

Exponential Capacity in an Autoencoder Neural Network with a Hidden Layer

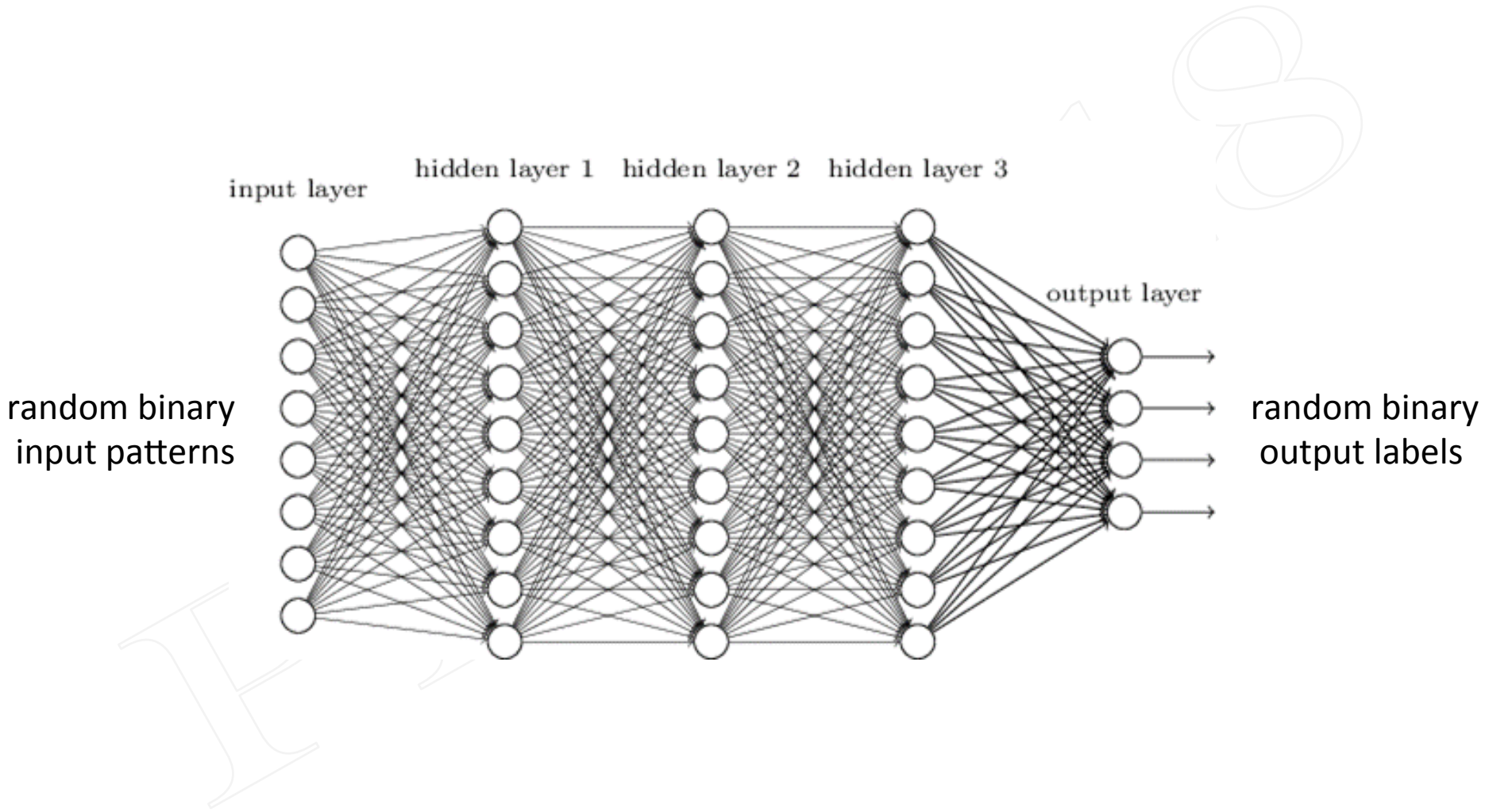
Alireza Alemi*, Alia Abbara



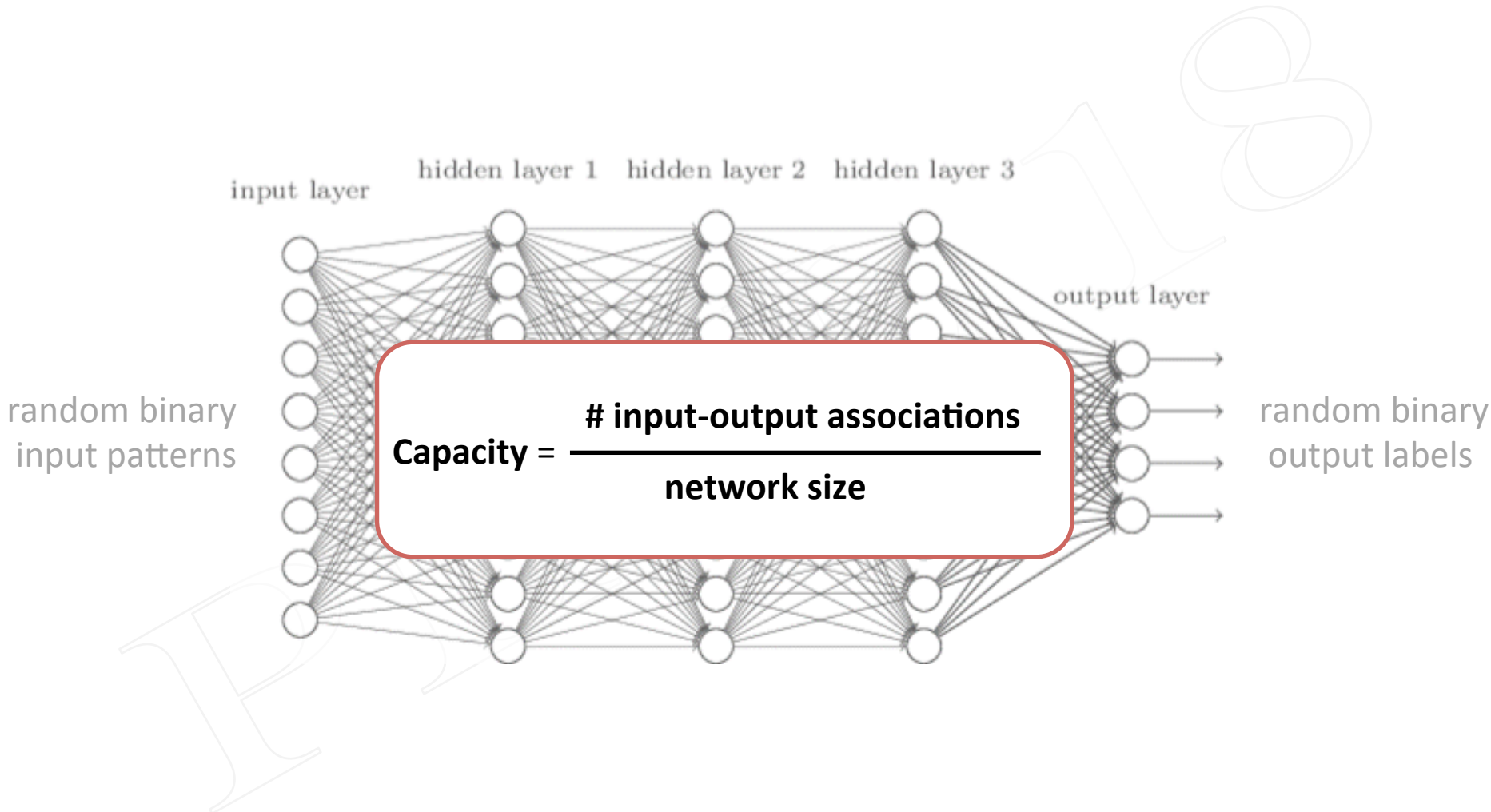
