

Sample complexity of learning Convolutional and Recurrent NNs

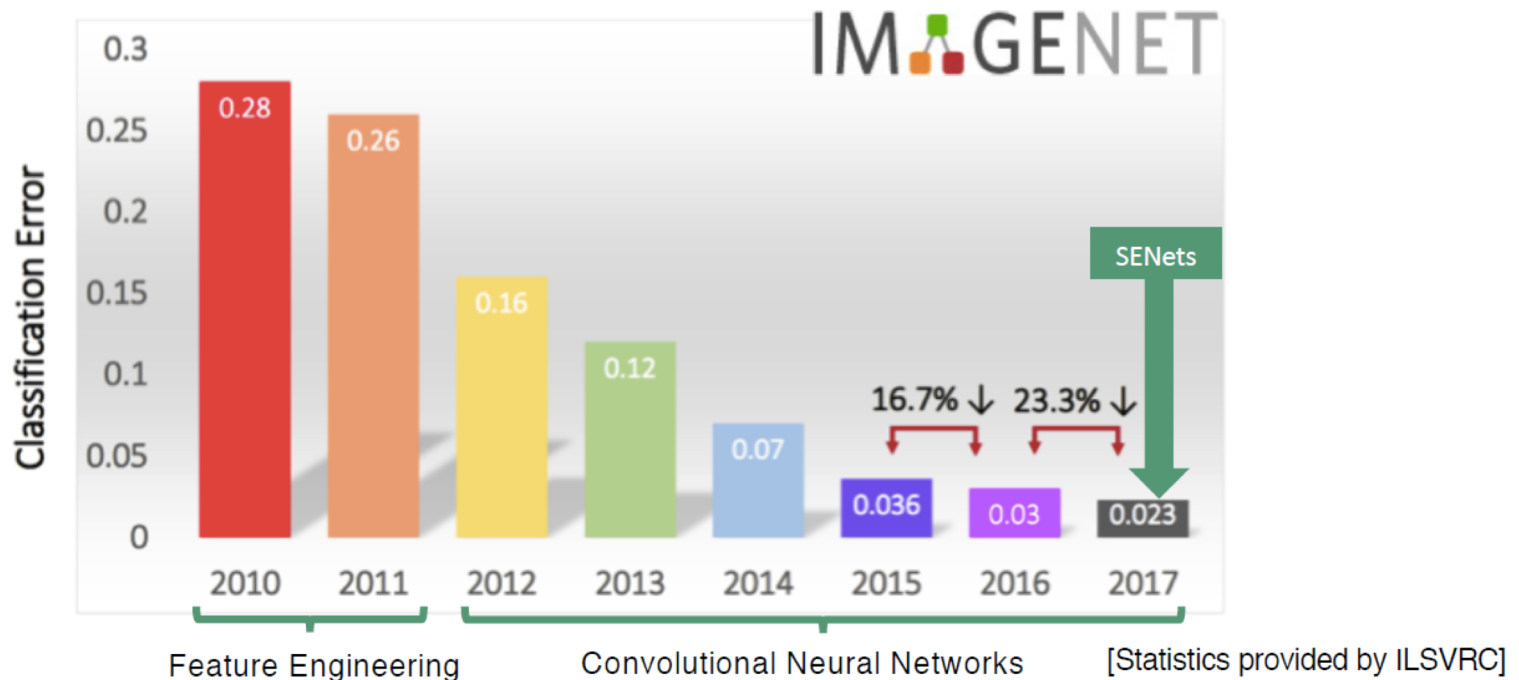
Aarti Singh
Associate Professor

Joint work with Y. Wang, S. Du, X. Zhai,
S. Balakrishnan, R. Salakhutdinov

Simons workshop on Frontiers of Deep Learning
July 2019

CNNs and RNNs

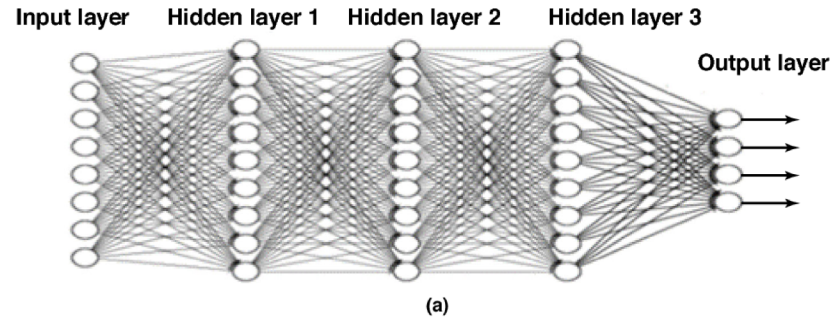
- Large part of the recent success of NNs, particularly for spatial image data, is due to Convolution Neural Network (CNN) architectures (LeNet, AlexNet, VGG, GoogLeNet, ResNet, ...)



