

Euler's characterization (two-variables)

Suppose f is twice differentiable.

Then we can partition the Hessian of f accordingly

$$\partial^2 f = \left[\begin{array}{c|c} \partial_1 \partial_1 f & \partial_2 \partial_1 f \\ \hline \partial_1 \partial_2 f & \partial_2 \partial_2 f \end{array} \right]$$

Euler's theorem applied block-wise gives us

$$f \in \text{hom}'(\Theta) \iff \begin{cases} \partial_1 \partial_1 f(\mathbf{x}) \mathbf{x}_1 = \Theta_{11} \partial_1 f(\mathbf{x}) & \partial_2 \partial_1 f(\mathbf{x}) \mathbf{x}_2 = \Theta_{12} \partial_1 f(\mathbf{x}) \\ \partial_1 \partial_2 f(\mathbf{x}) \mathbf{x}_1 = \Theta_{21} \partial_2 f(\mathbf{x}) & \partial_2 \partial_2 f(\mathbf{x}) \mathbf{x}_2 = \Theta_{22} \partial_2 f(\mathbf{x}) \end{cases} \text{ for all } \mathbf{x} \succ 0$$

Euler's characterization of multihomogeneous mappings

Definition.

The gradient of $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is multihomogeneous if for some m there exists a partition of the variable $\mathbf{x} \in \mathbb{R}^n$

$$\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_m), \quad \mathbf{x}_i \in \mathbb{R}^{n_i}, \sum_i n_i = m$$

and a matrix $\Theta \in \mathbb{R}^{m \times m}$ such that

$$\partial_i \partial_j f(\mathbf{x}) \mathbf{x}_i = \Theta_{ij} \mathbf{x}_i$$

for all $i, j = 1, \dots, m$ and all positive vectors $\mathbf{x} \succ 0$

Sub-multihomogeneity

$f : \mathbb{R}^n \rightarrow \mathbb{R}$ twice differentiable is sub-multihomogeneous if there exists a partition $\mathbf{x} = (\mathbf{x}_1, \dots, \mathbf{x}_m)$ and a matrix Θ such that

$$|\Theta_{ij}| = \min \left\{ \lambda \geq 0 : |\partial_i \partial_j f(\mathbf{x})| \mathbf{x}_i \leq \lambda |\partial_i f(\mathbf{x})| \right\}$$

for all $i, j = 1, \dots, m$ and $\mathbf{x} \succ 0$

where $|\cdot|$ denotes absolute value taken component-wise.

We write this compactly $f \in \text{subhom}'(\Theta)$.

Global optimization with nonlinear Perron eigenvectors

Consider the problem

$$(*) \quad \begin{cases} \max_{\mathbf{x} \in \mathbb{R}^n} & f(\mathbf{x}_1, \dots, \mathbf{x}_m) \\ \text{subject to} & g_1(\mathbf{x}_1) = \dots = g_m(\mathbf{x}_m) = 1, \end{cases}$$

where:

- $f \in \text{subhom}'(\Theta)$
- $g_i \in \text{hom}(1 + \alpha_i)$, $\alpha_i \neq 0$
- ∂g_i is invertible on $\mathbb{R}_{++}^{n_i}$
- Both $(\partial g_i)^{-1} : \mathbb{R}^{n_i} \rightarrow \mathbb{R}^{n_i}$ and $\partial_i f : \mathbb{R}^n \rightarrow \mathbb{R}^{n_i}$ are positive mappings, i.e. map positive vectors into positive vector

Perron–Frobenius theorem

Let $B = \text{Diag}(\alpha_1, \dots, \alpha_m)^{-1} \Theta$.

If $\rho(|B|) < 1/2$ or $\rho(|B|) < 1$ and $\partial^2 f$ is a positive map, then

- There exists a unique solution $\mathbf{x}^* \in \mathbb{R}^n$ to (*) and $\mathbf{x}^* \succ 0$
- \mathbf{x}_i^* are nonlinear singular vectors, solution to

$$M(\mathbf{x}^*)\mathbf{x}_i^* = \lambda_i \mathbf{x}_i^*$$

corresponding to the largest nonlinear singular values $(\lambda_1, \dots, \lambda_m)$

- The nonlinear power method

$$\begin{cases} \mathbf{y} = M(\mathbf{x}^{(k)})\mathbf{x}_i^{(k)} = (\partial g_i)^{-1} \circ \partial_i f(\mathbf{x}^{(k)}) \\ \mathbf{x}^{(k+1)} = \begin{bmatrix} \frac{\mathbf{y}_1}{g_1(\mathbf{y}_1)} & \dots & \frac{\mathbf{y}_m}{g_m(\mathbf{y}_m)} \end{bmatrix} \end{cases} \quad k = 0, 1, 2, \dots$$

converges to \mathbf{x}^* as $O(\rho(|B|)^k)$, for any $\mathbf{x}_0 \succ 0$.

