

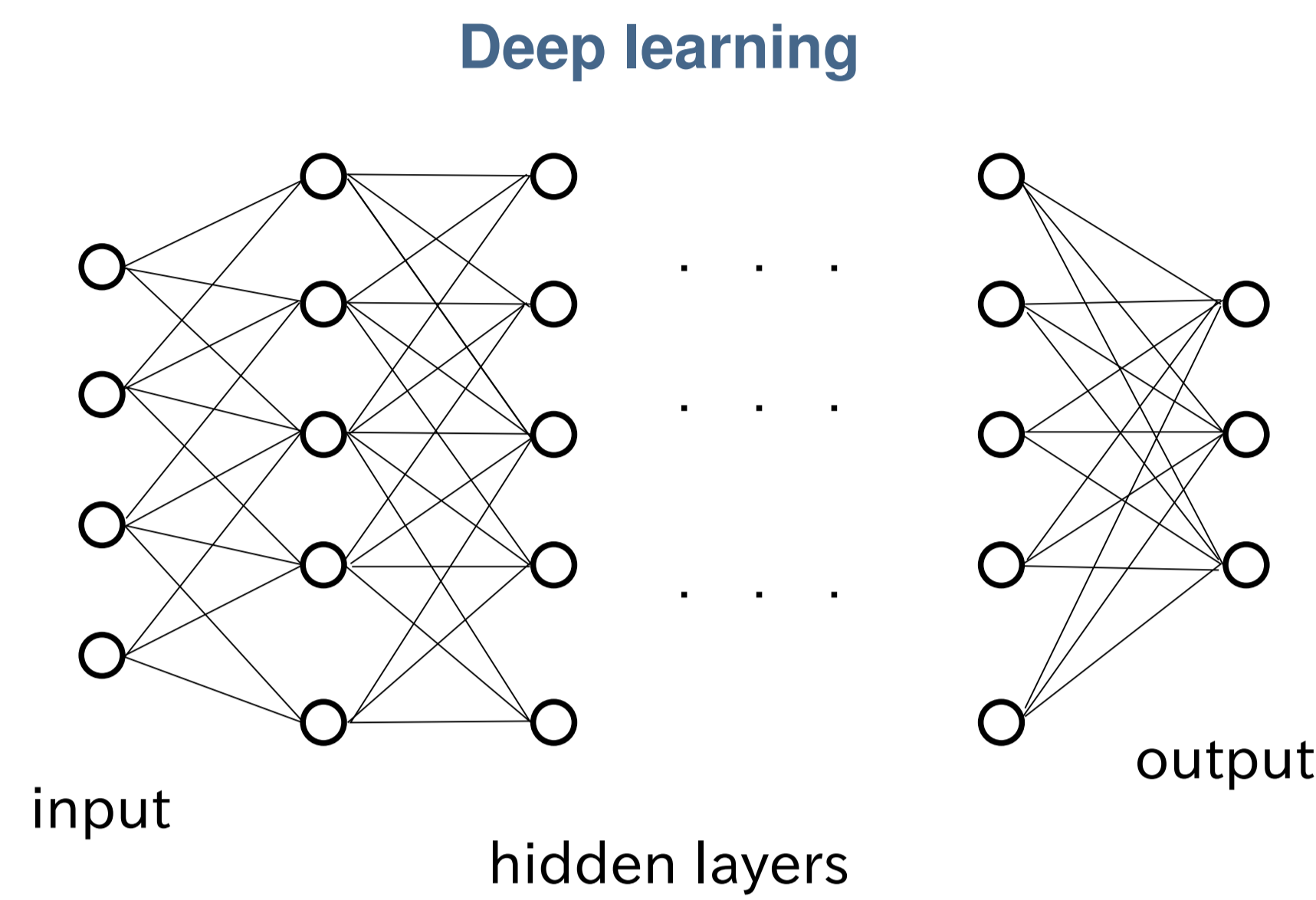
How Regularization Affects Critical Points in Linear Network