PHYSICS INFORMED MACHINE LEARNING

Santa Fe, NM, Jan 2018

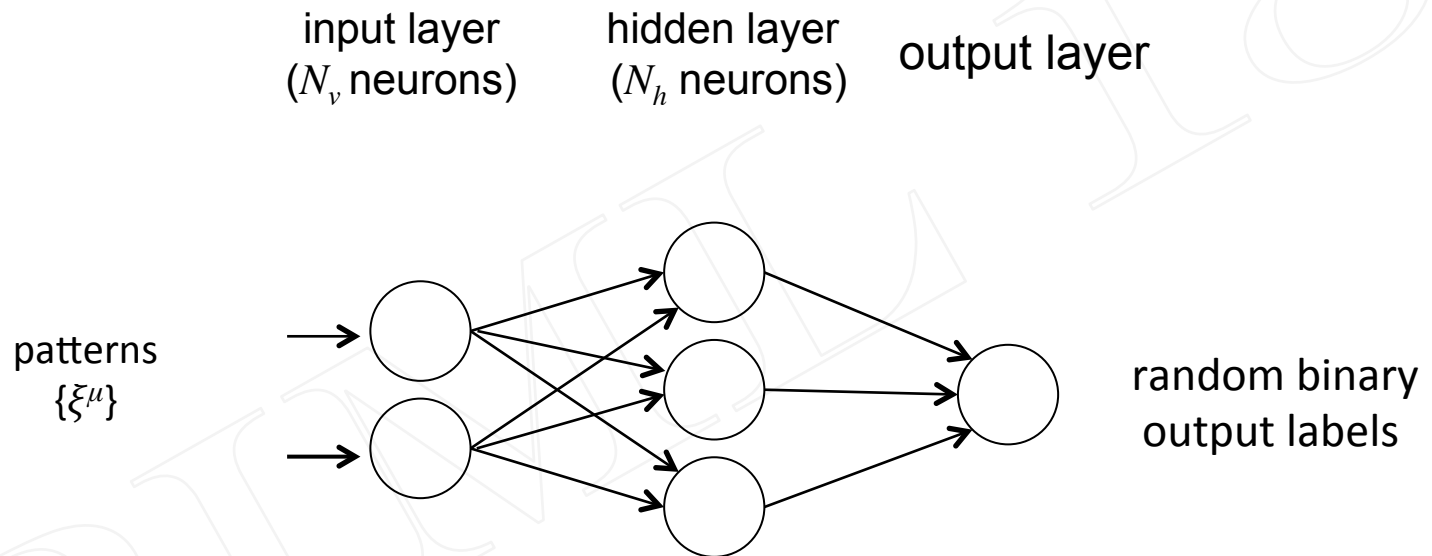
Limits of computations in deep nets?



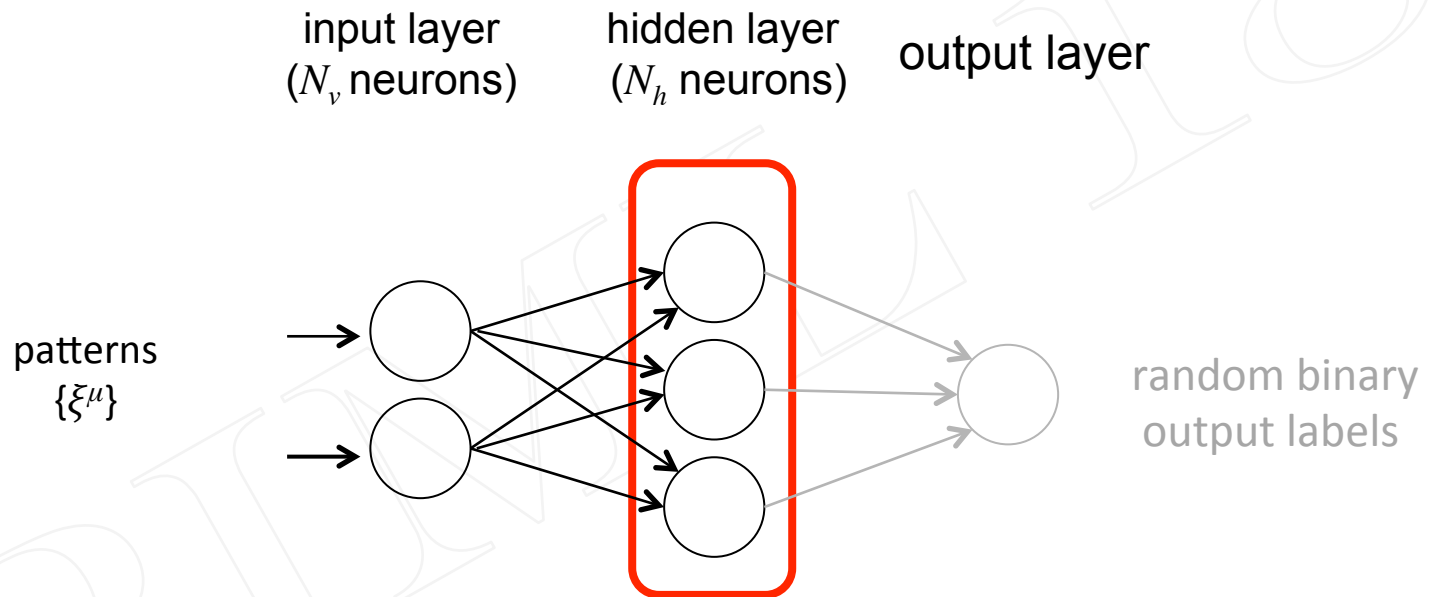
Limits of computations in deep nets?



