

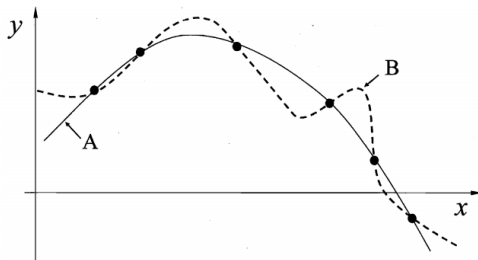
# Inductive Biases due to Dropout

Shihua Zhang

October 31, 2024

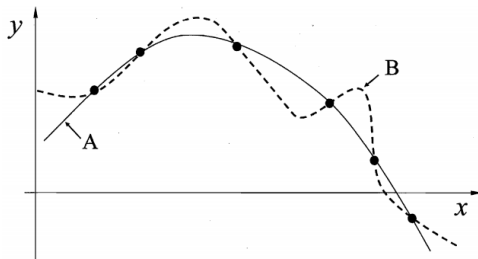
# Inductive Bias

- Given a finite dataset, there are **many possible solutions** to the learning problem.
- They exhibit **equally “good” performance** on the training points.



# Inductive Bias

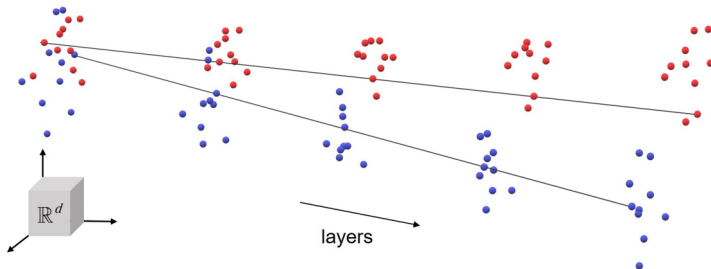
- Given a finite dataset, there are **many possible solutions** to the learning problem.
- They exhibit **equally “good” performance** on the training points.



- How to select the ones for **better generalization**?
- The **inductive bias** of a learning algorithm is **the set of assumptions** that the learner uses to predict unseen data.

# Inductive Biases in Deep Learning

- ResNet with weight decay learns the **geodesic curve in Wasserstein Space** [Gai and Zhang, 2021].
  - There are infinite curves connecting two distributions in  $\mathcal{P}(\mathbb{R}^d)$ .
  - This geodesic curve is induced by Optimal Transport map.
  - Data points are transported through **straight lines**.

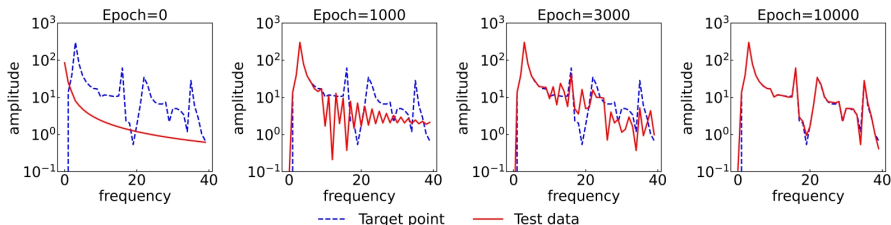


# Inductive Biases in Deep Learning

**Occam's Razor:** Entities should not be multiplied unnecessarily.

- **Frequency Principle** [Xu et al., 2024]

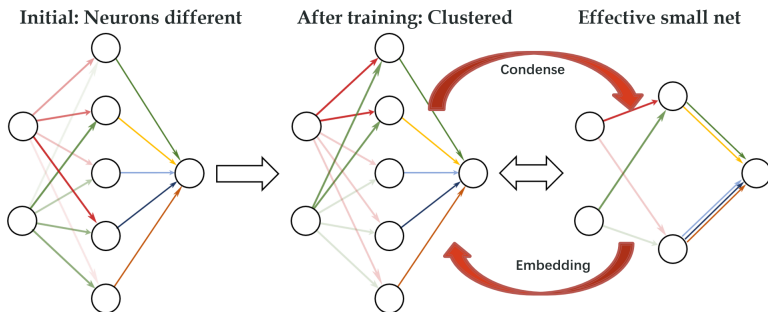
- DNNs often fit target functions **from low to high frequencies**.
- Red: FFT of the target function.
- Blue: FFT of DNN output.



# Inductive Biases in Deep Learning

**Occam's Razor:** Entities should not be multiplied unnecessarily.

- Frequency Principle [Xu et al., 2024]
- Condensation [Zhou et al., 2022]
  - Input weights of neurons in a group are same.
  - **Few effective** neurons.



