

Error and Entropy in the Function Space of Multi-layer Networks

Bo Li (黎勃) and David Saad

22 January 2018

Outline

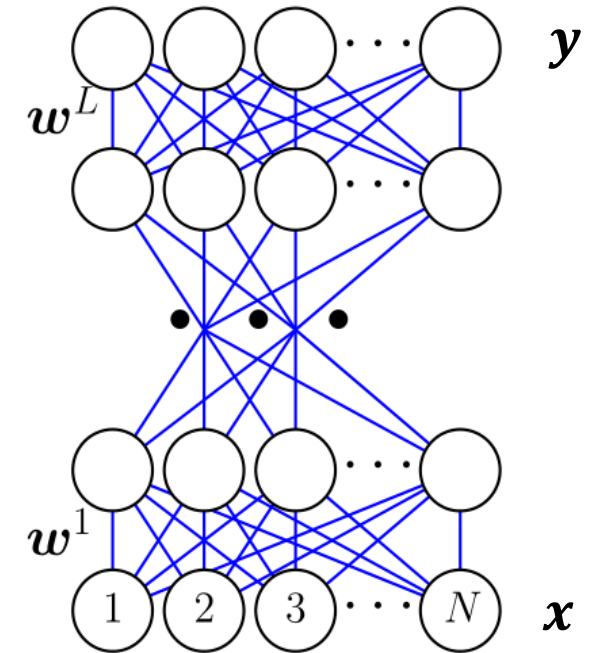
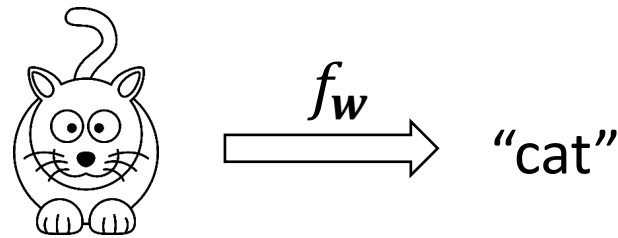
- Deep learning machines and why entropy in function-space matters
- Statistical mechanics of learning from examples
- Mapping multi-layer networks to dynamical systems - Generating functional analysis
- Continuous and discrete weights – framework and results
- Summary and future work

Deep Learning Machines

Implement an input-output mapping

$$y = f_w(x),$$

where the parameters w are to be estimated based on the training data $\{(\xi^\mu, \sigma^\mu)\}_{\mu=1,2,\dots,P}$ to perform a desired mapping.



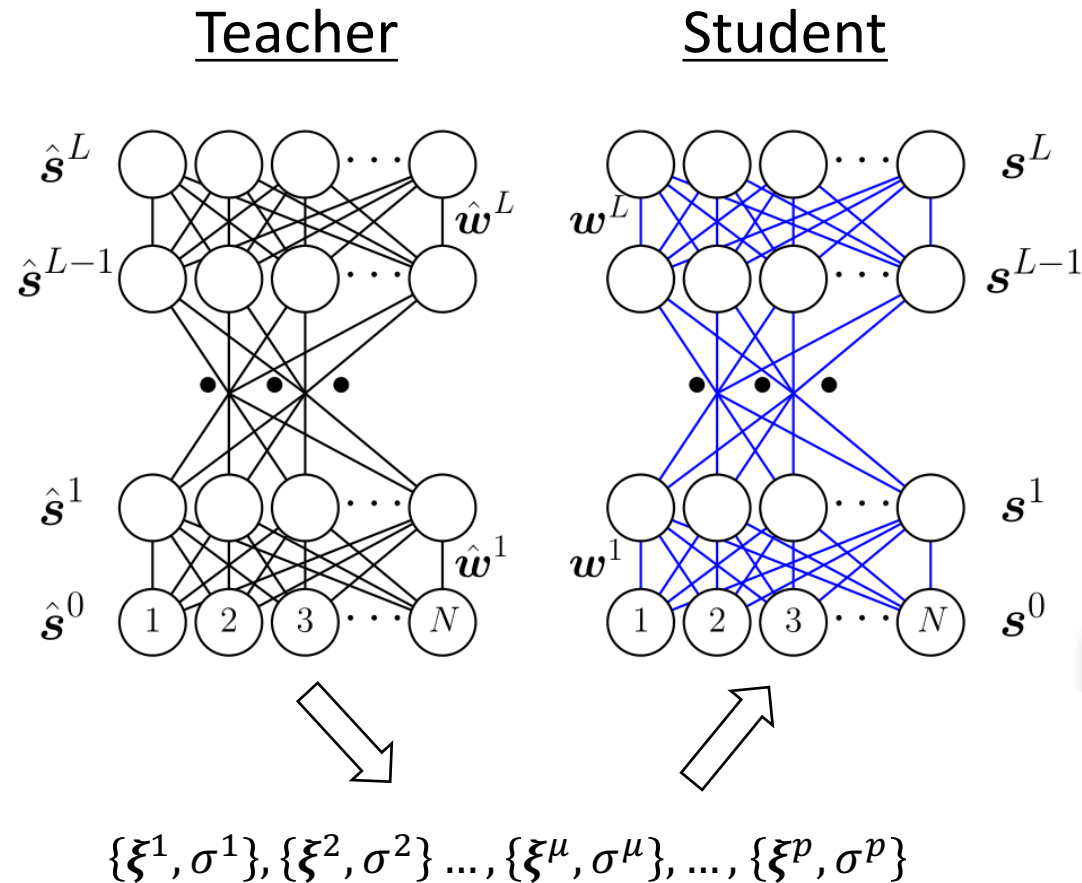
We want to understand:

- (i) Their **generalization ability** even with numerous parameters
- (ii) Nature of the **internal representations**

Macroscopic Analysis – Typical Behavior of Single Layer Machines

- Mapped to **disordered** systems of **infinite dimension**
- Use **replica analysis, cavity method** (batch learning); **dynamical** methods for on-line learning; exploit **high dimensionality**
- **Typical** behavior (in contrast to worse-case) of storage capacity and generalization curves
- Technically quite involved (single or two-layer systems)
- Input data structure and internal representations are rarely addressed

Teacher-student Scenario for DLM?