Amirhossein Taghvaei, Jin Kim, Prashant Mehta

Coordinated Science Laboratory, University of Illinois at Urbana-Champaign

Midwest Machine Learning Symposium, Chicago, 2017

Motivation



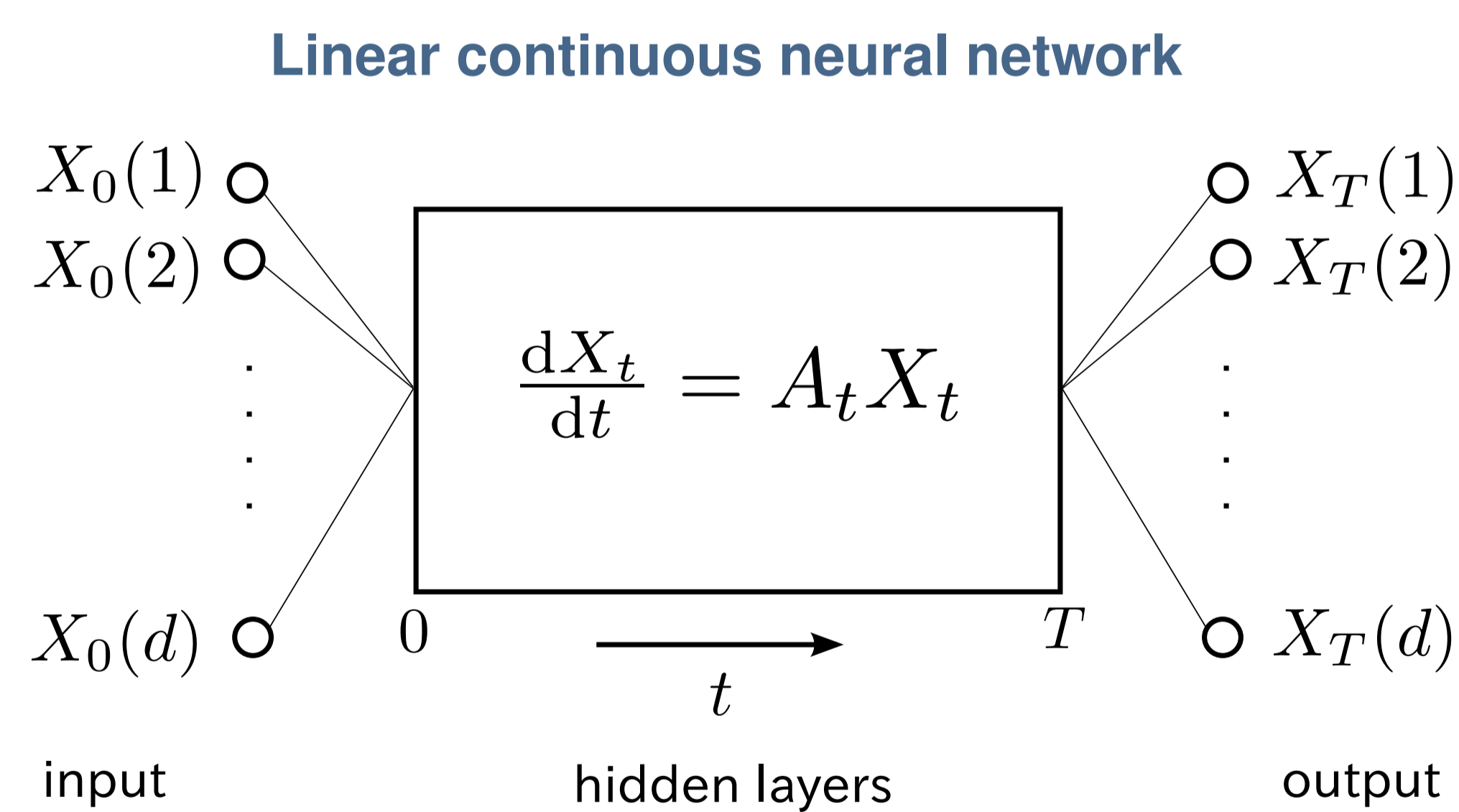
Objective: Analysis of the critical points of the associated non-convex optimization problem

This work: Analysis of the critical points of a linear network (with regularization)

Related work:

- ▶ A. M. Saxe, et. al. Exact solutions to the nonlinear dynamics of learning in deep linear neural networks (2013)
- ▶ M. Hardt and T. Ma. Identity matters in deep learning (2016).