Sketch of the proof

Let $H : \mathbb{R}^n \rightarrow \mathbb{R}^n$ be the iterator of the nonlinear power method

1. \mathbf{x}^* is solution of (*) if and only if $H(\mathbf{x}^*) = \mathbf{x}^*$
2. $H(\mathcal{K}) \subseteq \mathcal{K}$ where $\mathcal{K} = \{\mathbf{x} \succ 0 : g_i(\mathbf{x}_i) = 1, \forall i\}$
3. \mathcal{K} is a complete metric space with respect to the Thompson metric δ_T
4. $\delta_T(H(\mathbf{x}), H(\mathbf{y})) \leq C \delta_T(\mathbf{x}, \mathbf{y})$ with

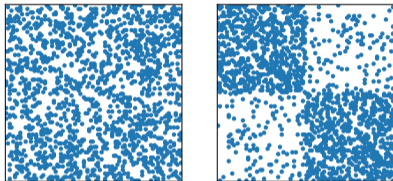
$$C = \sup_{\mathbf{x} \in \mathcal{K}} \left\| \text{diag}(H(\mathbf{x}))^{-1} |M(\mathbf{x})| \mathbf{x} \right\|_{\infty}$$

5. $C \leq \rho(|B|)$

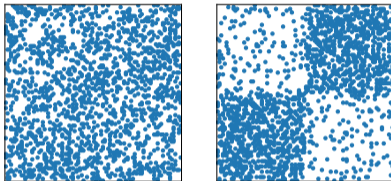
Core–Periphery detection in networks

Matrix reordering problems

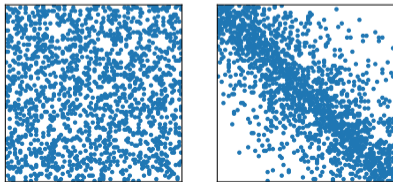
Clusters (communities)



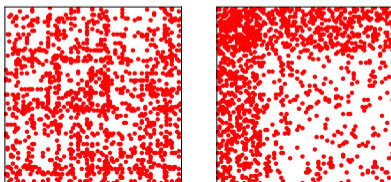
Bipartite (anti-communities)



Lattice (small-world)



Core-periphery



Core–periphery

 Borgatti, Everett, *Social Networks*, 1999

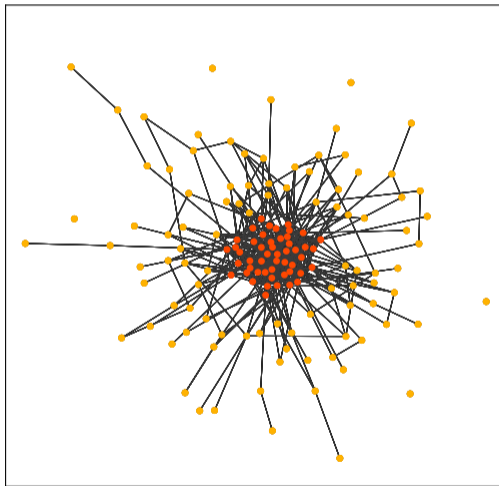
Core: nodes strongly connected across the whole network

Periphery: nodes strongly connected only to the core

 Csermely, London, Wu, Uzzi, *J. of Complex Networks*, 2013

 Rombach, Porter, Fowler, Mucha, *SIAM Review*, 2018

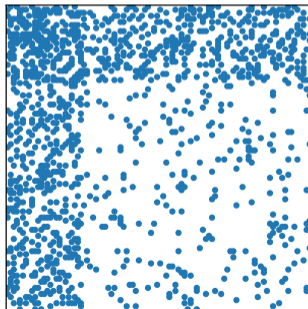
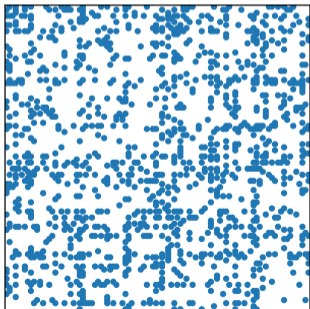
Core-periphery visualization