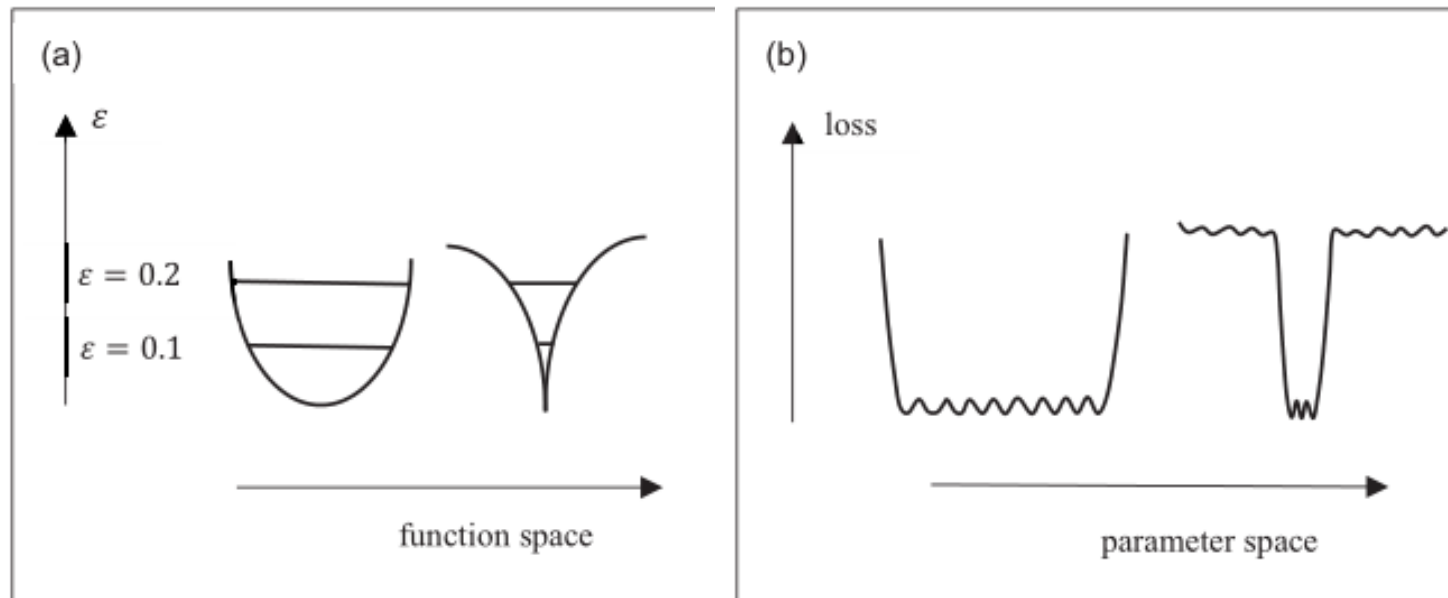
Difficulties:

- Constraints imposed by the examples (input-output pairs) on the hidden units are complex – **recursive nonlinear mapping**.
- Permutation, reflection and other **symmetries/invariances of hidden units**, no simple relation between teacher-student overlap and generalization error.

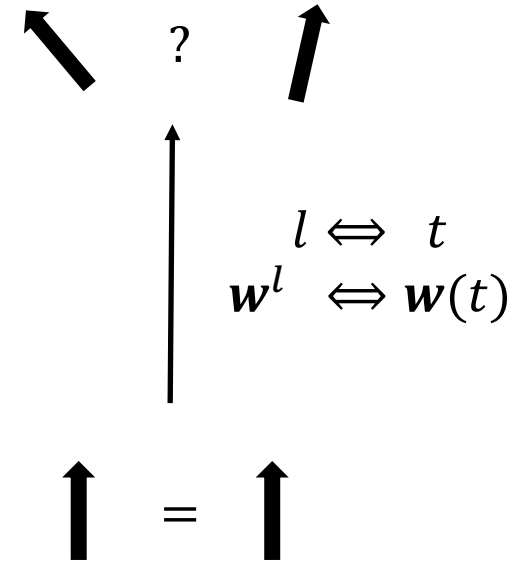
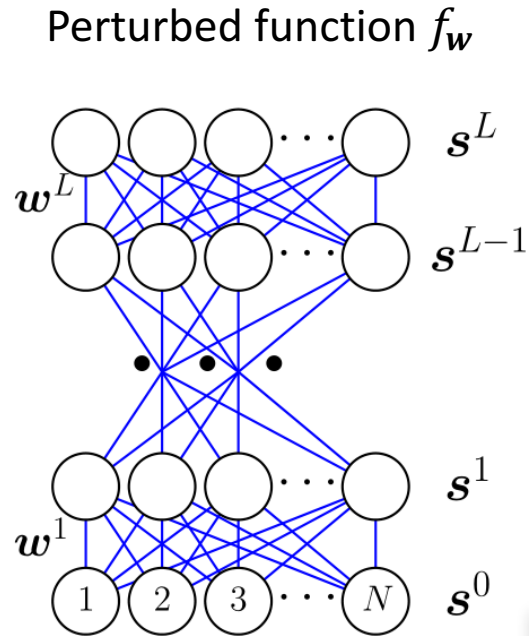
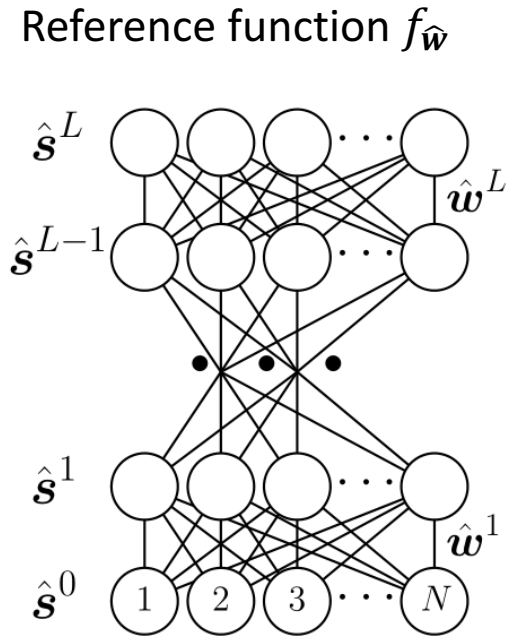
Especially interesting in the **over-parametrized** regime

Function Space, Error and Entropy

- We would like to approximate a reference/target function $f_{\hat{w}}$, as closely as possible from data.
- Given noisy data, sub-optimal training methods - more relevant to find *good approximations*. **How many such functions exist?**
- The entropy (log-volume) of functions at distance- ε away from $f_{\hat{w}}$ indicates how easy it is obtaining them.