# Inductive Biases due to Algorithmic Regularization

Several regularization strategies help to **generalize** in deep learning:

- Explicit regularization on objectives
  - $\ell_1$  regularization
  - $\ell_2$  regularization/weight decay

# Inductive Biases due to Algorithmic Regularization

Several regularization strategies help to **generalize** in deep learning:

- Explicit regularization on objectives
  - $\ell_1$  regularization
  - $\ell_2$  regularization/weight decay
- Heuristic techniques
  - Early stopping of back-propagation [Caruana et al., 2001]
  - Batch normalization [Ioffe and Szegedy, 2015]
  - Layer normalization [Ba et al., 2016]
  - Dropout [Srivastava et al., 2014]



# Inductive Biases due to Algorithmic Regularization

Several regularization strategies help to **generalize** in deep learning:

- Explicit regularization on objectives
  - $\ell_1$  regularization
  - $\ell_2$  regularization/weight decay
- Heuristic techniques
  - Early stopping of back-propagation [Caruana et al., 2001]
  - Batch normalization [Ioffe and Szegedy, 2015]
  - Layer normalization [Ba et al., 2016]
  - Dropout [Srivastava et al., 2014]

Today, we focus on the inductive biases due to **Dropout**.

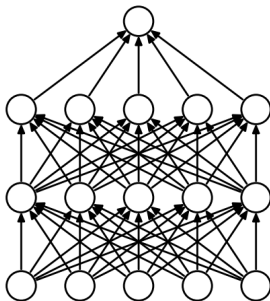
- Introduction to Dropout
- Matrix Sensing with Dropout
  - Gaussian sensing matrices
  - Matrix completion
- Dropout: Explicit Forms and Capacity Control
- Dropout Effects on Loss Landscape of the Optimization Problem
  - Implicit bias in local optima
  - Landscape properties

- 1 Introduction to Dropout
- 2 Matrix Sensing with Dropout
- 3 Dropout: Explicit Forms and Capacity Control
- 4 Dropout Effects on Loss Landscape of the Optimization Problem

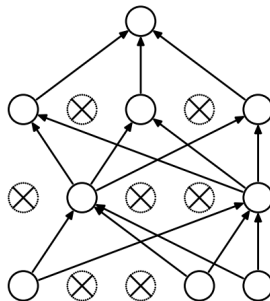
# Dropout

- A popular algorithmic heuristic with **limited formal understanding**.
- **Key idea**: **randomly drop units** of DNN during training.
- **Motivation**: as a way to **break “co-adaptation”** [Srivastava et al., 2014].

SRIVASTAVA, HINTON, KRIZHEVSKY, SUTSKEVER AND SALAKHUTDINOV



(a) Standard Neural Net



(b) After applying dropout.

# Training with Dropout

- With dropout, the feed-forward operation becomes

$$B_{ij} \sim \frac{1}{1-p} \text{Bernoulli}(1-p) \quad \text{i.i.d.}$$

$$z_i^{(l+1)} = W_i^{(l+1)} B y^{(l)} + b_i^{(l+1)}$$

$$y_i^{(l+1)} = \sigma(z_i^{(l+1)})$$

## Stochastic gradient descent

- For each training case in a mini-batch, **sample a thinned network**
- Do forward and back propagation on this thinned network
- The gradients for each parameter are **averaged** over the training cases in each mini-batch

# Experiments on Image Data Sets

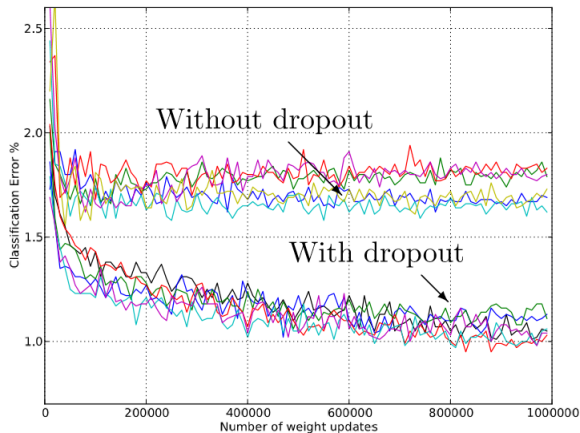


Figure: Test error for different architectures with dropout [Srivastava et al., 2014].

- 1 Introduction to Dropout
- 2 Matrix Sensing with Dropout**
  - Induced Regularizer for Matrix Sensing
  - Gaussian Matrix Sensing
  - Matrix Completion
- 3 Dropout: Explicit Forms and Capacity Control
- 4 Dropout Effects on Loss Landscape of the Optimization Problem

# Matrix Sensing

- Recover a matrix  $M_* \in \mathbb{R}^{d_2 \times d_0}$ , with rank  $r_* := \text{Rank}(M_*)$
- Given  $y_i = \langle M_*, A^{(i)} \rangle$ , for matrices  $A^{(1)}, \dots, A^{(n)}$ ,  $n \ll d_2 d_0$
- Represent  $M$  in **the factorized form** and solve:

$$\underset{U \in \mathbb{R}^{d_2 \times d_1}, V \in \mathbb{R}^{d_0 \times d_1}}{\text{minimize}} \quad \hat{L}(U, V) := \frac{1}{n} \sum_{i=1}^n \left( y_i - \langle UV^\top, A^{(i)} \rangle \right)^2 \quad (1)$$