Core–periphery detection problem

Tasks:

1. **Reorder** nodes to reveal core–periphery structure
2. assign **coreness score** to nodes



Core-periphery kernel optimization

Core-score vector \mathbf{u} is such that :

if $u_i > u_j \implies i$ is closer to the core than j



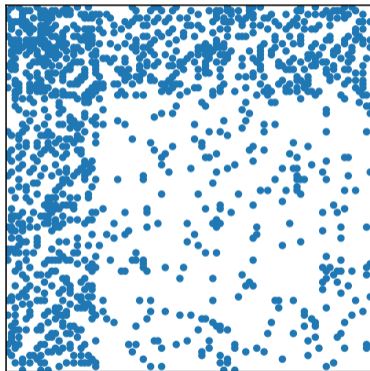
F T, D J Higham, SIAM Math of Data Science, 2019

Core-score vector as solution of the following constrained optimization

$$(cp) \quad \begin{cases} \text{maximize} & f_\alpha(\mathbf{u}) = \sum_{i,j=1}^n A_{ij} \kappa_\alpha(u_i, u_j) \\ \text{subject to} & \|\mathbf{u}\|_p = 1, \mathbf{u} \succeq 0 \end{cases}$$

with $A =$ adjacency matrix and $\kappa_\alpha(x, y) = \left(\frac{x^\alpha + y^\alpha}{2}\right)^{1/\alpha}$

Core-periphery kernel

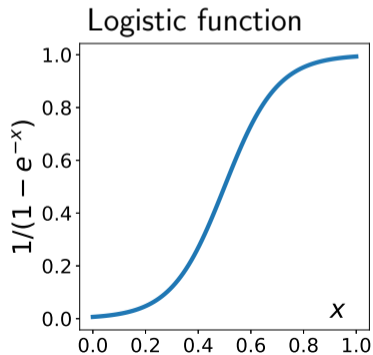


α large $\Rightarrow \kappa_\alpha(x, y) \approx \max\{x, y\}$

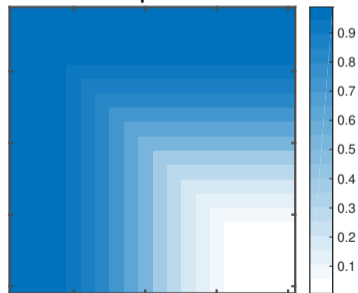
$f_\alpha(\mathbf{u}) = \sum_{ij} A_{ij} \kappa_\alpha(u_i, u_j)$ is large when edges $A_{ij} = 1$ involve at least one node with large core-score

Logistic core–periphery (LCP) random model

$$\text{Random graph: } \Pr(i \sim j) = \frac{1}{1 + e^{-\kappa_{\alpha}(u_i, u_j)}} =: p_{ij}(\mathbf{u})$$



Matrix of probabilities



Connection to CP kernel

Suppose we have a sample from the LCP random graph model, with nodes in arbitrary order.

Find \mathbf{u} that maximizes the **likelihood** $\lambda(\mathbf{u}) = \prod_{i \sim j} p_{ij}(\mathbf{u}) \prod_{i \not\sim j} (1 - p_{ij}(\mathbf{u}))$

Connection to CP kernel

Suppose we have a sample from the LCP random graph model, with nodes in arbitrary order.

Find \mathbf{u} that maximizes the **likelihood** $\lambda(\mathbf{u}) = \prod_{i \sim j} p_{ij}(\mathbf{u}) \prod_{i \not\sim j} (1 - p_{ij}(\mathbf{u}))$

Theorem:

If \mathbf{u} is a node labeling (permutation) then

$$\mathbf{u} \text{ solves (cp)} \iff \mathbf{u} \text{ maximizes the likelihood } \lambda$$

(Useful for testing core-periphery detection algorithms)