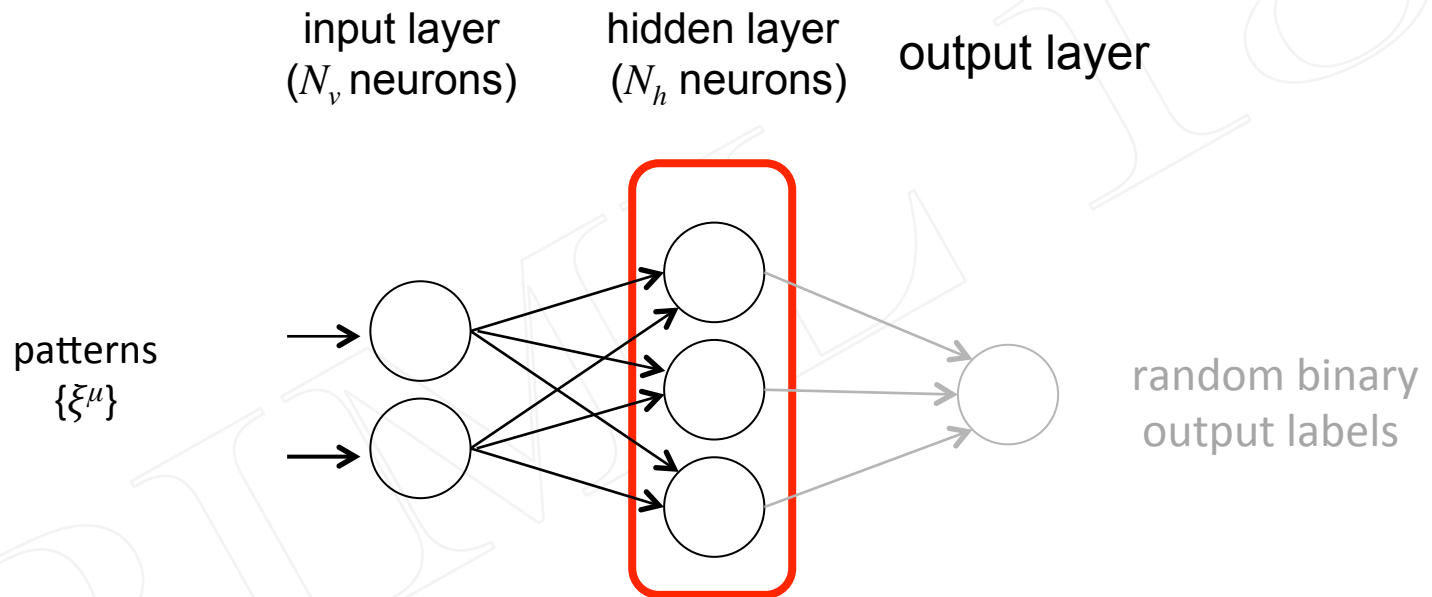
Capacity of multi-layer perceptron (MLP) with one hidden layer?