Exploring Function Space in DLM



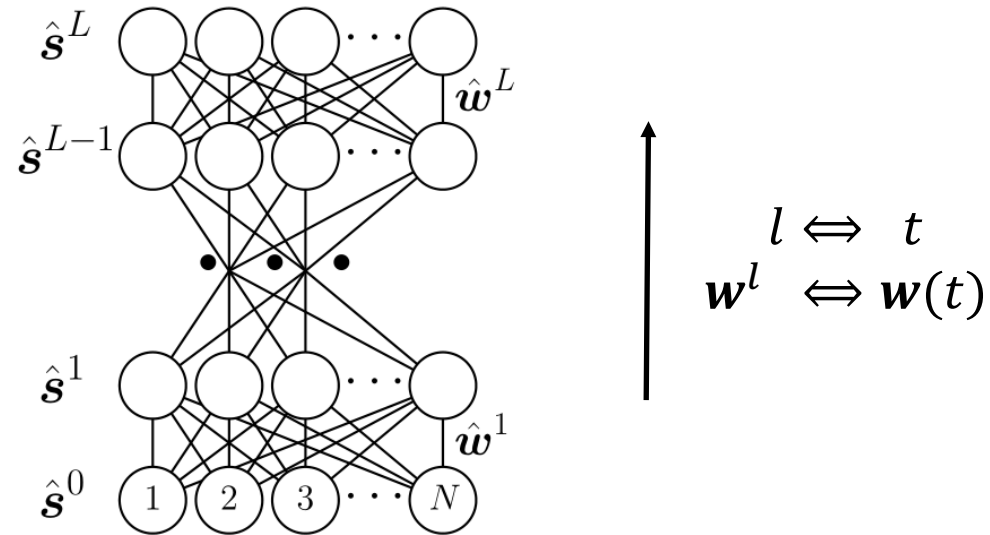
We map the DLM to **disordered spin systems** with discrete dynamics, $\hat{s}_i^l, s_i^l \in \{1, -1\}$, activation function is **sign function** $sgn(x)$.

The framework can be generalized to **real variables** and **other activation functions**.

Investigate the function sensitivity under small perturbations

$$w^l = \text{Perturb}(\hat{w}^l)$$

Deep Learning Machines as Dynamical Systems



Related work

- Poole et al., NIPS 2016 - Mean field theory to study input sensitivity and expressivity
- Li et al., arXiv:1710.09513, 2017 - Optimal control theory (Pontryagin's maximum principle) to devise new training algorithms

DLM as a Stochastic Dynamical System

- The layer evolution of two coupled DLMs:

$$P(\hat{\mathbf{s}}^l | \hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}, \beta) = \prod_i \frac{\exp \beta \hat{s}_i^l \hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1})}{2 \cosh \beta \hat{s}_i^l \hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1})}, P(\mathbf{s}^l | \mathbf{w}^l, \mathbf{s}^{l-1}, \beta) = \dots,$$

$\hat{h}_i^l(\hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}) = \sum_j \hat{w}_{ij}^l \hat{s}_j^{l-1} / \sqrt{N}$, β is the inverse-temperature quantifying the noise level; deterministic rule in the zero-noise limit $\beta \rightarrow \infty$.

- Any observable is given by

$$\langle O \rangle := \sum_{\{\hat{\mathbf{s}}^l, \mathbf{s}^l\}} O \cdot P(\hat{\mathbf{s}}^0) \delta_{\hat{\mathbf{s}}^0, \mathbf{s}^0} \prod_l P(\hat{\mathbf{s}}^l | \hat{\mathbf{w}}^l, \hat{\mathbf{s}}^{l-1}, \beta) \cdot P(\mathbf{s}^l | \mathbf{w}^l, \mathbf{s}^{l-1}, \beta),$$

summed over all the trajectories subject to the path measure.

- For discrete spins, the **overlap** between activities of the two systems is of interest

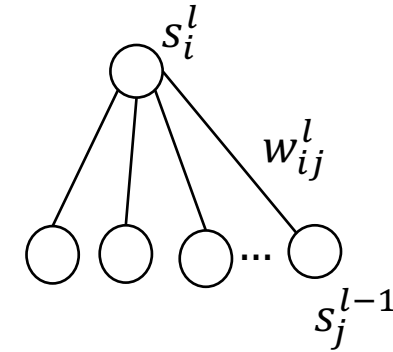
$$q^l(\hat{\mathbf{w}}, \mathbf{w}, \beta) = \frac{1}{N} \sum_i \langle \hat{s}_i^l s_i^l \rangle$$

Generating Functional Analysis

- Generating functional (characteristic function)

$$\Gamma[\hat{\boldsymbol{\psi}}, \boldsymbol{\psi}] := \left\langle \exp \left\{ -i \sum_{l,i} (\hat{\psi}_i^l \hat{s}_i^l + \psi_i^l s_i^l) \right\} \right\rangle,$$

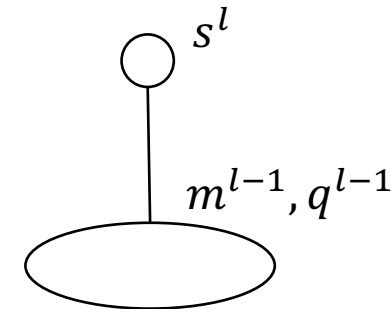
moments such as magnetization $\langle \hat{s}_i^l \rangle$ and overlap $\langle \hat{s}_i^l s_i^l \rangle$ can be obtained by differentiating $\Gamma[\hat{\boldsymbol{\psi}}, \boldsymbol{\psi}]$; angled brackets – average over all paths.



- Interested in the **typical** behavior of an ensemble of networks $\hat{\boldsymbol{w}} \sim P(\hat{\boldsymbol{w}})$, overbar – quenched average

$$\begin{aligned} \overline{\Gamma[\hat{\boldsymbol{\psi}}, \boldsymbol{\psi}]} &:= \sum_{\{\hat{\boldsymbol{w}}^l, \boldsymbol{w}^l\}} \Gamma[\hat{\boldsymbol{\psi}}, \boldsymbol{\psi}] P(\hat{\boldsymbol{w}}) P(\boldsymbol{w}) \\ &= \int \prod_l \frac{dQ^l dq^l}{2\pi/N} e^{N\Psi[\boldsymbol{q}, \boldsymbol{Q}]} \approx e^{N\Psi[\boldsymbol{q}_e, \boldsymbol{Q}_e]}, \text{ in the limit } N \rightarrow \infty, \end{aligned}$$

Represented by macroscopic order parameters; the saddle point $\boldsymbol{q}_e, \boldsymbol{Q}_e = \text{extr}_{\boldsymbol{q}, \boldsymbol{Q}} \Psi(\boldsymbol{q}, \boldsymbol{Q})$ satisfies certain self-consistent **mean-field** equation.