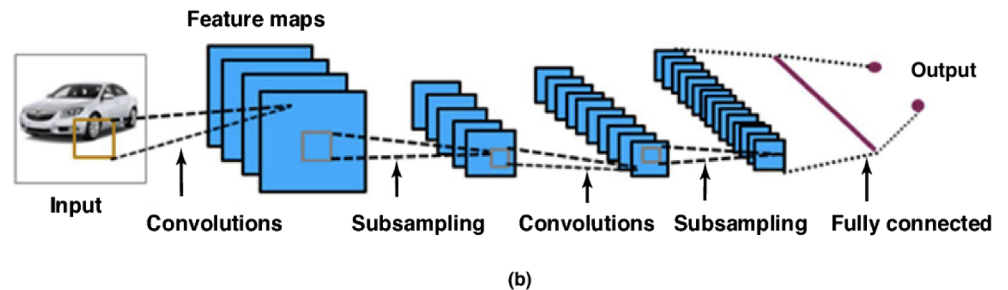
- Corresponding analogue for temporal or sequential data is the Recurrent Neural Network (RNN) architecture

FNN, CNN and RNN architectures

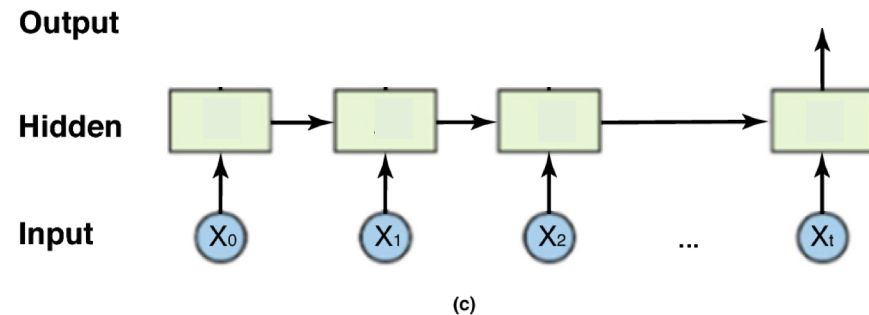
- F(Fully-connected)NN



- C(Convolutional)NN



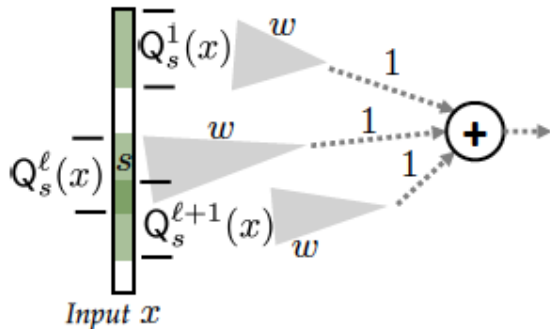
- R(Recurrent)NN



CNN generative models

$$Y^i = F(X^i; \theta) + \xi_i \quad X^i \in \mathbb{R}^d \quad \{X^i\}_{i=1}^n \stackrel{\text{i.i.d}}{\sim} \mu$$

CNN with Average pooling $F^{\text{CA}}(X^i; w) = \sum_{\ell=0}^{\lfloor (d-m)/s \rfloor} w^\top Q_s^\ell(x^i)$

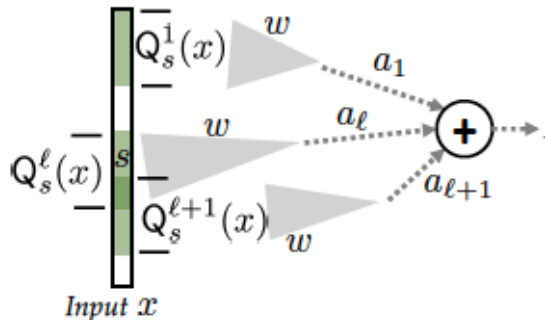


m – size of filter s – stride of filter

$$Q_s^\ell(x^i) = (x_{\ell s+1}^i, \dots, x_{\ell s+m}^i)$$

length m segment of input

CNN with Weighted pooling $F^{\text{CW}}(X^i; w, a) = \sum_{\ell=0}^{\lfloor (d-m)/s \rfloor} a_\ell w^\top Q_s^\ell(x^i)$



