

# THE GEOMETRIC STRUCTURE OF GENERALIZED SOFTMAX LEARNING

---

Xiangxiang Xu<sup>1</sup>

Collaborators: Shao-Lun Huang<sup>2</sup>   Lizhong Zheng<sup>3</sup>   Lin Zhang<sup>2</sup>

ITW 2018

<sup>1</sup>Dept. of Electronic Engineering, Tsinghua University

<sup>2</sup>Tsinghua-Berkeley Shenzhen Institute, Tsinghua University

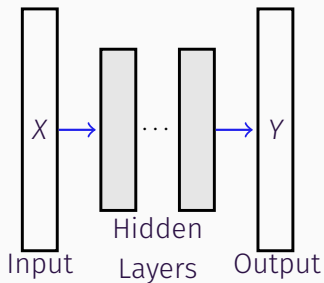
<sup>3</sup>Dept. of EECS, Massachusetts Institute of Technology

# INTRODUCTION

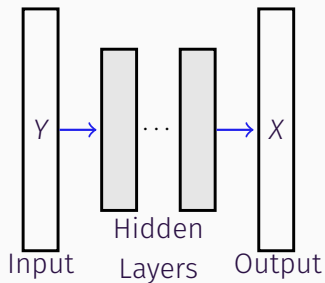
---

# REVERSE A NEURAL NETWORK

Use  $X$  to predict  $Y$

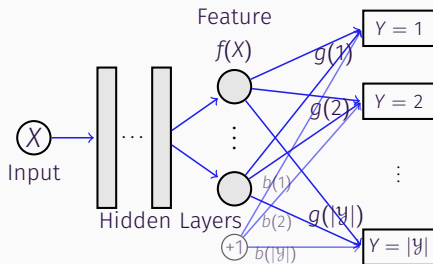


Use  $Y$  to predict  $X$



What do they have in common?

# NEURAL NETWORK AND M-PROJECTION



Softmax Output

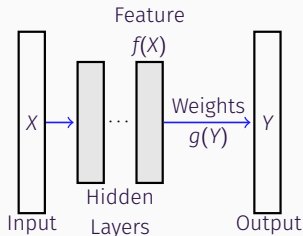
$$\tilde{P}_{Y|X}(y|X) \triangleq \frac{e^{f^T(x)g(y)+b(y)}}{\sum_{y' \in \mathcal{Y}} e^{f^T(x)g(y')+b(y')}}$$

## Maximum likelihood estimation

$$\begin{aligned}(f, g, b)^* &= \arg \max_{f, g, b} \mathbb{E}_{P_{X, Y}} \left[ \log \tilde{P}_{Y|X}(Y|X) \right] \\ &= \arg \min_{f, g, b} D(P_{X, Y} \| P_X \tilde{P}_{Y|X})\end{aligned}$$

# GENERALIZED SOFTMAX LEARNING

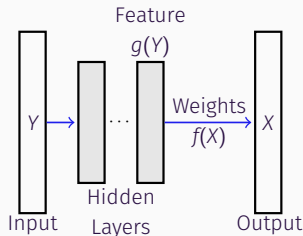
Use  $X$  to predict  $Y$



$$\tilde{P}_{Y|X}(y|x) \propto \exp(f^T(x)g(y) + b(y))$$

$$\text{minimize } D(P_{X,Y} \| P_X \tilde{P}_{Y|X})$$

Use  $Y$  to predict  $X$



$$\tilde{P}_{X|Y}(x|y) \propto \exp(g^T(y)f(x) + a(x))$$

$$\text{minimize } D(P_{X,Y} \| P_Y \tilde{P}_{X|Y})$$

Generalized softmax learning

$$\text{minimize } D(P_{X,Y} \| Q_{X,Y})$$

$$Q_{X,Y}(x, y) \propto \exp(f^T(x)g(y) + a(x) + b(y))$$

# GENERALIZED SOFTMAX LEARNING

---

The M-projection of  $P_{X,Y}$  onto

$$\mathcal{E}_k \triangleq \{Q_{X,Y}: Q_{X,Y}(x,y) \propto \exp(f^T(x)g(y) + a(x) + b(y))\},$$

where  $f: \mathcal{X} \rightarrow \mathbb{R}^k, g: \mathcal{Y} \rightarrow \mathbb{R}^k, a: \mathcal{X} \rightarrow \mathbb{R}, b: \mathcal{Y} \rightarrow \mathbb{R}$ .