# Matrix Sensing

- Recover a matrix  $M_* \in \mathbb{R}^{d_2 \times d_0}$ , with rank  $r_* := \text{Rank}(M_*)$
- Given  $y_i = \langle M_*, A^{(i)} \rangle$ , for matrices  $A^{(1)}, \dots, A^{(n)}$ ,  $n \ll d_2 d_0$
- Represent  $M$  in **the factorized form** and solve:

$$\underset{U \in \mathbb{R}^{d_2 \times d_1}, V \in \mathbb{R}^{d_0 \times d_1}}{\text{minimize}} \quad \hat{L}(U, V) := \frac{1}{n} \sum_{i=1}^n \left( y_i - \langle UV^\top, A^{(i)} \rangle \right)^2 \quad (1)$$

- Dropout as an instance of SGD:

$$\hat{L}_{\text{drop}}(U, V) = \frac{1}{n} \sum_{i=1}^n \mathbb{E}_B \left( y_i - \langle UB V^\top, A^{(i)} \rangle \right)^2 \quad (2)$$

where diagonal matrix  $B$  has  $B_{jj} \sim \frac{1}{1-p} \text{Ber}(1-p)$

# Explicit Regularizer

- **Key:** Dropout **explicitly regularizes** the empirical objective

$$\hat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \hat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1-p} \hat{R}(\mathbf{U}, \mathbf{V}) \quad (3)$$

where  $\hat{R}(\mathbf{U}, \mathbf{V}) = \sum_{j=1}^{d_1} \frac{1}{n} \sum_{i=1}^n \left( \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right)^2$ , data dependent

# Explicit Regularizer

- **Key:** Dropout **explicitly regularizes** the empirical objective

$$\hat{L}_{\text{drop}}(\mathbf{U}, \mathbf{V}) = \hat{L}(\mathbf{U}, \mathbf{V}) + \frac{\rho}{1-\rho} \hat{R}(\mathbf{U}, \mathbf{V}) \quad (3)$$

where  $\hat{R}(\mathbf{U}, \mathbf{V}) = \sum_{j=1}^{d_1} \frac{1}{n} \sum_{i=1}^n \left( \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right)^2$ , data dependent

- **Proof.** Consider one of the summands in the Dropout objective.

$$\begin{aligned} \mathbb{E}_{\mathbf{B}} \left[ \left( y_i - \langle \mathbf{U} \mathbf{B} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle \right)^2 \right] &= \left( \mathbb{E}_{\mathbf{B}} \left[ y_i - \langle \mathbf{U} \mathbf{B} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle \right] \right)^2 \\ &\quad + \text{Var} \left( y_i - \langle \mathbf{U} \mathbf{B} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle \right) \end{aligned}$$

- Note that  $\mathbb{E} [\mathbf{B}_{jj}] = 1$  and  $\text{Var} (\mathbf{B}_{jj}) = \frac{\rho}{1-\rho}$ , the first term on the right side is equal to  $\left( y_i - \langle \mathbf{U} \mathbf{V}^\top, \mathbf{A}^{(i)} \rangle \right)^2$ .

# Explicit Regularizer

- For the second term we have

$$\begin{aligned}\text{Var} \left( y_i - \left\langle \text{UBV}^\top, \mathbf{A}^{(i)} \right\rangle \right) &= \text{Var} \left( \left\langle \text{UBV}^\top, \mathbf{A}^{(i)} \right\rangle \right) \\ &= \text{Var} \left( \left\langle \mathbf{B}, \mathbf{U}^\top \mathbf{A}^{(i)} \mathbf{V} \right\rangle \right) \\ &= \text{Var} \left( \sum_{j=1}^{d_1} \mathbf{B}_{jj} \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right) \\ &= \sum_{j=1}^{d_1} \left( \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right)^2 \text{Var} (\mathbf{B}_{jj}) \\ &= \frac{p}{1-p} \sum_{j=1}^{d_1} \left( \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right)^2\end{aligned}$$

# Induced Regularizer

- Thus,

$$\begin{aligned}\hat{L}_{\text{drop}} &= \frac{1}{n} \sum_{i=1}^n \left( y_i - \langle \mathbf{U}\mathbf{V}^\top, \mathbf{A}^{(i)} \rangle \right)^2 + \frac{1}{n} \sum_{i=1}^n \frac{p}{1-p} \sum_{j=1}^{d_1} \left( \mathbf{u}_j^\top \mathbf{A}^{(i)} \mathbf{v}_j \right)^2 \\ &= \hat{L}(\mathbf{U}, \mathbf{V}) + \frac{p}{1-p} \hat{R}(\mathbf{U}, \mathbf{V})\end{aligned}$$

- **Expected regularizer:**  $R(\mathbf{U}, \mathbf{V}) := \mathbb{E}_{\mathbf{A}}[\hat{R}(\mathbf{U}, \mathbf{V})]$
- **Induced regularizer:** Consider the factors with the minimal value of  $R(\mathbf{U}, \mathbf{V})$  among all that yield the same empirical loss

$$\Theta(\mathbf{M}) := \min_{\mathbf{U}\mathbf{V}^\top = \mathbf{M}} R(\mathbf{U}, \mathbf{V})$$