Size of output layer, $J = \lfloor (d-m)/s \rfloor + 1$

RNN generative model

$$Y^i = F(X^i; \theta) + \xi_i \quad X^i \in \mathbb{R}^d \quad \{X^i\}_{i=1}^n \stackrel{\text{i.i.d}}{\sim} \mu$$

RNN $F^R(X^i, A, B) = \mathbf{1}^\top h_L^i$

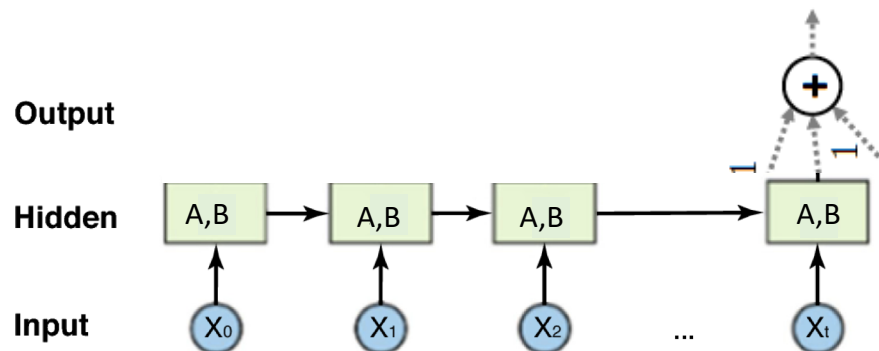
$$h_t^i = Ah_{t-1}^i + Bx_t^i, \quad t = 1, 2, \dots, L, \quad \text{Initial hidden state, } h_0^i = \mathbf{0}$$

$$A \in \mathbb{R}^{r \times r}$$

L – length of input

$$B \in \mathbb{R}^{r \times d}$$

r – hidden state dim



Minimax analysis

- Model may be non-identifiable (parameters not unique)
E.g. w , a scaling for CNN or exchange hidden units in RNN

Focus on mean-square prediction error

$$\text{err}(\hat{\theta}, \theta) := \mathbb{E}_{\mu} |F(x; \theta) - F(x; \hat{\theta})|^2$$

Goal: Upper and lower bound **Minimax risk**

$$\mathfrak{M}(n; F) := \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}_{\mu, \theta} [\text{err}(\hat{\theta}, \theta)]$$

Estimator and Assumptions

Least Squares Estimator

$$\hat{\theta} \in \arg \min_{\theta \in \Theta} \sum_{i=1}^n |Y^i - F(X^i; \theta)|^2$$

- May be non-unique, guarantees apply to any global minimizer
- Ignore computational considerations

Assumptions

A1) Noise is independent centered sub-gaussian (σ^2)

A2) Input distribution μ is centered sub-gaussian with

$$cI \preceq \mathbb{E}_{\mu}[xx^{\top}] \preceq CI$$

Main results (Informal)

$$\mathfrak{M}(n; F) := \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}_{\mu, \theta} [\text{err}(\hat{\theta}, \theta)]$$

CNN with Average pooling

$$\mathfrak{M}(n; F^{\text{CA}}) = \tilde{\Theta} \left(\frac{m}{n} \right)$$

Independent of
input dimension d
FNN $\sim d/n$

CNN with Weighted pooling

$$\mathfrak{M}(n; F^{\text{CW}}) = \tilde{\Theta} \left(\frac{m + J}{n} \right)$$

RNN

$$\mathfrak{M}(n; F^{\text{R}}) = \tilde{\Theta} \left(\frac{rd}{n} \right)$$

Independent of
sequence length L
FNN $\sim Ld/n$

Main results (Informal)

$$\mathfrak{M}(n; F) := \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}_{\mu, \theta} [\text{err}(\hat{\theta}, \theta)]$$