$$\mathcal{M}_{\mathcal{E}_k}(P_{X,Y}) \triangleq \arg \min_{Q_{X,Y} \in \mathcal{E}_k} D(P_{X,Y} \| Q_{X,Y}),$$

Stationary distributions  $\mathcal{E}_k^0: \nabla = 0$

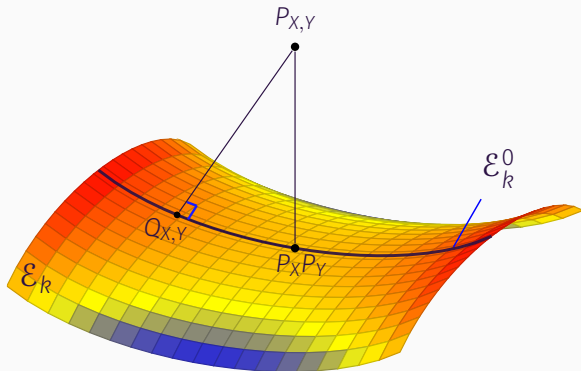
$$Q_X = P_X, \quad Q_Y = P_Y,$$

$$\mathbb{E}_{Q_{X|Y}} [f(X) | Y] = \mathbb{E}_{P_{X|Y}} [f(X) | Y],$$

$$\mathbb{E}_{Q_{Y|X}} [g(Y) | X] = \mathbb{E}_{P_{Y|X}} [g(Y) | X].$$

$$P_X P_Y \in \mathcal{E}_k^0$$

## Pythagorean theorem



$$\forall Q_{X,Y} \in \mathcal{E}_k^0, D(P_{X,Y} \| Q_{X,Y}) + D(Q_{X,Y} \| P_X P_Y) = D(P_{X,Y} \| P_X P_Y).$$

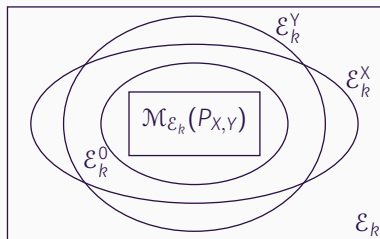
$$P_X P_Y \in \mathcal{M}_{\mathcal{E}_k}(P_{X,Y}) \iff X \perp Y$$



Stationary distributions  $\mathcal{E}_k^0$

$$Q_{X,Y} \in \mathcal{E}_k^0 \implies Q_X = P_X, Q_Y = P_Y$$

Different distribution families

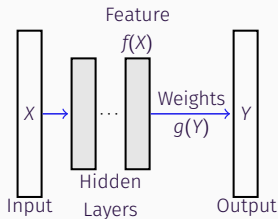


- $\mathcal{E}_k^0$ : stationary distributions
- $\mathcal{M}_{\mathcal{E}_k}(P_{X,Y})$ : GSL solutions
- $\mathcal{E}_k^X \leftrightarrow Q_{X,Y} \in \mathcal{E}_k, Q_X = P_X$
- $\mathcal{E}_k^Y \leftrightarrow Q_{X,Y} \in \mathcal{E}_k, Q_Y = P_Y$