Investigating internal representation



Expansive (or *overcomplete*) representation



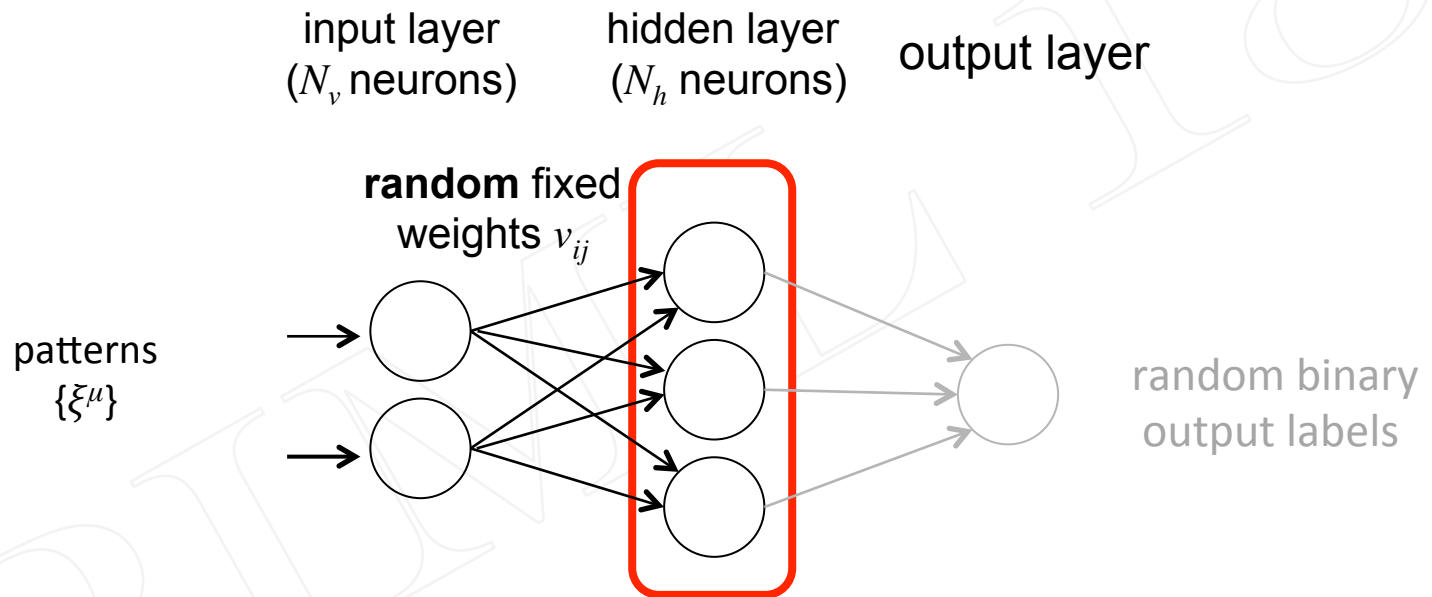
$$\Lambda = \frac{N_h}{N_v} > 1$$

Expansive (or *overcomplete*) representation

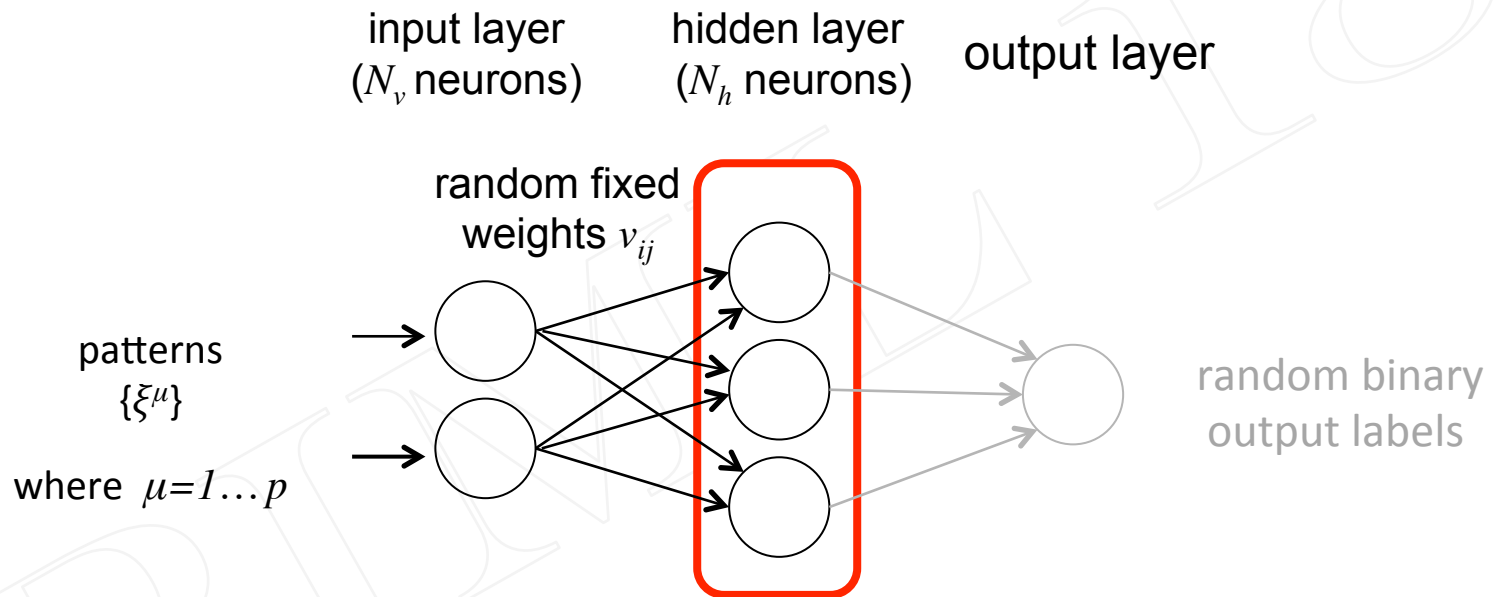


Nervous system	Expansion ratio (Λ)
Entorhinal cortex \rightarrow Dentate gyrus	~ 10
LGN \rightarrow V1	~ 25
in fly olfactory system: antennal lobe \rightarrow mushroom body	~ 40
in the cerebellum, mossy fibers \rightarrow granule cells	100 \sim 200
rodent's olfactory bulb \rightarrow piriform cortex	$\sim 10^3$

Expansive representation with random weights

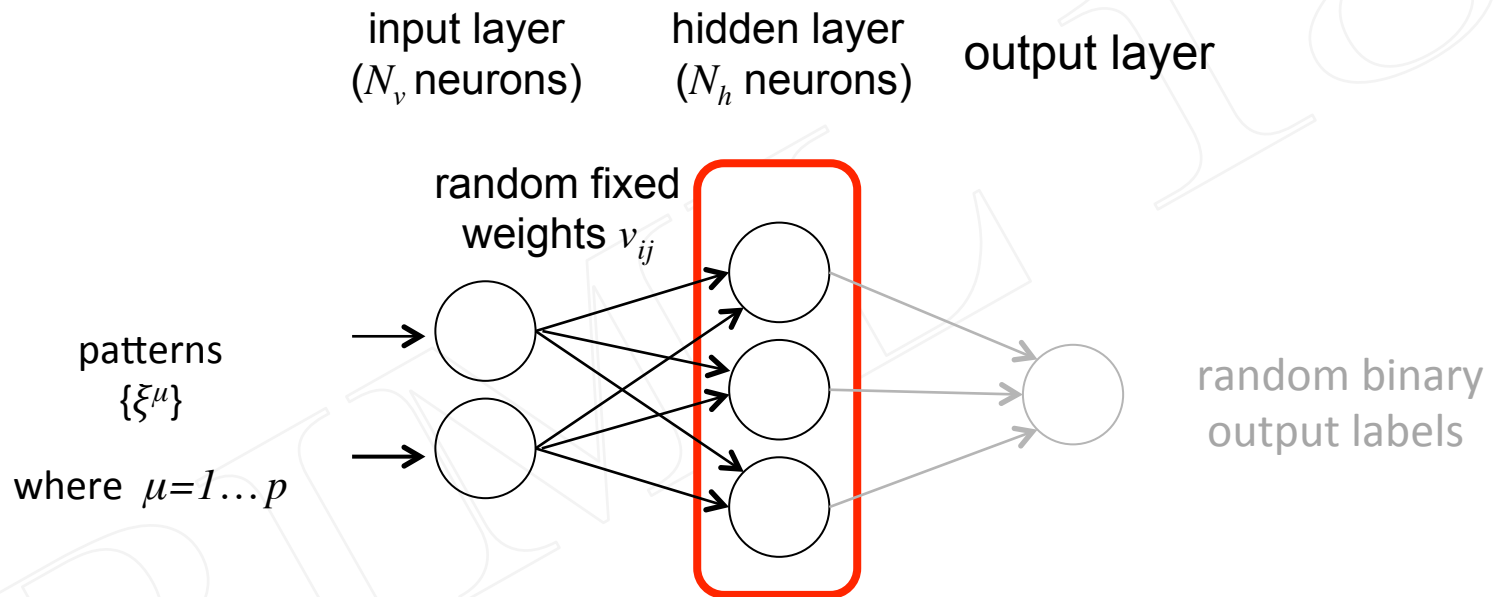


Expansive representation with random weights



How good is random-weight expansive representations?

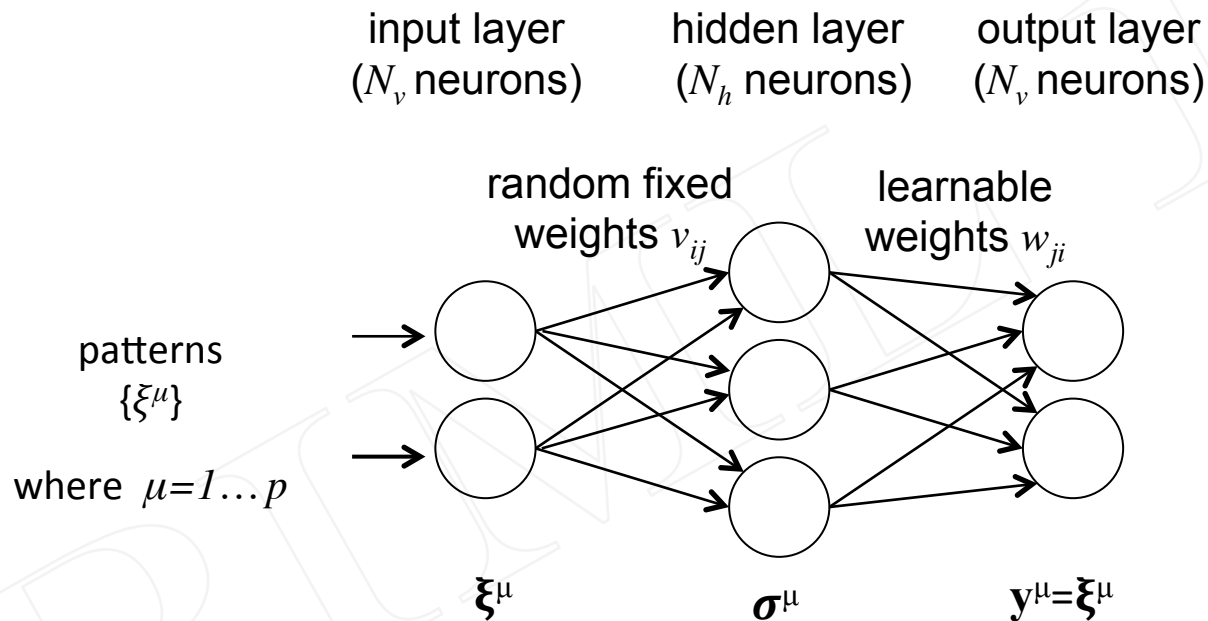
Expansive representation with random weights