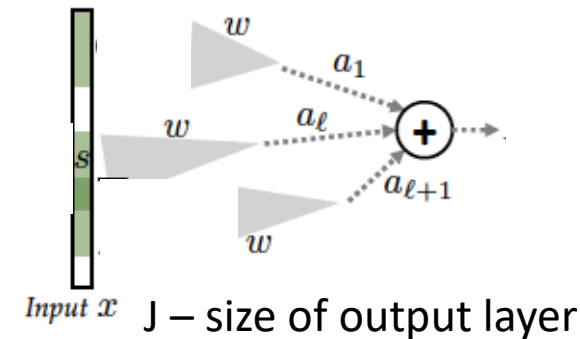
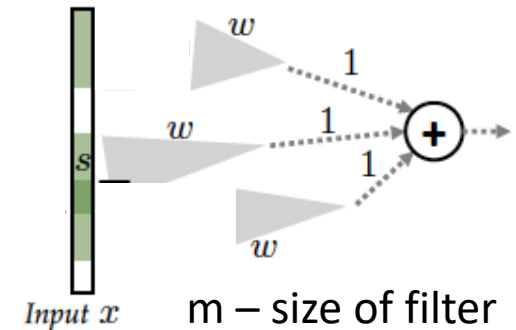
CNN with Average pooling

$$\mathfrak{M}(n; F^{\text{CA}}) = \tilde{\Theta} \left(\frac{m}{n} \right)$$

CNN with Weighted pooling

$$\mathfrak{M}(n; F^{\text{CW}}) = \tilde{\Theta} \left(\frac{m + J}{n} \right)$$

Match parameter count



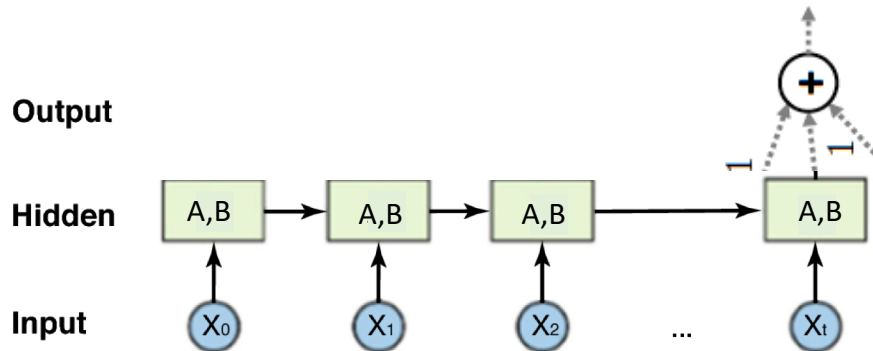
Main results (Informal)

$$\mathfrak{M}(n; F) := \inf_{\hat{\theta}} \sup_{\theta} \mathbb{E}_{\mu, \theta} [\text{err}(\hat{\theta}, \theta)]$$

Match parameter
count

RNN

$$\mathfrak{M}(n; F^{\mathbb{R}}) = \tilde{\Theta} \left(\frac{rd}{n} \right)$$



$$h_t^i = Ah_{t-1}^i + Bx_t^i$$

$$A \in \mathbb{R}^{r \times r}$$

$$B \in \mathbb{R}^{r \times d}$$

Related work

Generalization bounds for NNs; some also apply to CNNs

Arora et al' 18, Anthony and Bartlett'09, Bartlett et al'17, Neyshabur et al'17, Konstantinos et al'17, Zhou and Feng'18, Li et al'18, Long-Sedghi'19...

$$L(\theta) - L_{\text{tr}}(\theta) \leq D/\sqrt{n}$$

- Fast rate – we show $1/n$ rates (under some assumptions)
- Scale independence – model complexity D typically depends on norm of parameters

High-dimensional linear regression ($d > n$) – above issues akin to sparsity based analysis e.g. using lasso

Related work

RNN model special case of classical (Kalman, 1960) problem of learning a linear dynamical system

Recent statistical and computational analysis (Hazan et al'17; Hardt et al'18; Simchowitz et al'18; Oymak and Ozay'18)

- Sample complexity not tight (to best of our knowledge)

Upper bounds (formal)

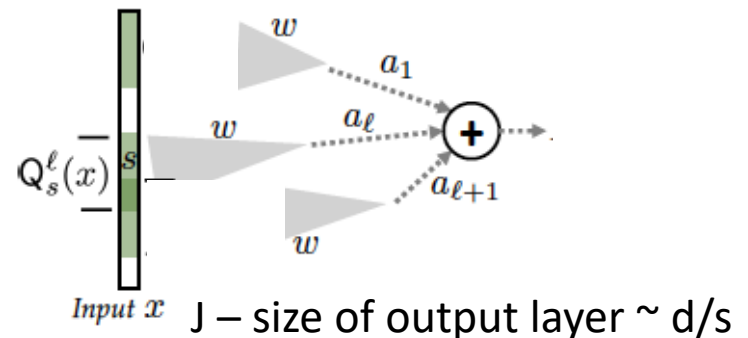
With probability $1-\delta$, for sufficiently large n ,

$$\mathfrak{M}(n; F^{\text{CA}}) \lesssim \frac{\sigma^2 m \log d}{n} \quad \sim m+J \quad \text{when } m/s \sim O(1)$$

$$\mathfrak{M}(n; F^{\text{CW}}) \lesssim \frac{\sigma^2 \min\{d, m + (d/s) \times (m/s)\} \cdot \log d}{n}$$

$$\mathfrak{M}(n; F^{\text{R}}) \lesssim \frac{\sigma^2 (d+L) \min\{r, d\} \log(Ld)}{n} \quad \sim rd \quad \text{when } r, L \ll d$$