Connection with node degrees

If $p = 2$ and $\alpha = 1$ then $\kappa_1 =$ arithmetic mean

$$(cp) \iff \max_{\mathbf{u} \geq 0} \frac{\|A\mathbf{u}\|_1}{\|\mathbf{u}\|_2} = \|A\|_{2 \rightarrow 1}$$

and the maximizer is

$$\mathbf{u} = \text{degree vector}$$

Connection with eigenvector centrality

If $p = 1$ and $\alpha = 0$ then $\kappa_0 =$ geometric mean

$$(cp) \iff \max_{\mathbf{u} \geq 0} \frac{\mathbf{u}^T A \mathbf{u}}{\mathbf{u}^T \mathbf{u}} = \rho(A)$$

and the maximizer is

$$\mathbf{u} = \text{Perron eigenvector of } A$$

What about the general case

Using Perron–Frobenius

$$(cp) \quad \begin{cases} \text{maximize} & f_\alpha(\mathbf{u}) = \sum_{i,j=1}^n A_{ij} \kappa_\alpha(u_i, u_j) \\ \text{subject to} & \|\mathbf{u}\|_p = 1, \mathbf{u} \succeq 0 \end{cases}$$

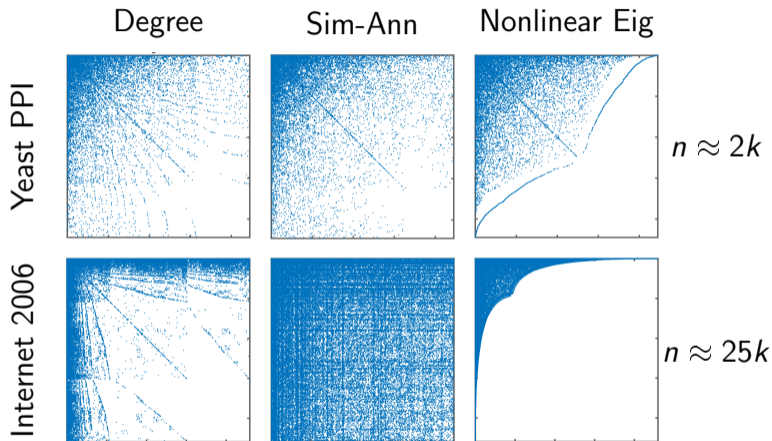
A direct computation reveals that

- $f_\alpha \in \text{subhom}'(\alpha - 1)$ and $\partial^2 f$ is positive
- $g(\mathbf{u}) = \|\mathbf{u}\|_p^p \in \text{hom}(p)$ and ∂g is invertible on \mathbb{R}_{++}^n

Then, if $\rho(|B|) = |\alpha - 1|/|p - 1| < 1$ we have that (cp) has a **unique positive solution**, which we can **compute to an arbitrary precision** and which coincides with the **nonlinear Perron eigenvector**

$$M(\mathbf{u})\mathbf{u} = \lambda\mathbf{u} \quad \text{with} \quad M(\mathbf{u})_{ij} = \frac{1}{u_i^{1-\alpha}} \frac{A_{ij} u_j^{\alpha-2}}{(u_i^\alpha/u_j^\alpha + 1)^{\frac{1-\alpha}{\alpha}}}$$

Qualitative results



Degree coincides with (cp) for $\alpha = 1$ and $p = 2$

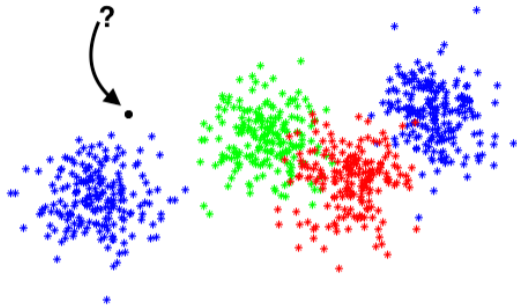
Convergence in a **few seconds** vs **several minutes** with Sim-Ann.

Beyond matrix-like forms

$f(x) = \sum_{ij} A_{ij} \kappa_{\alpha}(x_i, x_j)$ still looks like a nonlinear variation of a matrix form....