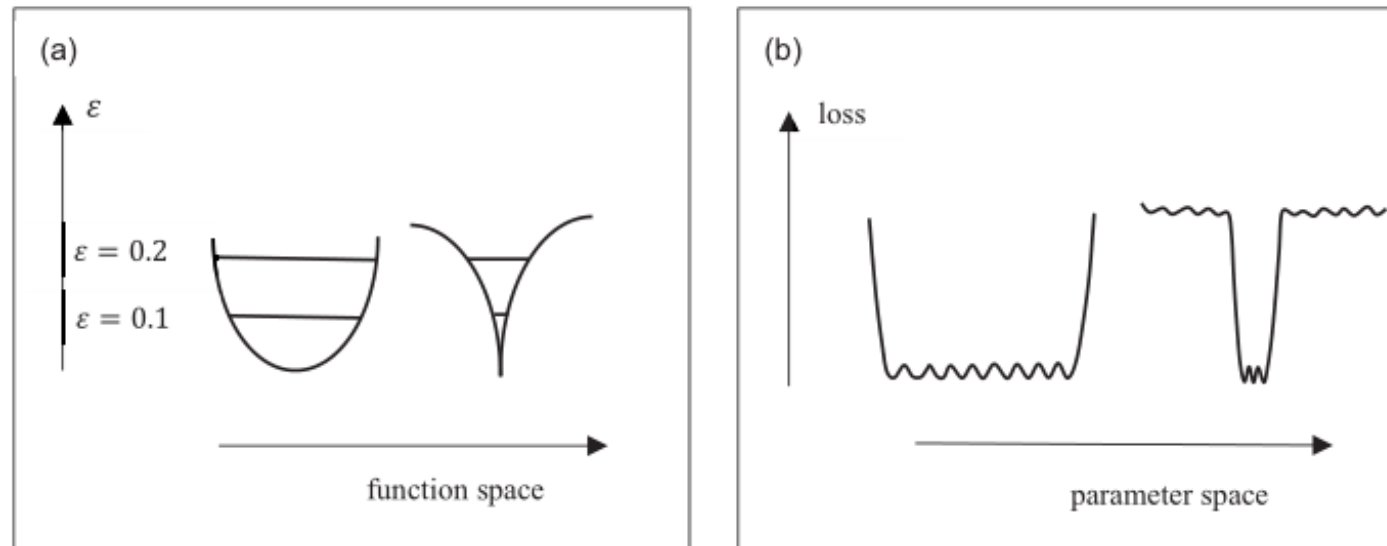
Function Error and Entropy

- Function error is defined as the **expected Hamming distance** of output layers between $f_{\hat{w}}$ and f_w

$$\varepsilon := \frac{1}{2N} \sum_{i=1}^N \overline{|\hat{s}_i^L - s_i^L|} = \frac{1}{2} (1 - q^L),$$

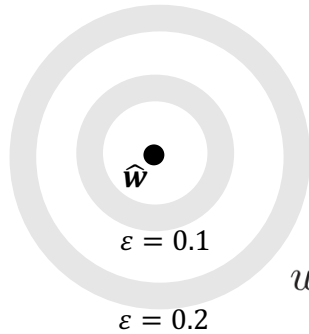
which provides a distance measure between $f_{\hat{w}}$ and f_w .

- Also interested in the entropy (log-volume) of f_w at distance- ε away from the reference function $f_{\hat{w}}$, e.g.,



Fully-connected Networks – Continuous/Binary weights

Consider fully-connected networks, with $P(\hat{\mathbf{w}}^l) = \prod_j P(\hat{w}_{ij}^l)$



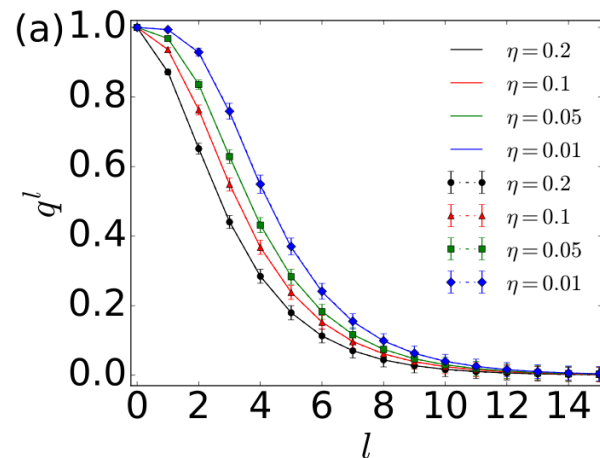
Continuous weights:

$$\hat{w}_{ij}^l \sim \mathcal{N}(0, \sigma^2) \quad \text{Perturbation strength at layer } l$$

$$w_{ij}^l = \sqrt{1 - (\eta^l)^2} \hat{w}_{ij}^l + \eta^l \delta w_{ij}^l$$

$$\delta w_{ij}^l \sim \mathcal{N}(0, \sigma^2)$$

Typical overlaps: $q^l = \frac{2}{\pi} \sin^{-1} \left(\sqrt{1 - (\eta^l)^2} q^{l-1} \right)$

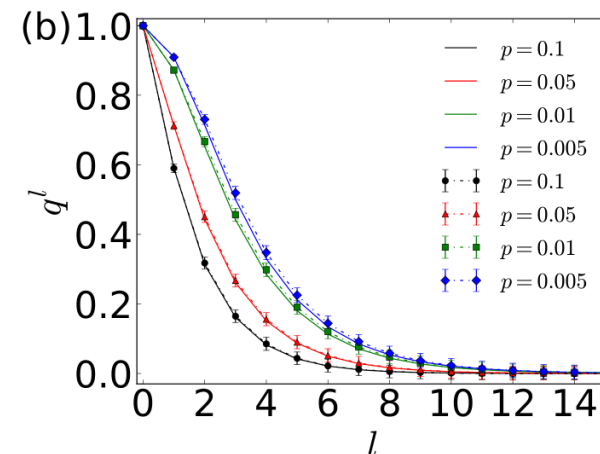


Binary weights:

$$P(\hat{w}_{ij}^l) = \frac{1}{2} \delta_{\hat{w}_{ij}^l, 1} + \frac{1}{2} \delta_{\hat{w}_{ij}^l, -1}$$

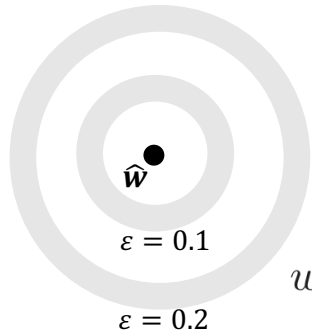
$$P(w_{ij}^l) = (1 - p^l) \delta_{w_{ij}^l, \hat{w}_{ij}^l} + p^l \delta_{w_{ij}^l, -\hat{w}_{ij}^l}$$

$$q^l = \frac{2}{\pi} \sin^{-1} \left((1 - 2p^l) q^{l-1} \right)$$



Entropy of Perturbed Functions

Consider fully-connected networks, with $P(\hat{\mathbf{w}}_i^l) = \prod_j P(\hat{w}_{ij}^l)$