- All bounds are achieved by the least squares estimator
- $n/\log^2 n \gtrsim$ numerator
- Match parameter counts



Proof sketch

For least-squares solution, $\frac{1}{n} \sum_{i=1}^n (Y_i - \langle X^i, \hat{\theta} \rangle)^2 \leq \frac{1}{n} \sum_{i=1}^n (Y_i - \langle X^i, \theta \rangle)^2$

Because of generative model $\|\hat{\theta} - \theta\|_X^2 \leq \frac{2}{n} \sum_{i=1}^n \xi_i \langle X^i, \hat{\theta} - \theta \rangle$

Self-normalized empirical process $\|\hat{\theta} - \theta\|_X \leq 2 \cdot \sup_{\phi \in \bar{\Theta}_X} \frac{1}{n} \sum_{i=1}^n \xi_i \langle X^i, \phi \rangle$

where $\bar{\Theta}_X := \{\phi = \theta - \theta' : \theta, \theta' \in \Theta, \|\phi\|_X \leq 1\}$

Dudley's integral upper bounds expectation of the process, and hence the error, in terms of covering number of $\bar{\Theta}_X$

relate to covering number of using restricted eigenvalues (ensured by A2)

$$\bar{\Theta}_2(\rho) := \{\phi = \theta - \theta' : \theta, \theta' \in \Theta, \|\phi\|_2 \leq \rho\}$$

$$\lambda_{\min}(\{X^i\}_{i=1}^n; \Phi) := \inf_{\phi \in \Phi} \|\phi\|_X^2 / \|\phi\|_2^2$$

$$\lambda_{\max}(\{X^i\}_{i=1}^n; \Phi) := \sup_{\phi \in \Phi} \|\phi\|_X^2 / \|\phi\|_2^2$$

Proof sketch

Lemma [Covering number of low-dim linear subspaces]: For any q , $k \leq q$, $\rho > 0$, and $\epsilon' \in (0, 1/2]$ there exists a finite set \mathcal{W} of k -dimensional subspaces in \mathbb{R}^q such that

for any k -dimensional subspace S in \mathbb{R}^q there exists a subspace $S' \in \mathcal{W}$ such that

$$\sup_{u \in S, \|u\|_2 \leq \rho} \inf_{v \in S', \|v\|_2 \leq \rho} \|u - v\|_2 \leq \epsilon'$$

And the size of the set $\log |\mathcal{W}| \lesssim kq \log(\rho q / \epsilon')$.

Example RNN: $\theta := (\mathbf{1}^\top A^{L-1} B \quad \mathbf{1}^\top A^{L-2} B \quad \dots \quad \mathbf{1}^\top B)$

L segments of d -dim, each of which lies in r -dim subspace

covering set of all $2r$ -dim subspaces in \mathbb{R}^d $O(rd \log d)$

covering of vectors in $2r$ -dim subspace for each $+ L \times O(r)$

Lower bounds (formal)

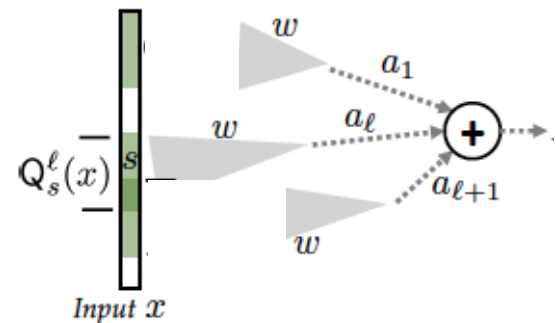
If input and noise distribution are standard normal, then there exists a universal constant $C > 0$ such that

$$\mathfrak{M}(n; F^{\text{CA}}) \geq C \frac{\sigma^2 m}{n}$$

$$\mathfrak{M}(n; F^{\text{CW}}) \geq C \frac{\sigma^2 (m + d/s)}{n} \quad \sim m+J$$

$$\mathfrak{M}(n; F^{\text{R}}) \geq C \frac{\sigma^2 \min\{rd, Ld\}}{n} \quad \sim rd \quad \text{since } r \ll L$$