# Gaussian Matrix Sensing

- Assume that the entries of the sensing matrices are iid as standard Gaussian, i.e.,  $A_{k\ell}^{(i)} \sim \mathcal{N}(0, 1)$ .
- **Hint:** The induced regularizer due to Dropout provides **the nuclear-norm** regularization:

$$\Theta(M) := \min_{UV^T=M} R(U, V) = \frac{1}{d_1} \|M\|_*^2$$

# Gaussian Matrix Sensing

- Assume that the entries of the sensing matrices are iid as standard Gaussian, i.e.,  $A_{k\ell}^{(i)} \sim \mathcal{N}(0, 1)$ .
- Hint:** The induced regularizer due to Dropout provides **the nuclear-norm** regularization:

$$\Theta(M) := \min_{UV^T=M} R(U, V) = \frac{1}{d_1} \|M\|_*^2$$

- For any pair of factors  $(U, V)$ , the **expected regularizer** is

$$R(U, V) = \sum_{i=1}^{d_1} \mathbb{E}_A \left[ \left( \mathbf{u}_i^T A \mathbf{v}_i \right)^2 \right] = \sum_{i=1}^{d_1} \|\mathbf{u}_i\|^2 \|\mathbf{v}_i\|^2$$

# Gaussian Matrix Sensing

- By Cauchy-Schwartz inequality

$$\begin{aligned} R(U, V) &= \sum_{i=1}^{d_1} \|u_i\|^2 \|v_i\|^2 \geq \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|u_i\| \|v_i\| \right)^2 \\ &= \frac{1}{d_1} \left( \sum_{i=1}^{d_1} \|u_i v_i^\top\|_* \right)^2 \\ &\geq \frac{1}{d_1} \left( \left\| \sum_{i=1}^{d_1} u_i v_i^\top \right\|_* \right)^2 = \frac{1}{d_1} \|UV^\top\|_*^2 \end{aligned}$$

- Here the equality follows because for any pair of vectors  $a, b$ , it holds that  $\|ab^\top\|_* = \|ab^\top\|_F = \|a\| \|b\|$
- Lower bound can be achieved for all  $(U, V)$  s.t.

$$\|u_i\| \|v_i\| = \frac{1}{d_1} \|UV^\top\|_*, \forall i$$



# Gaussian Matrix Sensing

Based on the following result on  $(U, V)$  [Mianjy et al., 2018]:

## Theorem 1

*For any pair of matrices  $U \in \mathbb{R}^{d_2 \times d_1}$ ,  $V \in \mathbb{R}^{d_0 \times d_1}$ , there exists a rotation matrix  $Q$  such that matrices  $\tilde{U} := UQ$ ,  $\tilde{V} := VQ$  satisfy  $\|\tilde{u}_i\| \|\tilde{v}_i\| = \frac{1}{d_1} \|UV^\top\|_*$ , for all  $i \in [d_1]$ .*

# Gaussian Matrix Sensing

Based on the following result on  $(U, V)$  [Mianjy et al., 2018]:

## Theorem 1

*For any pair of matrices  $U \in \mathbb{R}^{d_2 \times d_1}$ ,  $V \in \mathbb{R}^{d_0 \times d_1}$ , there exists a rotation matrix  $Q$  such that matrices  $\tilde{U} := UQ$ ,  $\tilde{V} := VQ$  satisfy  $\|\tilde{u}_i\| \|\tilde{v}_i\| = \frac{1}{d_1} \|UV^\top\|_*$ , for all  $i \in [d_1]$ .*

$$\begin{aligned} R(UQ, VQ) &= \sum_{i=1}^{d_1} \|U_{Q_i}\|^2 \|V_{Q_i}\|^2 \\ &= \sum_{i=1}^{d_1} \frac{1}{d_1^2} \|UV^\top\|_*^2 \\ &= \frac{1}{d_1} \|UV^\top\|_*^2 \end{aligned}$$

# Matrix Completion

- Matrix completion (MC) can be formulated as a special case of matrix sensing with sensing matrices being **random indicator matrices**.
- Let  $A^{(j)}$  be an indicator matrix whose  $(i, k)$ -th element is selected randomly with probability  $p(i), q(k)$ , then

$$\Theta(M) = \frac{1}{d_1} \left\| \sqrt{\text{diag}(p)} UV^T \sqrt{\text{diag}(q)} \right\|_*^2 \quad (\text{weighted trace-norm})$$

- The **weighted trace-norm or nuclear norm** has been studied by [Salakhutdinov and Srebro, 2010][Foygel et al., 2011]

# Matrix Completion

- Matrix completion (MC) can be formulated as a special case of matrix sensing with sensing matrices being **random indicator matrices**.
- Let  $A^{(j)}$  be an indicator matrix whose  $(i, k)$ -th element is selected randomly with probability  $p(i), q(k)$ , then

$$\Theta(M) = \frac{1}{d_1} \left\| \sqrt{\text{diag}(p)} UV^T \sqrt{\text{diag}(q)} \right\|_*^2 \quad (\text{weighted trace-norm})$$

- The **weighted trace-norm or nuclear norm** has been studied by [Salakhutdinov and Srebro, 2010][Foygel et al., 2011]
- Key:** A **generalization bound** for MC with dropout in terms of the value of the explicit regularizer at the minimum of the empirical problem [Arora et al., 2021].