How good is random
expansive representation?

How many of input patterns can be
reconstructed from this representation?

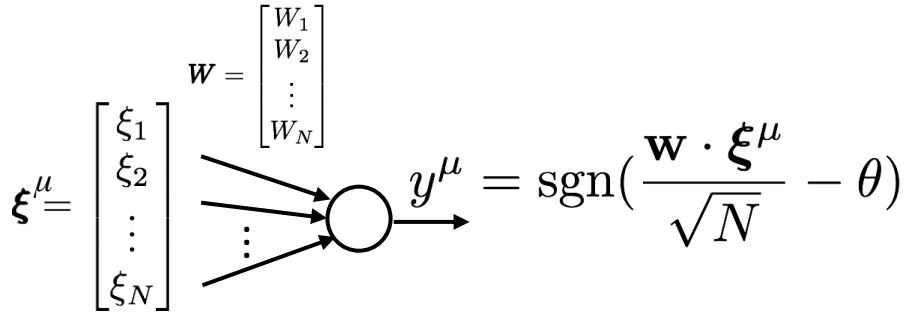
Expansive autoencoder with random weights



Focus of this talk:

Optimal capacity of expansive autoencoder?

Critical capacity of the perceptron: Gardner approach



- Find \mathbf{w} to store a set of random ensemble of input-output pairs $\{ (\xi^\mu, d^\mu) \}$, $\mu=1 \dots p$ with a robustness κ

- Storage criteria: $\forall \mu, i : d^\mu \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N w_i \xi_i^\mu - \theta \right) > \kappa$

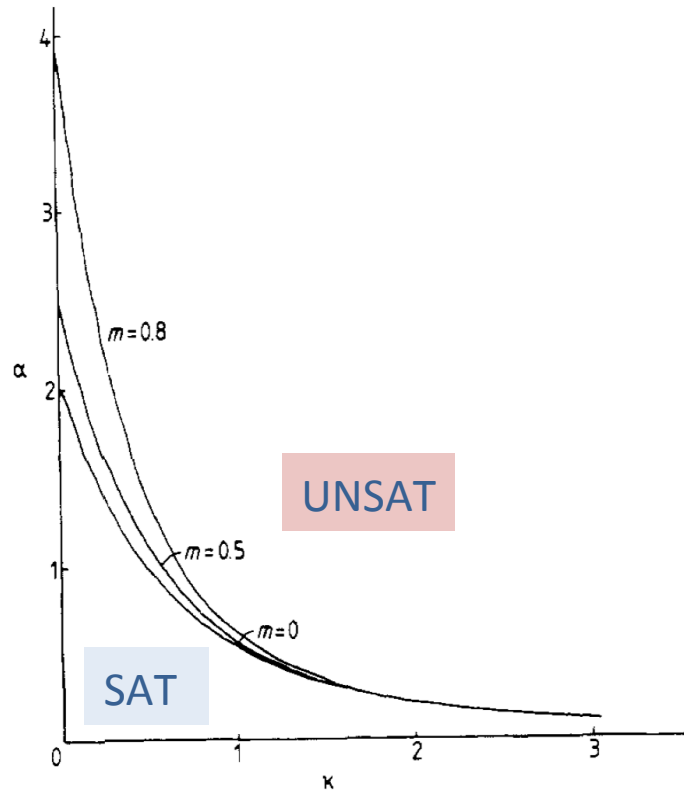
Gardner volume ($N \rightarrow \infty$):

$$\Omega = \int_{\|\mathbf{w}\|^2=N} d^N \mathbf{w} \prod_{\mu=1}^p \Theta \left(d^\mu \left(\frac{1}{\sqrt{N}} \sum_{i=1}^N w_i \xi_i^\mu - \theta \right) - \kappa \right)$$

$\langle\langle \log(\Omega) \rangle\rangle$ averaging over the (quenched) distribution of the patterns

$$\langle\langle \log \Omega \rangle\rangle = \lim_{n \rightarrow 0} \frac{\langle\langle \Omega^n \rangle\rangle - 1}{n}, \quad \text{critical capacity: } \alpha_c = \lim_{N \rightarrow \infty} \frac{p_c}{N}$$

- Perceptron learning rule ($\kappa=0$): $\Delta w_i = \eta \xi_i (d^\mu - y^\mu)$



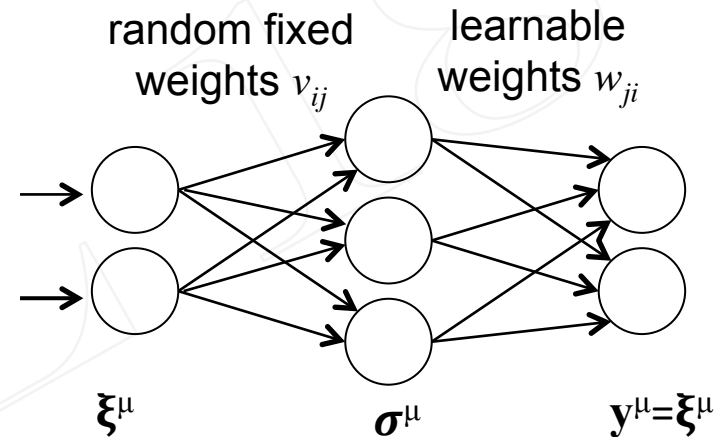
Gardner E. (1988) J. Phys. A: Math. Gen.