- Bound holds for *any* estimator
- Lower bound for standard normal implies lower bound for general case
- Match parameter count



J – size of output layer $\sim d/s$

Proof sketch

Tsybakov extension of Fano's Lemma for Gaussian case

Corollary For any finite subset $\Theta' = \{\theta_0, \theta_1, \dots, \theta_M\} \subseteq \Theta$, denote $\rho_{\min} := \min_{j>0} \|\theta_0 - \theta_j\|_2/2$ and $\rho_{\text{avg}}^2 := \frac{1}{M} \sum_{i=1}^M \|\theta_i - \theta_0\|_2^2$. Then for any n ,

$$\inf_{\hat{\theta}_n} \sup_{\theta \in \Theta} \mathbb{E}_{\mu} [\|\hat{\theta}_n - \theta\|_2] \geq \rho_{\min} \times \frac{\sqrt{M}}{1 + \sqrt{M}} \left(1 - \frac{n\rho_{\text{avg}}^2}{\sigma^2 \log M} - 2\sqrt{\frac{n\rho_{\text{avg}}^2}{2\sigma^2 \log^2 M}} \right).$$

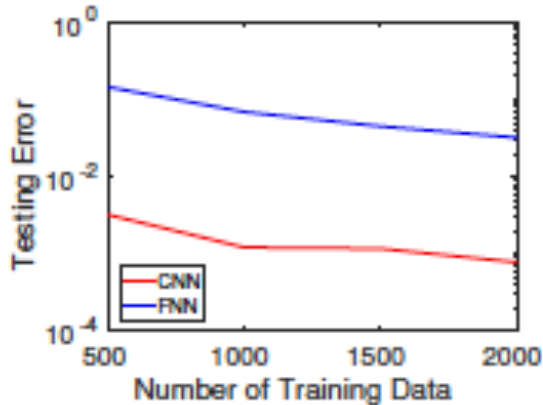
Characterization in terms of free parameters

Let $\Theta \subseteq \mathbb{R}^D, \mathcal{I} \subseteq [D]$. Suppose for any $u \in \mathbb{R}^{|\mathcal{I}|}$, there exists $\theta \in \Theta$ such that θ restricted to \mathcal{I} equals u . Then there exists a finite subset $\Theta' \subseteq \Theta$ as in Corollary , with $\log M \asymp |\mathcal{I}|$ and $\rho_{\min} \asymp \rho_{\text{avg}} \asymp \sqrt{|\mathcal{I}|}\epsilon$ for any $\epsilon > 0$.

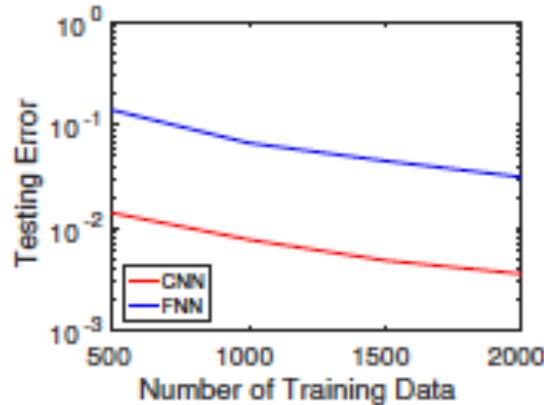
Experiments - CNN (average pooling) vs FNN

$s = m$ (non-overlapping)

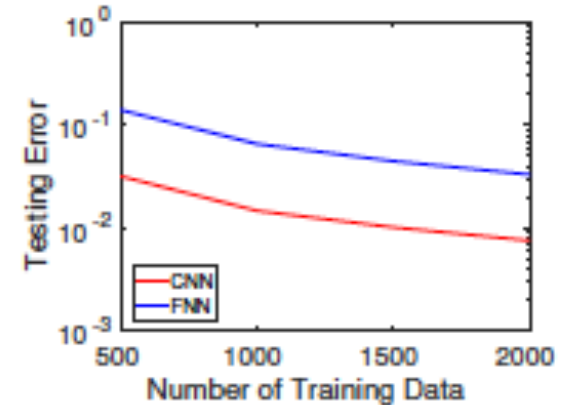
$d = 64$



(a) Filter size $m = 2$.