# A Generalization Bound for Matrix Completion

## Theorem 2 ([Arora et al., 2021])

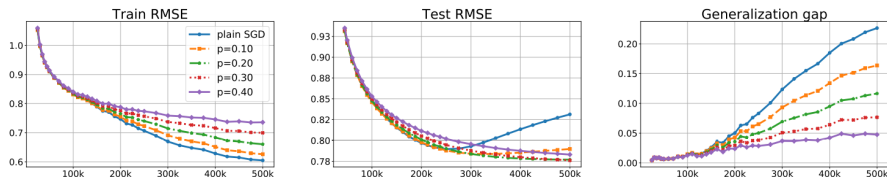
Assume that  $d_2 \geq d_0$  and  $\|M_*\| \leq 1$ . Furthermore, assume that  $\min_{i,k} p(i)q(k) \geq \frac{\log(d_2)}{n\sqrt{d_2d_0}}$ . Let  $(U, V)$  be a minimizer of the dropout objective in equation (3). Let  $\alpha$  be such that  $R(U, V) \leq \alpha/d_1$ . Then, for any  $\delta \in (0, 1)$ , the following generalization bounds holds with probability at least  $1 - \delta$  over a sample of size  $n$ :

$$L\left(g\left(UV^\top\right)\right) \leq \widehat{L}(U, V) + 8\sqrt{\frac{2\alpha d_2 \log(d_2) + \frac{1}{4} \log(2/\delta)}{n}}$$

where  $g(M)$  thresholds  $M$  at  $\pm 1$ , i.e.,  $g(M)(i, j) = \max\{-1, \min\{1, M(i, j)\}\}$ , and  $L(g(UV^\top)) := \mathbb{E}(y - \langle g(UV^\top), A \rangle)^2$  is the true risk of  $g(UV^\top)$

# Empirical Results on Matrix Completion

MovieLens dataset: 10M ratings for 11K movies by 72K users.



- The training error, test error, and generalization gap for plain SGD and dropout with different  $p$  as a function of the number of iterations.
- Intuitively, a larger dropout rate  $p$  results in a smaller  $\alpha$ .

# Empirical Results on Matrix Completion

MovieLens dataset: 10M ratings for 11K movies by 72K users.

width	plain SGD		dropout			
	last iterate	best iterate	$p = 0.1$	$p = 0.2$	$p = 0.3$	$p = 0.4$
$d_1 = 30$	0.8041	0.7938	0.7805	0.785	0.7991	0.8186
$d_1 = 70$	0.8315	0.7897	0.7899	0.7771	0.7763	0.7833
$d_1 = 110$	0.8431	0.7873	0.7988	0.7813	0.7742	0.7743
$d_1 = 150$	0.8472	0.7858	0.8042	0.7852	0.7756	0.7722
$d_1 = 190$	0.8473	0.7844	0.8069	0.7879	0.7772	0.772

**Figure:** Test RMSE of plain SGD and the dropout algorithm with various dropout rates for various factorization sizes [Arora et al., 2021].

- Dropout performance improves with the size of the parametrization.
- SGD has worse generalization even for best iterate on test data.

- 1 Introduction to Dropout
- 2 Matrix Sensing with Dropout
- 3 Dropout: Explicit Forms and Capacity Control**
- 4 Dropout Effects on Loss Landscape of the Optimization Problem



# Regression with Deep Neural Networks

- $\mathcal{X} \subseteq \mathbb{R}^{d_0}$ ,  $\mathcal{Y} \subseteq [-1, 1]^{d_2}$ ,  $\mathcal{D}$  is an (unknown) distribution on  $\mathcal{X} \times \mathcal{Y}$
- 2-layers neural networks parameterized by  $\mathbf{w}$

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{U} \sigma \left( \mathbf{V}^\top \mathbf{x} \right)$$

where  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{d_2}] \in \mathbb{R}^{d_2 \times d_1}$ ,  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_{d_1}] \in \mathbb{R}^{d_0 \times d_1}$ .

- Squared  $\ell_2$  loss,  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , with  $\ell(\mathbf{y}, \mathbf{y}') = \|\mathbf{y} - \mathbf{y}'\|^2$

# Regression with Deep Neural Networks

- $\mathcal{X} \subseteq \mathbb{R}^{d_0}$ ,  $\mathcal{Y} \subseteq [-1, 1]^{d_2}$ ,  $\mathcal{D}$  is an (unknown) distribution on  $\mathcal{X} \times \mathcal{Y}$
- 2-layers neural networks parameterized by  $\mathbf{w}$

$$f_{\mathbf{w}}(\mathbf{x}) = \mathbf{U} \sigma \left( \mathbf{V}^\top \mathbf{x} \right)$$

where  $\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_{d_2}] \in \mathbb{R}^{d_2 \times d_1}$ ,  $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_{d_1}] \in \mathbb{R}^{d_0 \times d_1}$ .