Polynomial neural networks

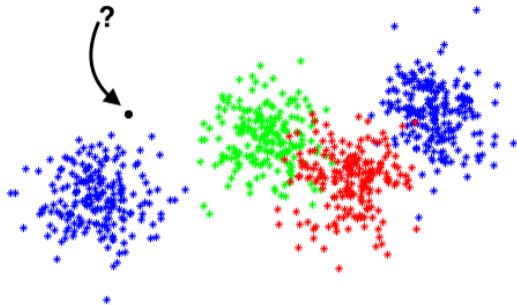
Supervised learning: training a nnet



Training points: $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(d)} \in \mathbb{R}^n$, $c_i \in \{1, 2, 3\}$ = class of $\mathbf{a}^{(i)}$

Activation matrices: (our variable) $X = (X_1, \dots, X_m)$

Supervised learning: training a nnet



Training points: $\mathbf{a}^{(1)}, \dots, \mathbf{a}^{(d)} \in \mathbb{R}^n$, $c_i \in \{1, 2, 3\}$ = class of $\mathbf{a}^{(i)}$

Activation matrices: (our variable) $X = (X_1, \dots, X_m)$

$$f(X) = \frac{1}{d} \sum_{i=1}^d \left[L\left(c_i, \varphi(X)(\mathbf{a}^{(i)})\right) + \mathbb{1}^\top \varphi(X)(\mathbf{a}^{(i)}) \right]$$

$$L(j, \mathbf{z}) = z_j - \log \left(e^{z_1} + e^{z_2} + e^{z_3} \right) \text{ (cross-entropy loss)}$$

Polynomial activation function

 A Gautier, Q Nguyen, M Hein, NeurIPS16

Train a classifier via

$$\begin{cases} \min_{X_1, \dots, X_m} f(X_1, \dots, X_m) \\ \text{subject to } \|X_1\|_{p_1} = \dots = \|X_m\|_{p_m} = 1, \end{cases}$$

with **activation functions** $\varphi_j(\mathbf{u}) = \mathbf{u}^{\mathbf{b}_j} = (u_1^{(\mathbf{b}_j)_1}, \dots, u_{n_j}^{(\mathbf{b}_j)_{n_j}})$

Polynomial activation function

 A Gautier, Q Nguyen, M Hein, NeurIPS16

Train a classifier via

$$\begin{cases} \min_{X_1, \dots, X_m} f(X_1, \dots, X_m) \\ \text{subject to} \quad \|X_1\|_{p_1} = \dots = \|X_m\|_{p_m} = 1, \end{cases}$$

with **activation functions** $\varphi_j(\mathbf{u}) = \mathbf{u}^{\mathbf{b}_j} = (u_1^{(\mathbf{b}_j)_1}, \dots, u_{n_j}^{(\mathbf{b}_j)_{n_j}})$

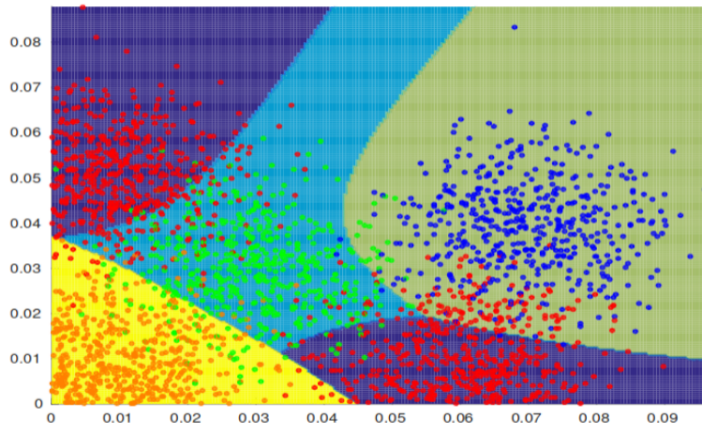
Theorem

There exists $\Theta = \Theta(\mathbf{b})$ such that

$$f(\mathbf{X}) = \frac{1}{d} \sum_i \left[L\left(c_i, \varphi(\mathbf{X})(\mathbf{a}^{(i)})\right) + \mathbf{1}^\top \varphi(\mathbf{X})(\mathbf{a}^{(i)}) \right]$$

is subhom'(Θ) and $\partial^2 f$ is positive. Conditions on \mathbf{b} to get $\rho(\Theta) < 1$.

Example 2D decision boundary



Conclusions

Conclusions

If you have a “nonnegative problem” – check out the nonlinear PF theory

- Website – `ftudisco.github.io/siam-nonlinear-pf-tutorial`
(feedback is very welcome)
- Book – soon (next year) to come

Some important questions concern:

- Convergence rate of nonlinear power method (e.g. for $Af(x) = \lambda x$)
- More advanced (faster) numerical eigensolvers

Thank you very much for your attention!