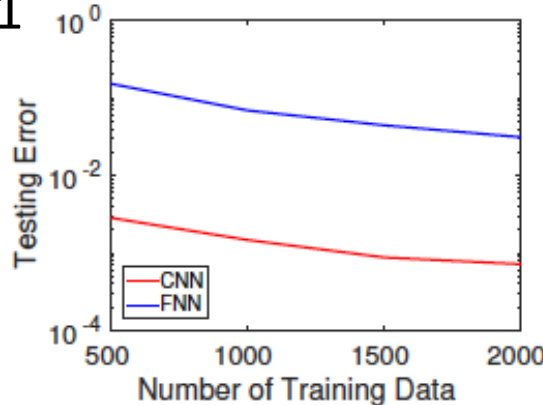


(b) Filter size $m = 8$.

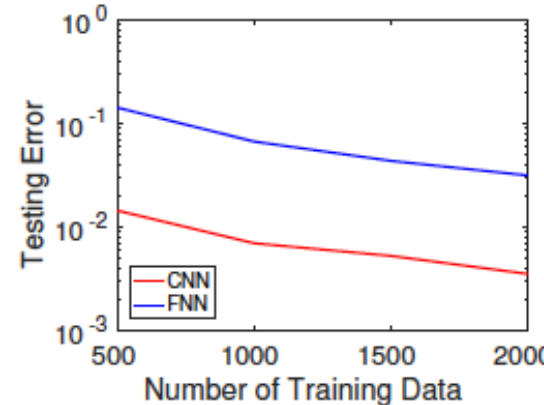


(c) Filter size $m = 16$.

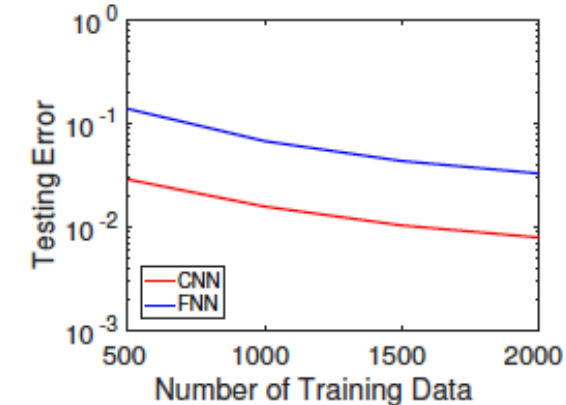
$s = 1$



(a) Filter size $m = 2$.



(b) Filter size $m = 8$.



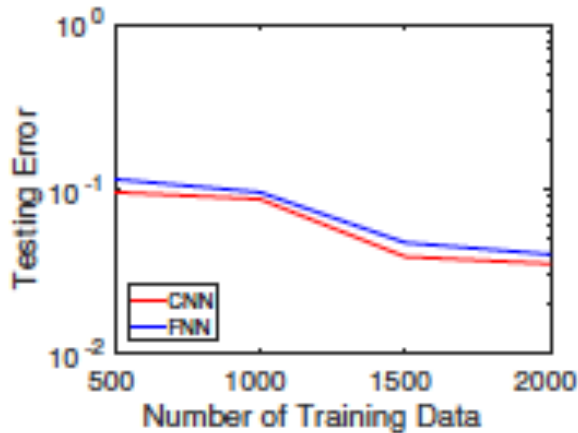
(c) Filter size $m = 16$.

$\tilde{\Theta} \binom{m}{n}$ Error decreases with n , increases with m and does not change with s