- Squared  $\ell_2$  loss,  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , with  $\ell(y, y') = \|y - y'\|^2$
- **Goal:** find a hypothesis  $f_{\mathbf{w}} : \mathcal{X} \rightarrow \mathcal{Y}$ , with a small

$$L(\mathbf{w}) := \mathbb{E}_{\mathcal{D}} [\ell(f_{\mathbf{w}}(\mathbf{x}), y)] \quad (\text{population risk})$$

- Given  $n$  samples  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$  drawn i.i.d. from  $\mathcal{D}$

# Regression with Deep Neural Networks

- $\mathcal{X} \subseteq \mathbb{R}^{d_0}$ ,  $\mathcal{Y} \subseteq [-1, 1]^{d_2}$ ,  $\mathcal{D}$  is an (unknown) distribution on  $\mathcal{X} \times \mathcal{Y}$
- 2-layers neural networks parameterized by  $w$

$$f_w(x) = U\sigma(V^\top x)$$

where  $U = [u_1, \dots, u_{d_2}] \in \mathbb{R}^{d_2 \times d_1}$ ,  $V = [v_1, \dots, v_{d_1}] \in \mathbb{R}^{d_0 \times d_1}$ .

- Squared  $\ell_2$  loss,  $\ell : \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ , with  $\ell(y, y') = \|y - y'\|^2$
- **Goal**: find a hypothesis  $f_w : \mathcal{X} \rightarrow \mathcal{Y}$ , with a small

$$L(w) := \mathbb{E}_{\mathcal{D}} [\ell(f_w(x), y)] \quad (\text{population risk})$$

- Given  $n$  samples  $\{(x_i, y_i)\}_{i=1}^n \sim \mathcal{D}^n$  drawn i.i.d. from  $\mathcal{D}$
- **ERM**: minimize

$$\hat{L}(w) := \frac{1}{n} \sum_{i=1}^n [\|y_i - f_w(x_i)\|^2] \quad (\text{empirical risk})$$

# Dropout in Deep Neural Networks

- Dropout as SGD iterates – the dropout objective:

$$\hat{L}_{\text{drop}}(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{B}} \left\| y_i - \text{UB} \sigma \left( \mathbf{V}^\top \mathbf{x}_i \right) \right\|^2$$

where  $B_{ij} \sim \frac{1}{1-p} \text{Bern}(1-p), i \in [d_1]$ .

- We seek to understand the **explicit regularizer** due to dropout:

$$\hat{R}(\mathbf{w}) := \hat{L}_{\text{drop}}(\mathbf{w}) - \hat{L}(\mathbf{w}) \quad (\text{explicit regularizer})$$

- Denote the output of the  $i$ -th hidden node on input  $\mathbf{x}$  by  $a_i(\mathbf{x})$ ;  
 $a(\mathbf{x}) \in \mathbb{R}^{d_1}$  denotes the activation of the hidden layer on input  $\mathbf{x}$ .
- Rewrite the Dropout objective as

$$\hat{L}_{\text{drop}}(\mathbf{w}) := \frac{1}{n} \sum_{i=1}^n \mathbb{E}_{\mathbf{B}} \left\| y_i - \text{UB} a(\mathbf{x}_i) \right\|^2.$$

# Dropout Regularizer in Deep Regression

- The explicit regularizer due to dropout is

$$\widehat{R}(\mathbf{w}) = \lambda \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \widehat{a}_j^2, \quad \widehat{a}_j = \sqrt{\frac{1}{n} \sum_{i=1}^n a_j(x_i)^2}$$

where  $\lambda = \frac{\rho}{1-\rho}$  is the regularization parameter.

# Dropout Regularizer in Deep Regression

- The explicit regularizer due to dropout is

$$\widehat{R}(\mathbf{w}) = \lambda \sum_{j=1}^{d_1} \|\mathbf{u}_j\|^2 \widehat{a}_j^2, \quad \widehat{a}_j = \sqrt{\frac{1}{n} \sum_{i=1}^n a_j(x_i)^2}$$

where  $\lambda = \frac{\rho}{1-\rho}$  is the regularization parameter.

- Consider ReLU activations and input distributions that are symmetric and isotropic, i.e.,  $\mathbb{P}_{\mathcal{X}}(\mathbf{x}) = \mathbb{P}_{\mathcal{X}}(-\mathbf{x})$  and  $\mathbf{C} = \mathbb{E}[\mathbf{x}\mathbf{x}^\top] = \mathbf{I}$ . Then the **expected regularizer** due to dropout is given as

$$R(\mathbf{w}) := \mathbb{E}[\widehat{R}(\mathbf{w})] = \frac{\lambda}{2} \sum_{i_0, i_1, i_2=1}^{d_0, d_1, d_2} U(i_2, i_1)^2 V(i_1, i_0)^2$$

- It is a data-dependent variant of **the  $\ell_2$  path-norm** of the network [Neyshabur et al., 2015].
- It can yield **capacity control** in deep learning.