Problem formulation



Model:

Data: $X_0 \in \mathbb{R}^d$ with $X_0 \sim p_0$

Measurement: $Z = RX_0 + \xi \in \mathbb{R}^d$

- ▶ R is a $d \times d$ matrix
- ▶ ξ is noise with mean zero
- ▶ $\Sigma := E[X_0 X_0^T]$ is invertible

Problem: Learn the linear transformation R with the linear continuous NN

Optimization problem

Optimal control formulation

$$\text{Minimize: } J[A] = \underbrace{\frac{\lambda}{2} \int_0^T \text{tr}(A_t^T A_t) dt}_{\text{regularization}} + \frac{1}{2} \underbrace{E[|X_T - Z|^2]}_{\text{mean-squared loss}}$$

$$\text{Subject to: } \frac{dX_t}{dt} = A_t X_t, \quad X_0 \sim p_0$$

- ▶ ($\lambda = 0$): No regularization.
- ▶ ($\lambda > 0$): Explicit regularization
- ▶ ($\lambda = 0^+$): The limit as $\lambda \rightarrow 0$. Models the dissipation in the learning.

Hamilton's formulation

Hamiltonian function:

$$H(x, y, B) = y^T Bx - \frac{\lambda}{2} \text{tr}(B^T B)$$

where $x, y \in \mathbb{R}^d$ and $B \in \mathbb{R}^{d \times d}$

Pontryagin's maximum principle: Suppose A_t is the minimizer. Then there exists a random process $Y : [0, T] \rightarrow \mathbb{R}^d$ such that

$$\begin{aligned} \frac{dX_t}{dt} &= +\frac{\partial H}{\partial y}(X_t, Y_t, A_t) = +A_t X_t, & X_0 &\sim p_0 \\ \frac{dY_t}{dt} &= -\frac{\partial H}{\partial x}(X_t, Y_t, A_t) = -A_t^T Y_t, & Y_T &= Z - X_T \end{aligned}$$

and A_t maximizes the expected value of the Hamiltonian

$$A_t = \arg \max_{B \in M_d(\mathbb{R})} E[H(X_t, Y_t, B)] = \frac{1}{\lambda} E[Y_t X_t^T]$$

Backpropagation (with dissipation)

First order variation:

$$\nabla J[A] := -E \left[\frac{\partial H}{\partial B}(X_t, Y_t, A_t) \right] = \lambda A_t - E[Y_t X_t^T]$$

where X_t and Y_t are obtained by solving the Hamilton's equations.

Stochastic gradient-descent:

$$A_t^{(k+1)} = A_t^{(k)} - \eta_k (\lambda A_t^{(k)} - Y_t^{(k)} X_t^{(k)T}),$$

where η_k is the step-size and $X_t^{(k)}$ and $Y_t^{(k)}$ are obtained by solving the Hamilton's equations:

$$\begin{aligned} \text{(Forward propagation)} \quad \frac{d}{dt} X_t^{(k)} &= +A_t^{(k)} X_t^{(k)}, & \text{with init. cond. } X_0^{(k)} \\ \text{(Backward propagation)} \quad \frac{d}{dt} Y_t^{(k)} &= -A_t^{(k)T} Y_t^{(k)}, & Y_T^{(k)} = \underbrace{Z^{(k)} - X_T^{(k)}}_{\text{error}} \end{aligned}$$

based on the sample $(X^{(k)}, Z^{(k)})$.

Critical points and characteristic equation

Main result:

▶ ($\lambda = 0$): Any A_t such that $\Phi_{0,T} = R$ where $\Phi_{0,t}$ is the state transition matrix corresponding to $\frac{dX_t}{dt} = A_t X_t$.

▶ ($\lambda > 0$): $A_t = e^{t(C-C^T)} C e^{-t(C-C^T)}$ where C is a solution of

$$\lambda C = F^T (R - F) \Sigma$$

where $F := e^{T(C-C^T)} e^{TC^T}$.