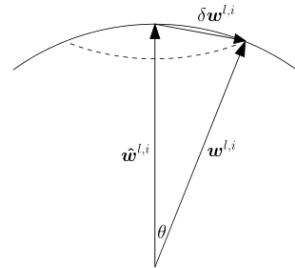
Continuous weights:

$$\hat{w}_{ij}^l \sim \mathcal{N}(0, \sigma^2)$$

Perturbation strength at layer l

$$w_{ij}^l = \sqrt{1 - (\eta^l)^2} \hat{w}_{ij}^l + \eta^l \delta w_{ij}^l$$

$$\delta w_{ij}^l \sim \mathcal{N}(0, \sigma^2)$$



Binary weights:

$$P(\hat{w}_{ij}^l) = \frac{1}{2} \delta_{\hat{w}_{ij}^l, 1} + \frac{1}{2} \delta_{\hat{w}_{ij}^l, -1}$$

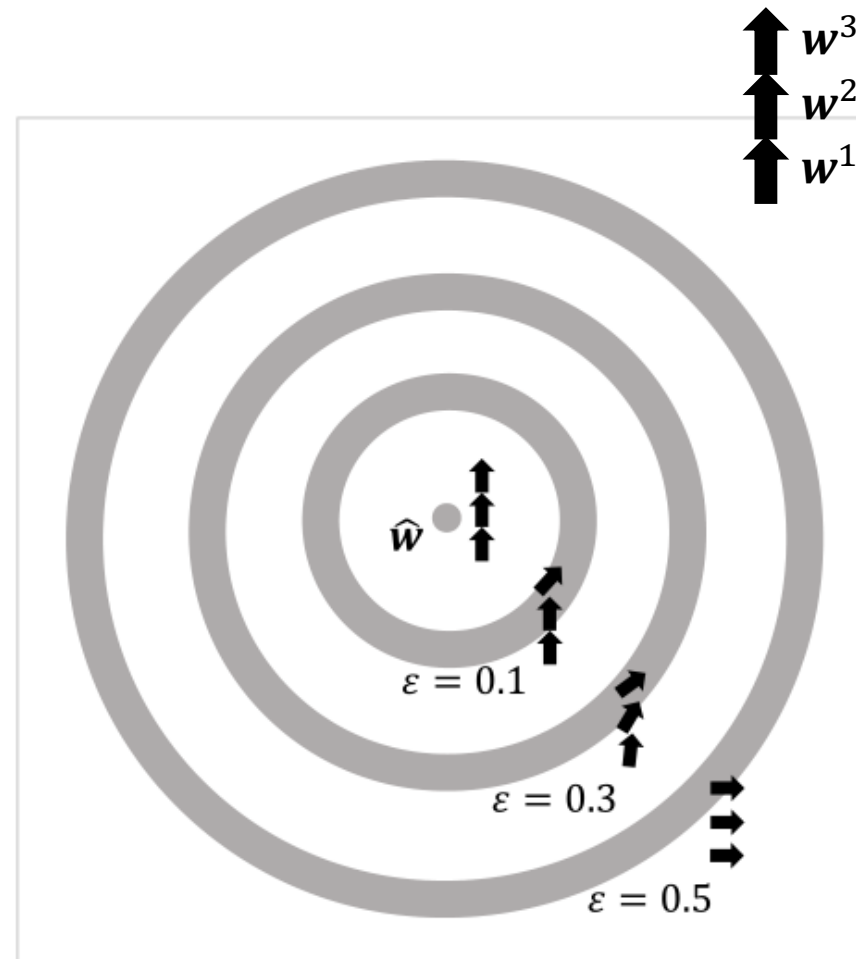
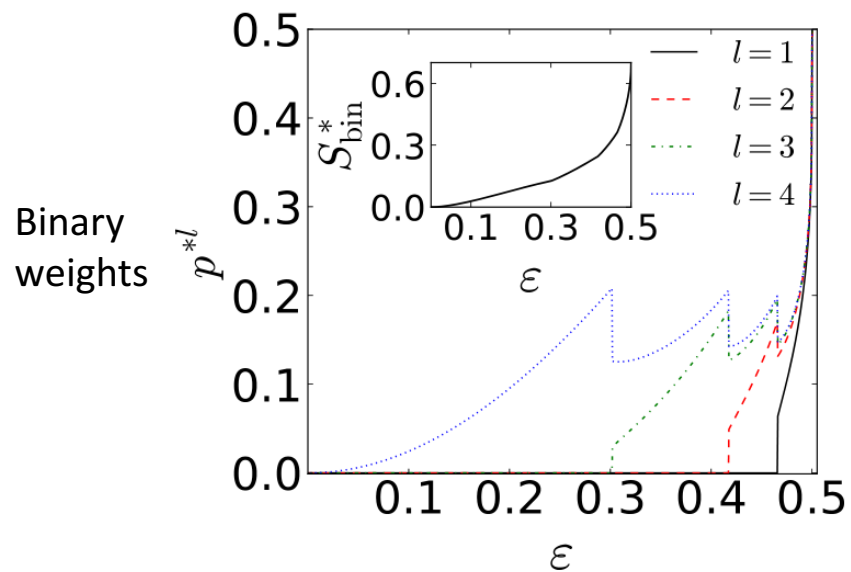
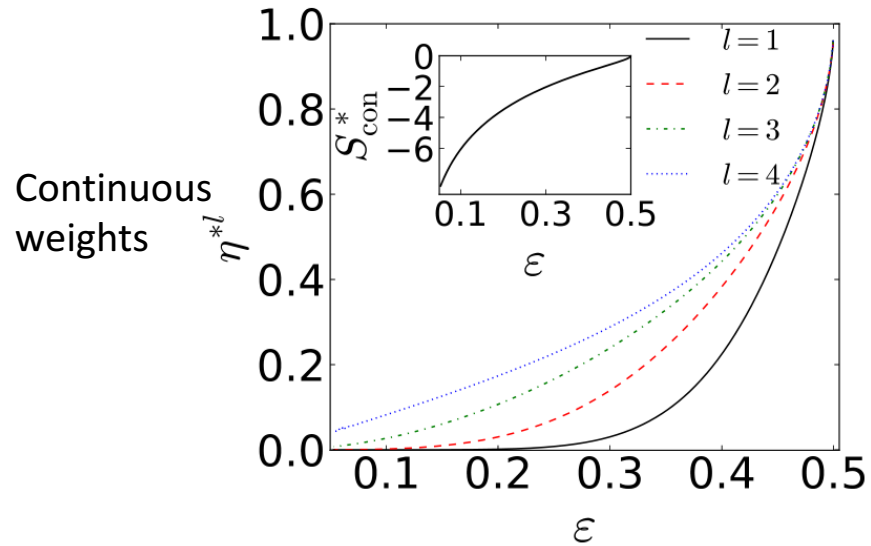
$$P(w_{ij}^l) = (1 - p^l) \delta_{w_{ij}^l, \hat{w}_{ij}^l} + p^l \delta_{w_{ij}^l, -\hat{w}_{ij}^l}$$