# Function Class Learned by Dropout

- Let  $d_2 = 1$ , we focus on the following distribution-dependent class

$$\mathcal{F}_\alpha := \left\{ f_w : x \mapsto u^\top \sigma(V^\top x), \sum_{i=1}^{d_1} |u_i| a_i \leq \alpha \right\}$$

where  $a_i^2 := \mathbb{E}_x [\hat{a}_i^2] = \mathbb{E}_x [a_i(x)^2]$

- We argue that networks trained with dropout belong to the class  $\mathcal{F}_\alpha$  (for a small value of  $\alpha$ ).

# Function Class Learned by Dropout

- Let  $d_2 = 1$ , we focus on the following distribution-dependent class

$$\mathcal{F}_\alpha := \left\{ f_w : x \mapsto u^\top \sigma(V^\top x), \sum_{i=1}^{d_1} |u_i| a_i \leq \alpha \right\}$$

where  $a_i^2 := \mathbb{E}_x [\hat{a}_i^2] = \mathbb{E}_x [a_i(x)^2]$

- We argue that networks trained with dropout belong to the class  $\mathcal{F}_\alpha$  (for a small value of  $\alpha$ ).
- By Cauchy-Schwartz inequality,

$$\sum_{i=1}^{d_1} |u_i| a_i \leq \sqrt{d_1 \sum_{i=1}^{d_1} |u_i|^2 a_i^2} = \sqrt{d_1 \frac{1}{\lambda} R(w)}$$

Thus, for a fixed width, dropout controls the function class  $\mathcal{F}_\alpha$ .



# Function Class Learned by Dropout

- This inequality is loose if a small subset of hidden nodes  $\mathcal{J} \subset [d_1]$  “co-adapt” in a way that the other hidden nodes (i.e., all  $j \in [d_1] \setminus \mathcal{J}$ ) are almost inactive, i.e.  $u_j a_j \approx 0$ .
- By minimizing the expected regularizer, dropout is biased towards networks where the gap between  $\frac{1}{d_1} \left( \sum_{i=1}^{d_1} |u_i| a_i \right)^2$  and  $R(\mathbf{w})$  is small, which in turn happens if

$$|u_i| a_i \approx |u_j| a_j, \forall i, j \in [d_1].$$

- Dropout breaks “co-adaptation” by promoting solutions with **nearly equal contribution** from hidden neurons.

# Function Class Learned by Dropout

- This inequality is loose if a small subset of hidden nodes  $\mathcal{J} \subset [d_1]$  “co-adapt” in a way that the other hidden nodes (i.e., all  $j \in [d_1] \setminus \mathcal{J}$ ) are almost inactive, i.e.  $u_j a_j \approx 0$ .
- By minimizing the expected regularizer, dropout is biased towards networks where the gap between  $\frac{1}{d_1} \left( \sum_{i=1}^{d_1} |u_i| a_i \right)^2$  and  $R(\mathbf{w})$  is small, which in turn happens if

$$|u_i| a_i \approx |u_j| a_j, \forall i, j \in [d_1].$$

- Dropout breaks “co-adaptation” by promoting solutions with **nearly equal contribution** from hidden neurons.
- Next, under mild condition on the input distribution, **a generalization bound** can be derived.

# Bound on the Rademacher Complexity

## Assumption 1 ( $\beta$ -retentive)

The marginal input distribution is  $\beta$ -retentive for some  $\beta \in (0, 1/2]$ , if for any non-zero vector  $\mathbf{v} \in \mathbb{R}^d$ , it holds that  $\mathbb{E} \sigma(\mathbf{v}^\top \mathbf{x})^2 \geq \beta \mathbb{E} (\mathbf{v}^\top \mathbf{x})^2$ .

- Mahalanobis norm:  $\|\mathbf{X}\|_{\mathbf{C}^\dagger}^2 = \sum_{i=1}^n \mathbf{x}_i^\top \mathbf{C}^\dagger \mathbf{x}_i$ .