A simplified expansive (over-complete) autoencoder

- Binary neurons (± 1)
- Expansion ratio: $\Lambda = N_h / N_v$
- Fixed encoding weights v_{ij} are randomly sampled from $\mathcal{N}(0, 1)$
- Goal: reconstruct random patterns $\{\xi^\mu\}_{\mu=1, \dots, p}$ by learning the decoding weights w_{ji}
- Perceptron learning rule: $\Delta w_{ji} = \eta (\xi_j^\mu - y_j^\mu) \sigma_i^\mu$
- Dense coding level for patterns: $P(\xi_i = +1) = 0.5$
- Coding level of hidden layer: $P(\sigma_i = +1) = f$
- Maximal/critical storage Capacity:
 $\alpha_{max} = p_{max} / N_h$ as $N_h \rightarrow \infty$
- Storage criteria: meeting fixed-point equations

$$\forall j, \mu: \xi_j^\mu = \text{sgn} \left(\sum_i w_{ji} \text{sgn} \left(\sum_l v_{il} \xi_l^\mu \right) \right)$$

input layer (N_v neurons) hidden layer (N_h neurons) output layer (N_v neurons)



Dynamics:

$$\sigma_i^\mu = \text{sgn} \left(\sum_{j=1}^{N_v} v_{ij} \xi_j^\mu - \theta \right)$$

$$y_j^\mu = \text{sgn} \left(\sum_{i=1}^{N_h} w_{ji} \sigma_i^\mu \right)$$

A mean-field approximation (MFA)

$$\xi_j^\mu = \text{sgn} \left(\sum_{i=1}^{N_h} w_{ji} \text{sgn} \left(\sum_{l=1}^{N_v} v_{il} \xi_l^\mu \right) \right)$$

Approximated by a *quenched* random variable with a Gaussian distribution : $\mathcal{N}(0, N_v)$

$$= \text{sgn} \left(\sum_{i=1}^{N_h} w_{ji} \text{sgn} \left(\sum_{l \neq j; l=1}^{N_v} v_{il} \xi_l^\mu + v_{ij} \xi_j^\mu \right) \right)$$

$$= \text{sgn} \left(\sum_{i=1}^{N_h} w_{ji} \text{sgn} \left(z_i^\mu + v_{ij} \xi_j^\mu \right) \right)$$

$$P(\sigma_i^\mu | \xi_j^\mu) \simeq \frac{1}{2} + \sigma_i^\mu \xi_j^\mu \frac{v_{ij}}{\sqrt{2\pi N_v}}$$

$$(\sigma_i^\mu \perp \sigma_k^\mu) | \xi_j^\mu \implies P(\sigma^\mu | \xi_j^\mu) = \prod_i P(\sigma_i^\mu | \xi_j^\mu)$$

A mean-field approximation (MFA)

Probability distribution

$$P(\sigma_i^\mu | \xi_j^\mu) \simeq \frac{1}{2} + \sigma_i^\mu \xi_j^\mu \frac{v_{ij}}{\sqrt{2\pi N_v}}$$

$$(\sigma_i^\mu \perp \sigma_k^\mu) | \xi_j^\mu \implies P(\sigma^\mu | \xi_j^\mu) = \prod_i P(\sigma_i^\mu | \xi_j^\mu)$$

Mean-Field Approximation (MFA) for the *replica* calculation

Gardner volume:

$$\Omega = \int_{\|\mathbf{w}\|^2=N_h} d^{N_h} \mathbf{w} \prod_{\mu=1}^p \Theta\left(\xi_j^\mu \frac{1}{\sqrt{N_h}} \sum_{i=1}^{N_h} w_{ji} \sigma_i^\mu - \kappa\right)$$

Replica theory: compute $\langle\langle \log(\Omega) \rangle\rangle$ over distribution of patterns using the replica method

$$\langle\langle \log \Omega \rangle\rangle = \lim_{n \rightarrow 0} \frac{\langle\langle \Omega^n \rangle\rangle - 1}{n}$$

Probability distribution

$$P(\sigma_i^\mu | \xi_j^\mu) \simeq \frac{1}{2} + \sigma_i^\mu \xi_j^\mu \frac{v_{ij}}{\sqrt{2\pi N_v}}$$

$$(\sigma_i^\mu \perp \sigma_k^\mu) | \xi_j^\mu \implies P(\sigma^\mu | \xi_j^\mu) = \prod_i P(\sigma_i^\mu | \xi_j^\mu)$$

Exponential critical capacity in the MFA

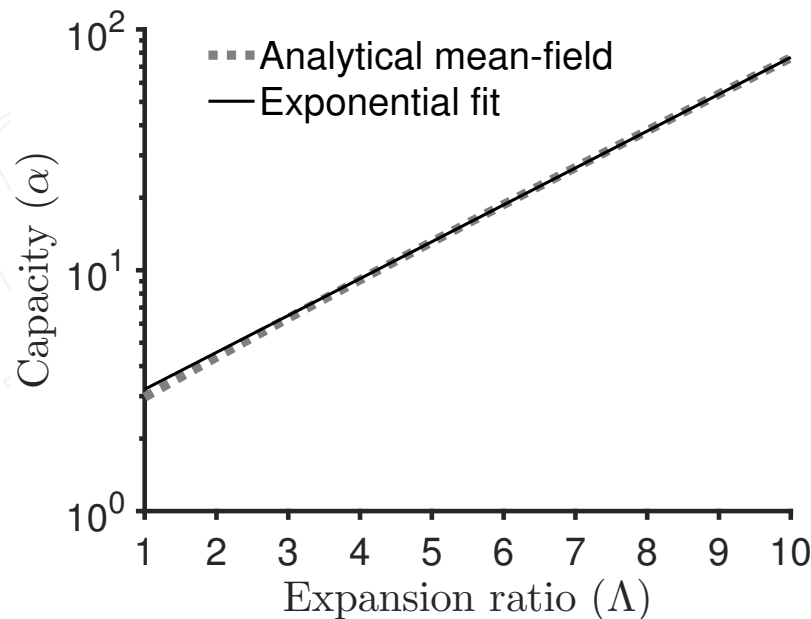
Solutions (replica-symmetric, locally stable):

$$\alpha_c^{\text{MFA}} \int_{M\sqrt{2\Lambda/\pi-\kappa}}^{\infty} Dt (\kappa + t - \sqrt{2\Lambda/\pi})^2 = 1 - M^2$$