Entropy of $f_{\mathbf{w}}$: $S_{\text{con}}(\{\eta^l\}) = \frac{1}{L} \sum_{l=1}^L \log \eta^l$ $S_{\text{bin}}(\{p^l\}) = \frac{1}{L} \sum_{l=1}^L -p^l \log p^l - (1 - p^l) \log(1 - p^l)$

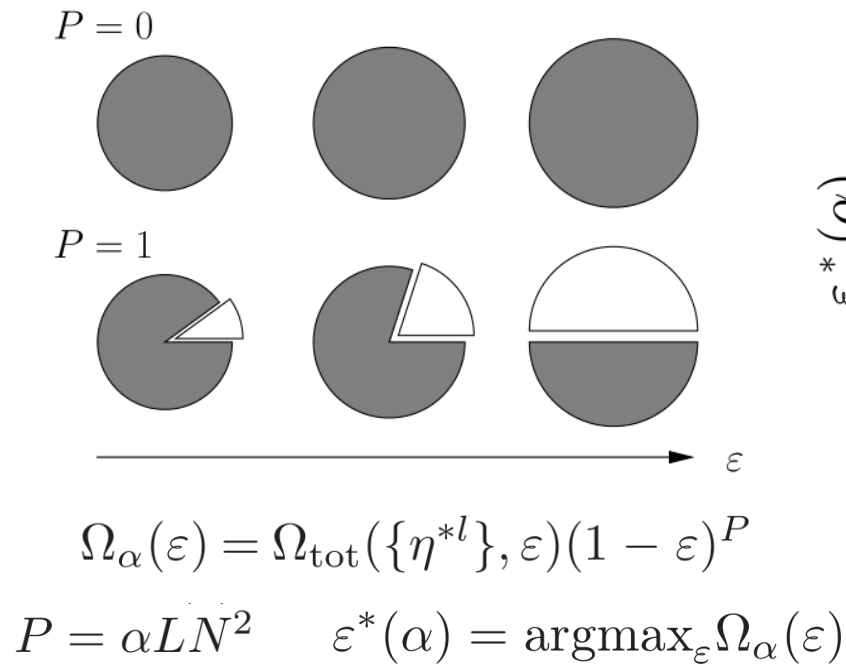
The distance- ε surface of $f_{\mathbf{w}}$ with volume $\Omega(\{\eta^l\}) = \exp N_p S_{\text{con}}(\{\eta^l\})$, is **exponentially** dominated by the maximum-entropy solutions when $N_p \rightarrow \infty$:

$$\eta^{*l} = \arg \max_{\eta^l} S_{\text{con}}(\{\eta^l\}), \quad \text{s.t. } q^L(\{\eta^l\}) = 1 - 2\varepsilon$$

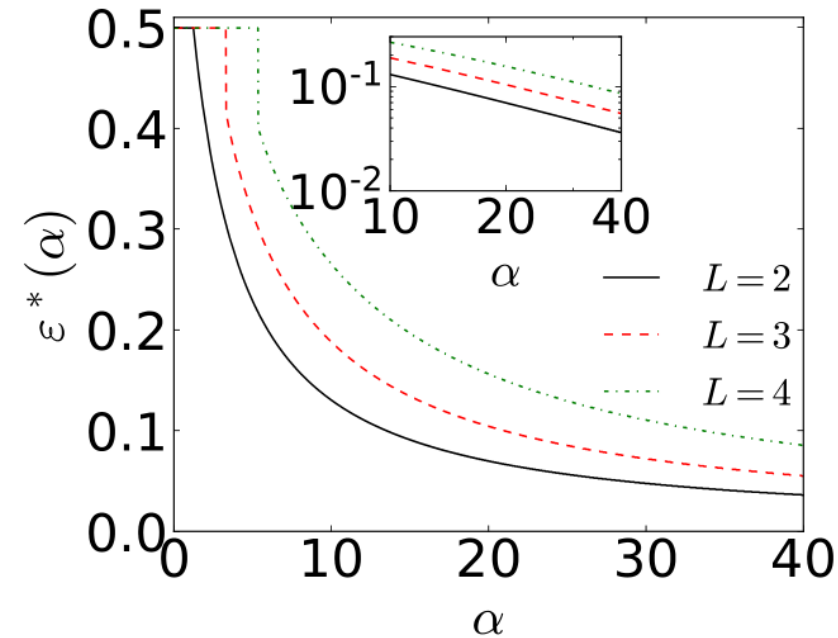
Earlier Layers Converge First When Decreasing ϵ



Approximated Generalization Curve (dense DLM with continuous weights)



Annealed theory of learning



Relevant in small ε (large α) limit.

Sparsely Connected Binary Networks

- Same architecture as before, except that each node is **randomly** connected to k units in the previous layer and $\widehat{w}_{ij}^l = 1$.

Such layered networks can implement a large class of Boolean functions.

[A. Mozeika and D. Saad PRL 2009]

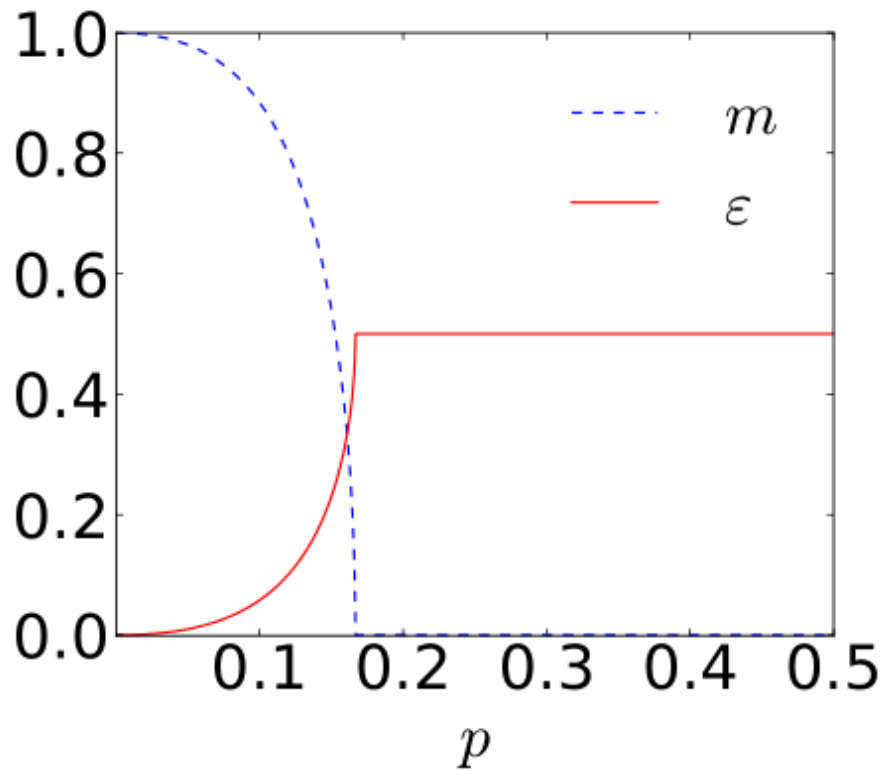
- In addition to the overlap q^l the magnetization $m^l := 1/N \sum_i \langle s_i^l \rangle$ order parameter characterizes the macroscopic dynamics.