▶ ($\lambda = 0^+$): $A_t = e^{t(C-C^T)} C e^{-t(C-C^T)}$ where C is a solution of

$$e^{T(C-C^T)} e^{TC^T} = R$$

Examples

Case I: $R \in \mathbb{R}$ is a positive scalar

▶ ($\lambda = 0$): $A_t = \frac{1}{T} \log(R) + B_t$ for any B_t s.t. $\int_0^T B_t = 0$

▶ ($\lambda > 0$): $A_t = C = \frac{1}{T} \log(R) + O(\lambda)$

▶ ($\lambda = 0^+$): $A_t = C = \frac{1}{T} \log(R)$ (minimum norm solution)

Case II: R is a normal matrix with $\det(R) > 0$.

▶ ($\lambda > 0$): All the constant solutions are

$$A_t = C = \frac{1}{T} \log(R) + O(\lambda)$$

where $\log(R)$ is multi-valued

Example: $R = \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, $\log(R)_n = (\pi/2 + 2n\pi) \begin{bmatrix} 0 & -1 \\ 1 & 0 \end{bmatrix}$, $n \in \mathbb{Z}$

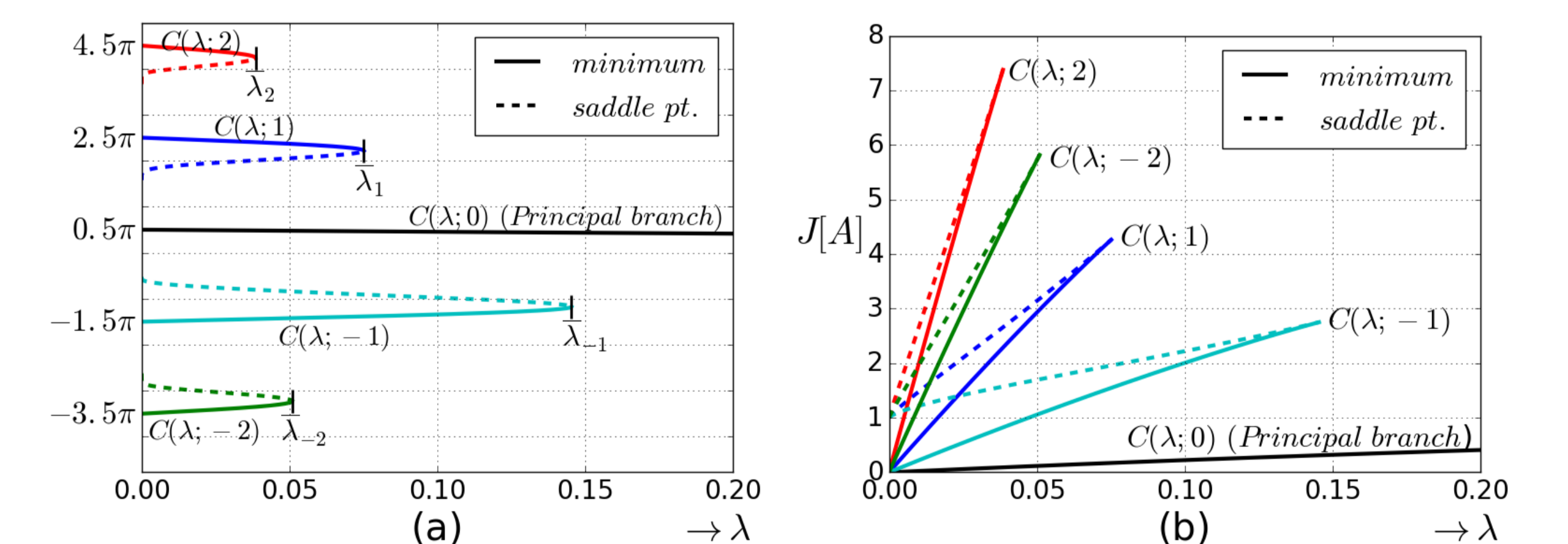


Figure : (a) Critical points (the $(2, 1)$ entry of the solution matrix $C(\lambda; n)$ is depicted for $n = 0, \pm 1, \pm 2$); (b) The cost $J[A]$ for these solutions.

Future work

1. Stability analysis of the critical points
2. Introducing nonlinearity to the network

Acknowledgment

National Science Foundation grants 1334987 and 1462773