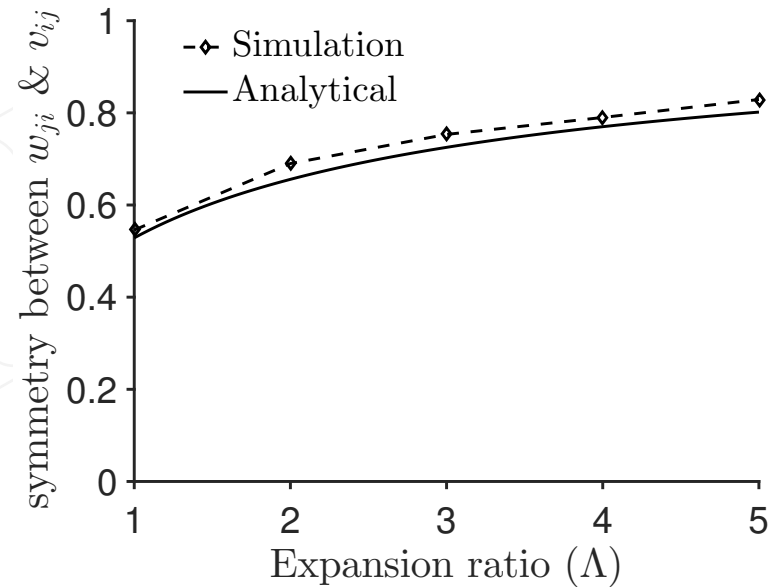
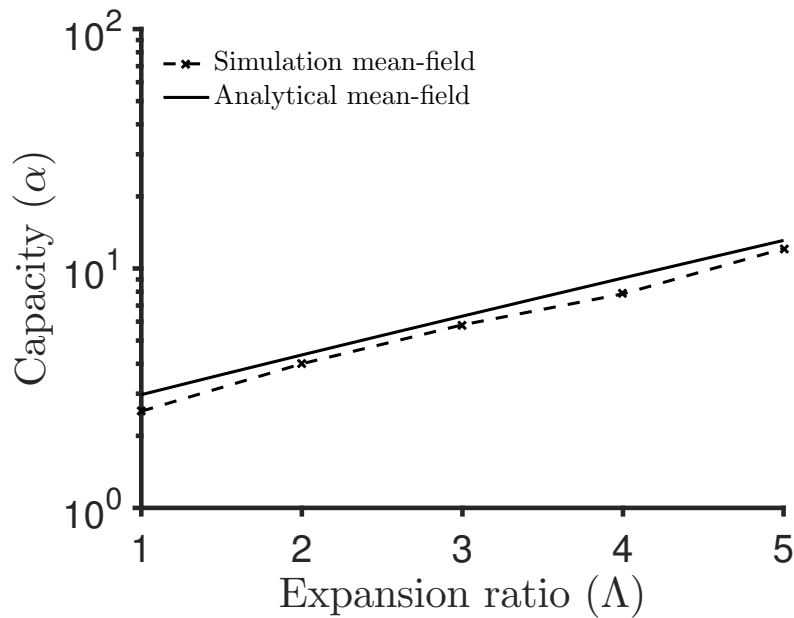
$$\alpha_c^{\text{MFA}} \int_{M\sqrt{2\Lambda/\pi-\kappa}}^{\infty} Dt (\kappa + t - \sqrt{2\Lambda/\pi}) \sqrt{2\Lambda/\pi} = M,$$

where $M = \sum_i \frac{v_i w_i}{N_h}$, $\Lambda = \frac{N_h}{N_v}$

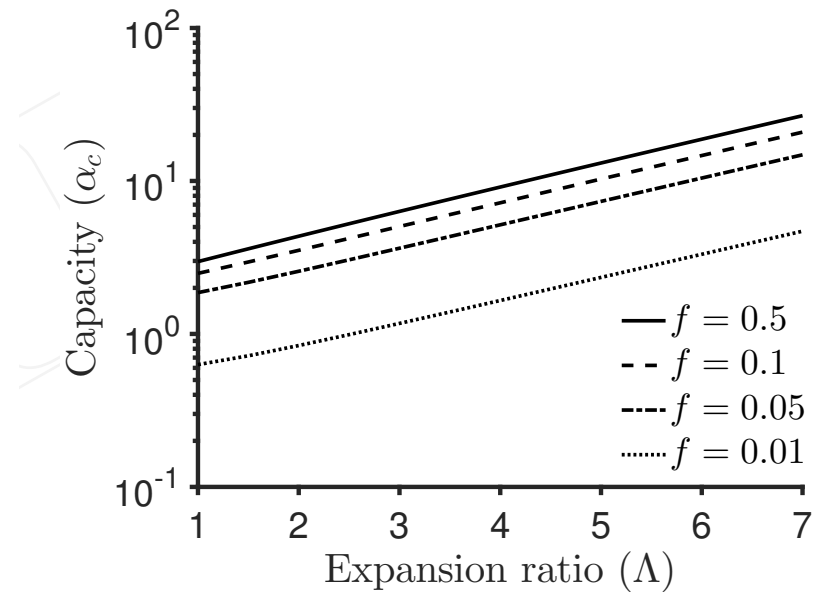
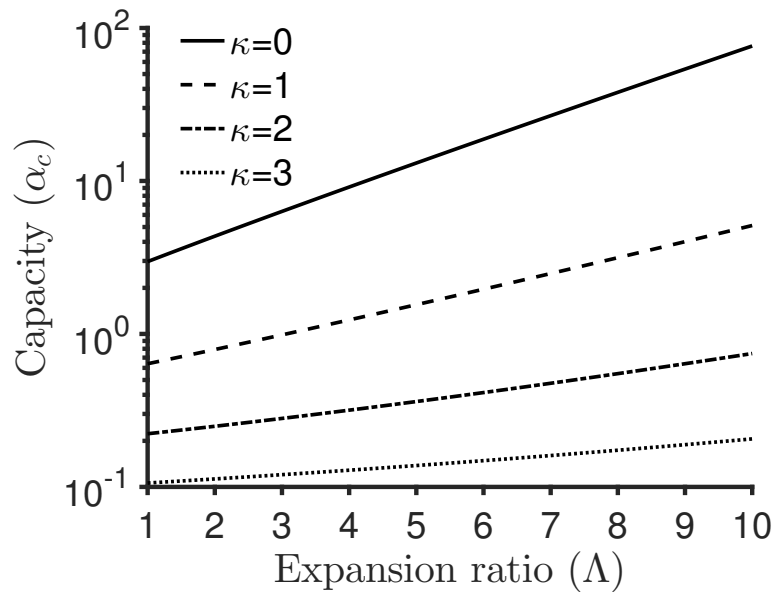
$$Dt \equiv \frac{dt}{\sqrt{2\pi}} e^{-\frac{t^2}{2}}$$