$$m^l = \sum_{\{s_j\}} \prod_{j=1}^k \frac{1}{2} [1 + s_j m^{l-1} (1 - 2p)] \operatorname{sgn} \left[\sum_{j=1}^k s_j \right]$$

$$q^l = \sum_{\{s_j, \hat{s}_j\}} \prod_{j=1}^k \frac{1}{4} [1 + \hat{s}_j \hat{m}^{l-1} + s_j m^{l-1} (1 - 2p) + s_j \hat{s}_j q^{l-1} (1 - 2p)] \operatorname{sgn} \left[\sum_{j=1}^k \hat{s}_j \right] \operatorname{sgn} \left[\sum_{j=1}^k s_j \right]$$

Sparsely Connected Binary Networks

For the same for perturbation at every layer $p^l = p$ (weight flipping probability) and at the infinite depth limit $L \rightarrow \infty$



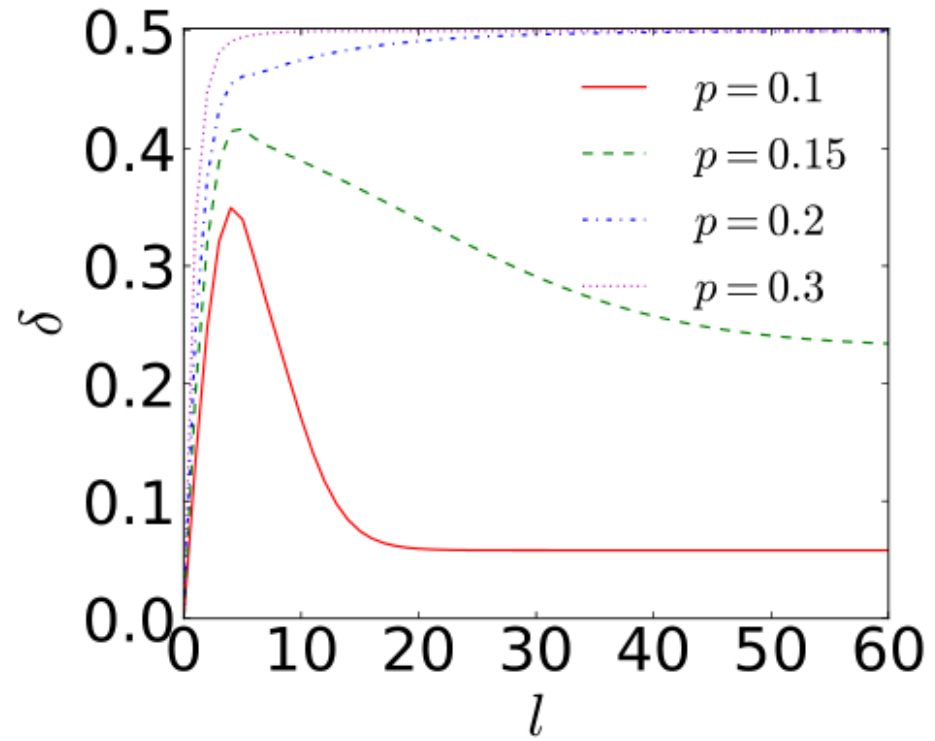
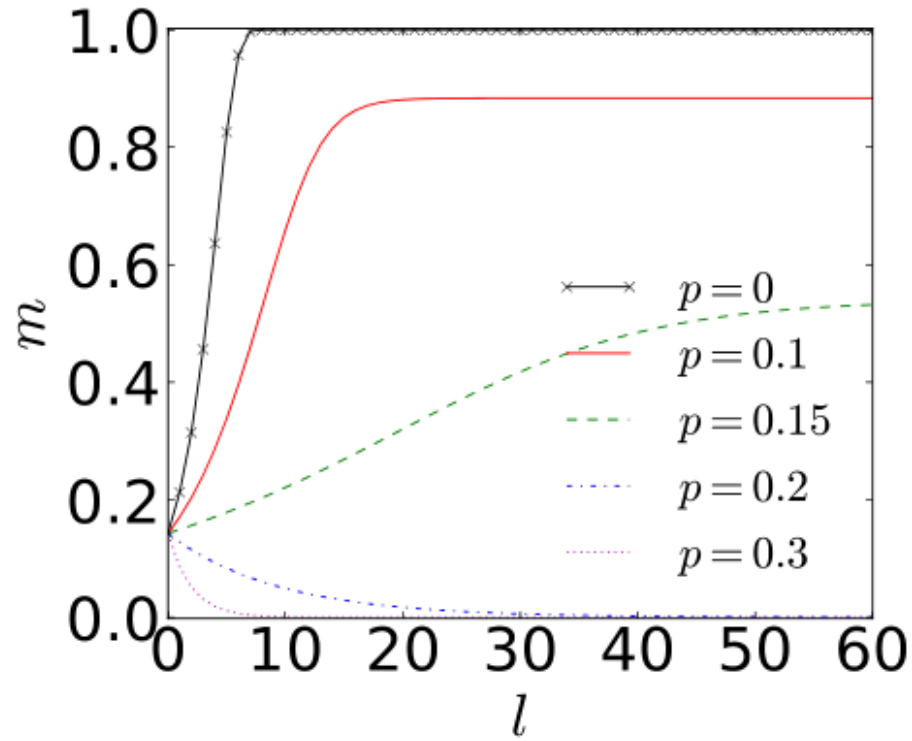
Reference function $f_{\hat{w}} = \text{sgn}(\sum \hat{s}_i^0)$ is majority vote of input.

Phase transition of **stationary states** as $L \rightarrow \infty$; $k = 3, m^0 > 0$.
($m^\infty < 0$ if $m^0 < 0$)

Similar to the phase transition of varying thermal noise β in the noisy computation setting.