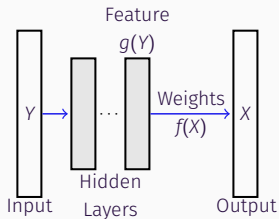
$$\mathcal{M}_{\mathcal{E}_k}(P_{X,Y}) = \mathcal{M}_{\mathcal{E}_k^X}(P_{X,Y}) = \mathcal{M}_{\mathcal{E}_k^Y}(P_{X,Y})$$

# EQUIVALENCE OF SOFTMAX LEARNING PROBLEMS

Use  $X$  to predict  $Y \leftrightarrow \mathcal{M}_{\mathcal{E}_k^X}(P_{X,Y})$



Use  $Y$  to predict  $X \leftrightarrow \mathcal{M}_{\mathcal{E}_k^Y}(P_{X,Y})$



$$\tilde{P}_{Y|X}(y|x) \propto \exp(f^T(x)g(y) + b(y))$$

$$\text{minimize } D(P_{X,Y} \| P_X \tilde{P}_{Y|X})$$

$$\tilde{P}_{X|Y}(x|y) \propto \exp(g^T(y)f(x) + a(x))$$

$$\text{minimize } D(P_{X,Y} \| P_Y \tilde{P}_{X|Y})$$

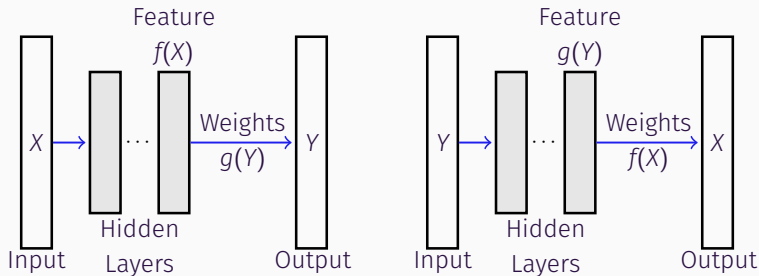
Generalized softmax learning  $\leftrightarrow \mathcal{M}_{\mathcal{E}_k}(P_{X,Y})$

$$\text{minimize } D(P_{X,Y} \| Q_{X,Y})$$

$$Q_{X,Y}(x, y) \propto \exp(f^T(x)g(y) + a(x) + b(y))$$

# SYMMETRY IN NEURAL NETWORKS

Symmetry between a network and its reverse network



## Local analysis regime

$$P_{X,Y} \approx P_X P_Y$$

## The $B$ matrix

$$B_{X,Y}(x,y) \triangleq \frac{P_{X,Y}(x,y)}{\sqrt{P_X(x)}\sqrt{P_Y(y)}}$$

## Solutions

$f^*(x) \leftrightarrow$  Left singular vectors of  $B_{X,Y}$

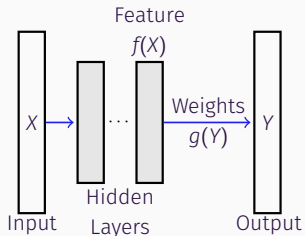
$g^*(y) \leftrightarrow$  Right singular vectors of  $B_{X,Y}$

## SIMULATION RESULTS

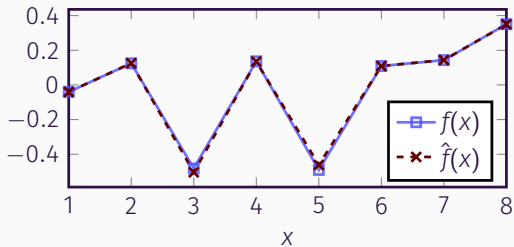
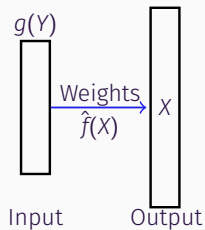
---

# SIMULATION

## Original network



## Reverse network\*



## SUMMARY

---

# SUMMARY

- Equivalence of softmax learning problems
- Discriminative  $\leftrightarrow$  Generative
- Asymmetric  $\leftrightarrow$  Symmetric
- Symmetry in neural networks