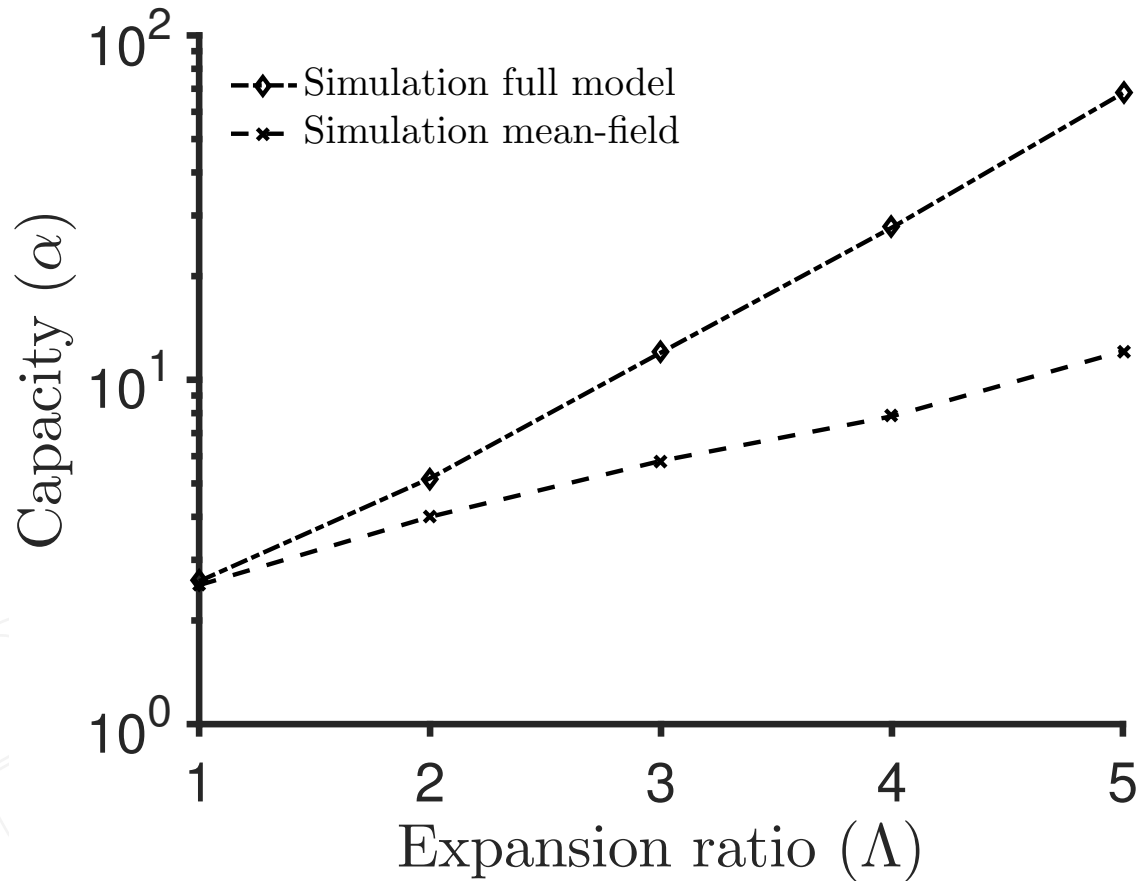
Comparison with simulation (perceptron learning rule)



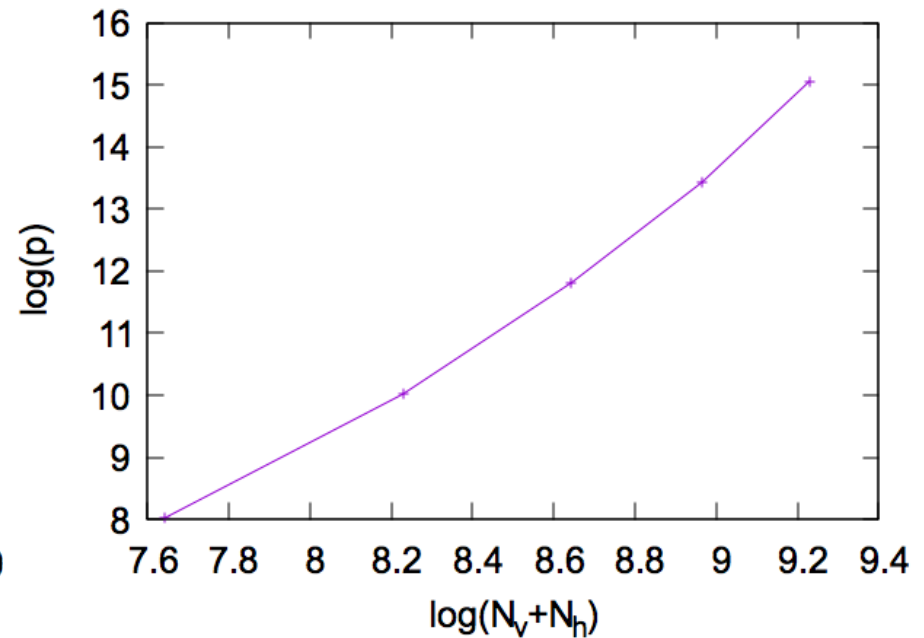
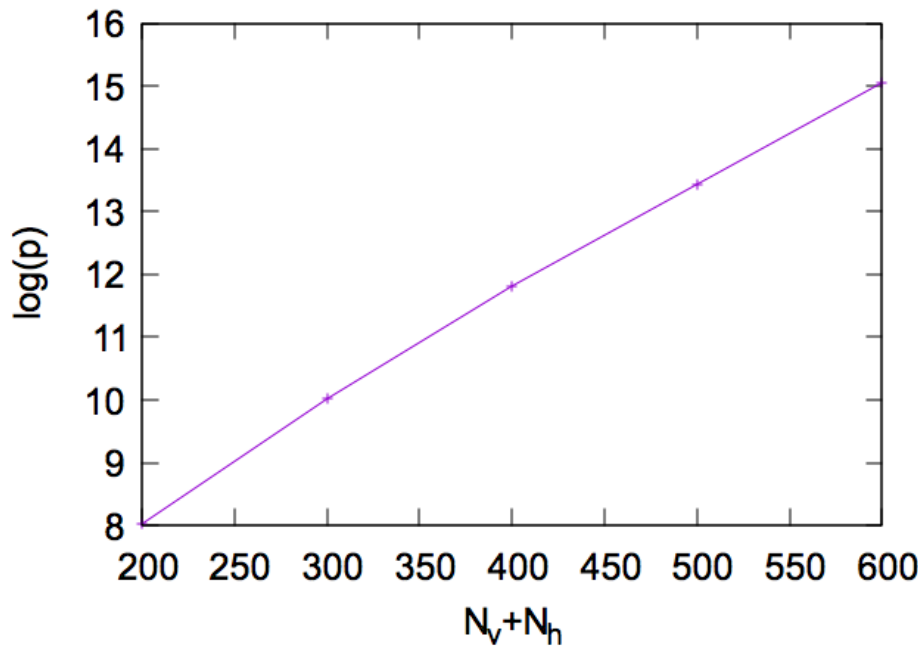
Effect of robustness of output, sparseness in the hidden layer



Comparison with the full-model



Log of number of patterns vs. total number of neurons



Discussion

- Expansive representation has important computational implications.
- *Exponential* scaling of the capacity with the expansion ratio in our autoencoder
 - Expansion makes up for the loss of information due to binarization
- In very good agreement with results of simulations using an online learning rule
- Future directions
 - Robustness to noise by training encoding weights
 - Computing the capacity of the full autoencoder model
 - Extension to the capacity of MLP and deep nets

Acknowledgement

Alia Abbara

Discussion with many colleagues ...



SIMONS FOUNDATION

Question?

PIML 2018