[A. Mozeika and D. Saad PRL 2009]

Deep Layers For Reliable Computation in Sparse Binary Deep Learning Machines

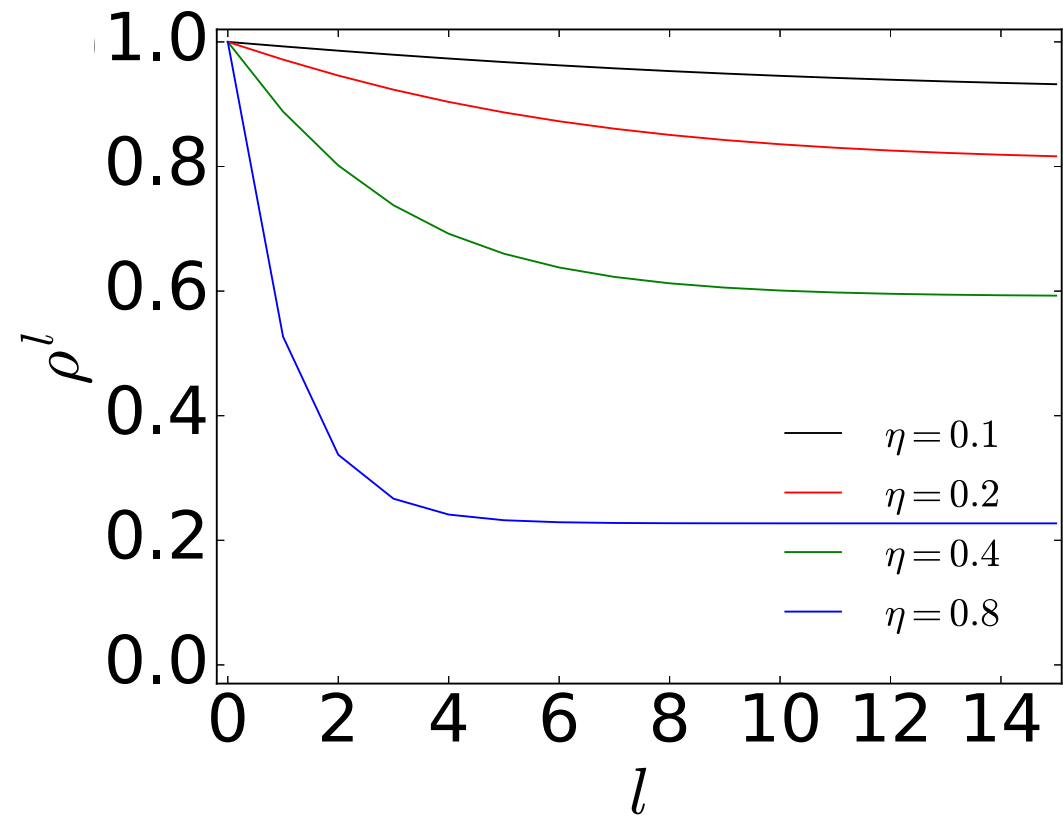
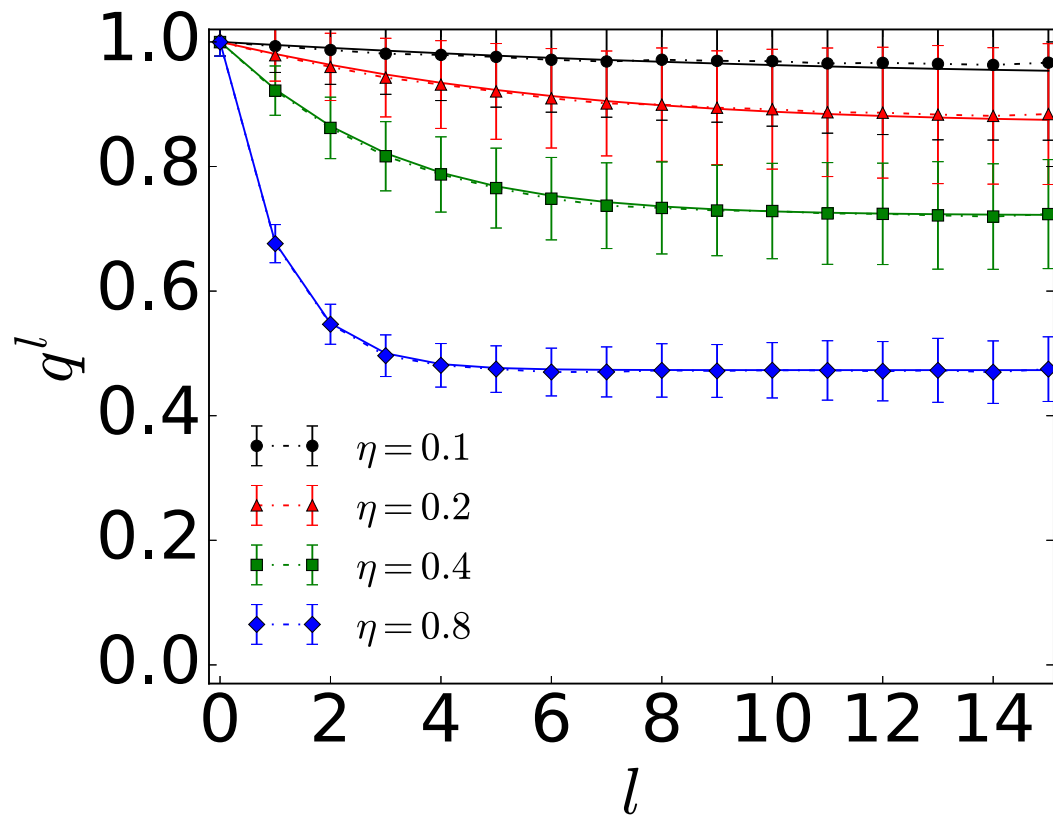


$$\delta^l := \frac{1}{2}(1 - q^l)$$

Internal error of activations

Continuous Variables and Activation Functions

We extended the framework for the ReLu function $\phi(x) = \begin{cases} x & x > 0 \\ 0 & x \leq 0 \end{cases}$



Summary and Future Work

Summary:

- The **generating functional analysis** is a principle method for deriving the **typical** behaviors of DLMs; it allows for non-trivial extensions.
- **Layer-by-layer matching** of weights is observed when getting closer to the reference function in densely-connected networks.
- Sparsely connected networks favor **deep layers** for a reliable representation.

Future work:

- A model learned from data $p(\mathbf{w}|D)$ – Typically many **redundant** weights.
- Exploring the function space for **correlated inputs**.
- The role of **over-parametrization** in the function landscape, error and generalization.
- Other models? Optimizing **variable hidden layer size**.
- **Noisy** training/computation for better generalization.