Experiments - CNN (weighted pooling) vs FNN

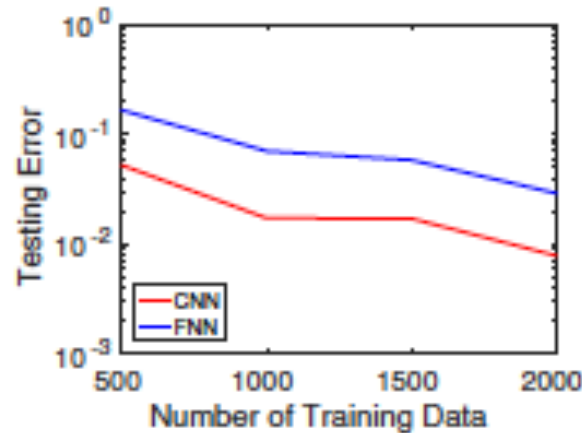
Filter size, $m = 8$

$d = 64$



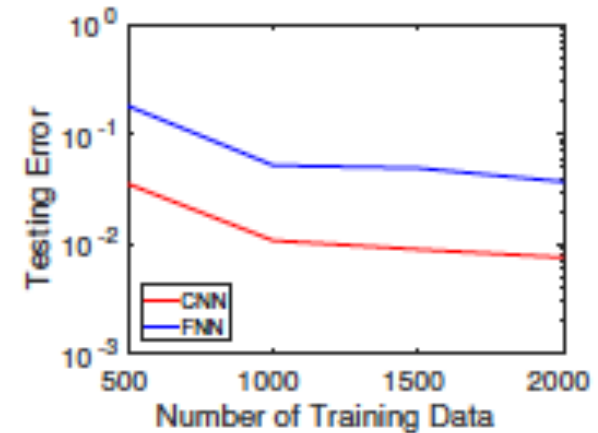
(a) Stride size $s = 1$.

$$m + J = 65$$



(b) Stride size $s = m/2$.

$$m + J = 23$$



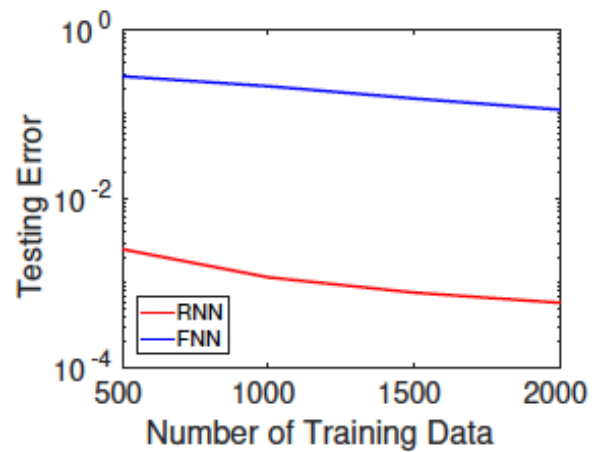
(c) Stride size $s = m$, i.e., non-overlapping.

$$m + J = 16$$

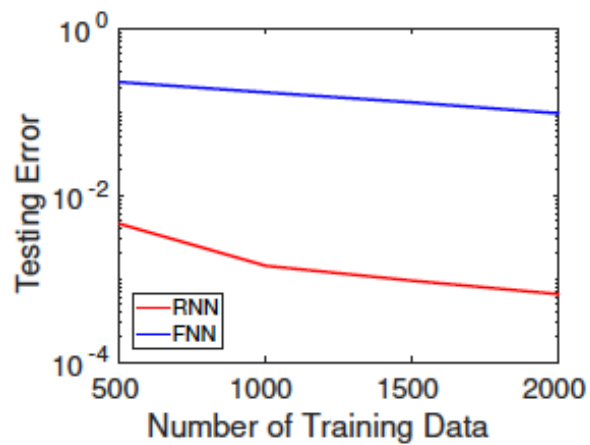
$\hat{\mathbb{I}} \left(\frac{m + J}{n} \right)$ Error decreases with n , increases with J (and m)
(larger stride s implies smaller output layer size $J \sim d/s$)

Experiments - RNN vs FNN

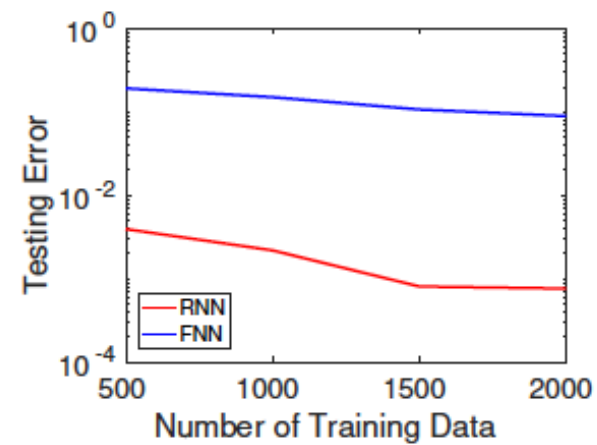
$d = 50, L = 50$



(a) Hidden units $r = 2$.



(b) Hidden units $r = 8$.



(c) Hidden units $r = 16$.

$$\tilde{\Theta} \left(\frac{rd}{n} \right)$$

Error decreases with n , increases with r (and d)

Open questions

- Is fast rate possible
 - without generative model assumption (i.e. non-realizable case)
 - without distributional assumptions in high dimensions ($d > n$)?
with computationally efficient estimators?

Nonlinear activations

Multiple filters

Deep models

- Role of Optimization

Similarities to sparse high-dimensional linear regression