## Theorem 3

For any sample  $S = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$  of size  $n$ ,

$$\mathfrak{R}_S(\mathcal{F}_\alpha) \leq \frac{2\alpha \|\mathbf{X}\|_{\mathbf{C}^\dagger}}{n\sqrt{\beta}}$$

Furthermore, it holds for the **expected Rademacher complexity** that

$$\mathfrak{R}_n(\mathcal{F}_\alpha) \leq 2\alpha \sqrt{\frac{\text{Rank}(\mathbf{C})}{\beta n}}.$$

# Generalization Bounds

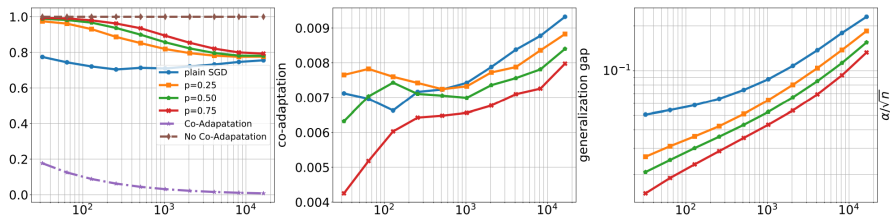
- Dropout regularizer **directly controls** the value of  $\alpha$ , thereby **controlling** the Rademacher complexity in Theorem 3.
- Let  $g_w(\cdot) := \max\{-1, \min\{1, f_w(\cdot)\}\}$  project the network output  $f_w$  onto the range  $[-1, 1]$ . We have the following generalization guarantees for  $g_w$  based on Theorem 3.

## Theorem 4

*For any  $f_w \in \mathcal{F}_\alpha$ , for any  $\delta \in (0, 1)$ , the following generalization bound holds with probability at least  $1 - \delta$  over a sample  $S$  of size  $n$*

$$L(g_w) \leq \hat{L}(g_w) + \frac{16\alpha\|X\|_{C^\dagger}}{\sqrt{\beta}n} + 12\sqrt{\frac{\log(2/\delta)}{2n}}$$

# Experimental Results



**Figure:** “co-adaptation”, generalization gap and  $\alpha/\sqrt{n}$  as a function of the width of networks trained with dropout on MNIST. The trained 2-layer networks achieve 100% training accuracy [Arora et al., 2021]

- Increasing the dropout rate results in **less co-adaptation** empirically.
- Increasing dropout rate **decreases the generalization gap**.
- The bound of the Rademacher complexity is predictive on the generalization gap.

- 1 Introduction to Dropout
- 2 Matrix Sensing with Dropout
- 3 Dropout: Explicit Forms and Capacity Control
- 4 Dropout Effects on Loss Landscape of the Optimization Problem**
  - Implicit bias in local optima
  - Landscape properties

# Goal

- Dropout is a **first-order method** and the landscape of the Dropout objective (e.g., Problem (4)) is highly **non-convex**.
- Can perhaps only hope to find a local minimum, that too provided if the problem has no degenerate saddle points [Ge et al., 2015].
- Therefore, the following questions are expected:
  - What is the **implicit bias of dropout** in terms of local minima?
  - Do local minima share anything with global minima structurally?
  - Can dropout **find a local optimum**?

# Problem Setup

- We focus on the case of single hidden layer linear autoencoders with **tied weights**, i.e.  $U = V$ .

$$\mathcal{H}_r := \left\{ h_U : x \mapsto UU^T x, U \in \mathbb{R}^{d_0 \times d_1} \right\}$$

- Assume that the input distribution is **isotropic**, i.e.  $C_x = \mathbb{E} [xx^T] = I$
- The population risk reduces to

$$\begin{aligned} \mathbb{E} \left[ \left\| y - UU^T x \right\|^2 \right] &= \text{Tr} (C_y) - 2 \left\langle C_{yx}, UU^T \right\rangle + \left\| UU^T \right\|_F^2 \\ &= \left\| M - UU^T \right\|_F^2 + \text{Tr} (C_y) - \|M\|_F^2 \end{aligned}$$

where  $M = \frac{C_{yx} + C_{xy}}{2}$ .



# Problem Setup

- Ignoring the terms that are independent of the weight matrix  $U$ , the goal is to minimize  $L(U) = \|M - UU^T\|_F^2$ .
- Solving the following problem with Dropout:

$$\min_{U \in \mathbb{R}^{d_0 \times d_1}} L_\theta(U) := \|M - UU^T\|_F^2 + \lambda \underbrace{\sum_{i=1}^{d_1} \|u_i\|^4}_{R(U)} \quad (4)$$

# Implicit Bias in Local Optima

- $L(U)$  is **rotation invariant**, i.e. for any rotation matrix  $Q$

$$L(UQ) = \left\| M - UQQ^\top U^\top \right\|_F^2 = L(U), \quad Q^\top Q = QQ^\top = I$$

But the regularizer is **not** rotation invariant.

# Implicit Bias in Local Optima

- $L(U)$  is **rotation invariant**, i.e. for any rotation matrix  $Q$

$$L(UQ) = \left\| M - UQQ^\top U^\top \right\|_F^2 = L(U), \quad Q^\top Q = QQ^\top = I$$

But the regularizer is **not** rotation invariant.

- By Cauchy-Schwartz inequality, we have

$$R(U) = \lambda \sum_{i=1}^{d_1} \|u_i\|^4 \geq \frac{\lambda}{d_1} \|U\|_F^4$$

with equality iff all the columns of  $U$  have **equal norms** (equalized).

- If the weight matrix  $U$  were not equalized, one can design a rotation matrix  $Q$  that  $UQ$  has a smaller regularizer, hence the objective.

# Implicit Bias in Local Optima – Theorem

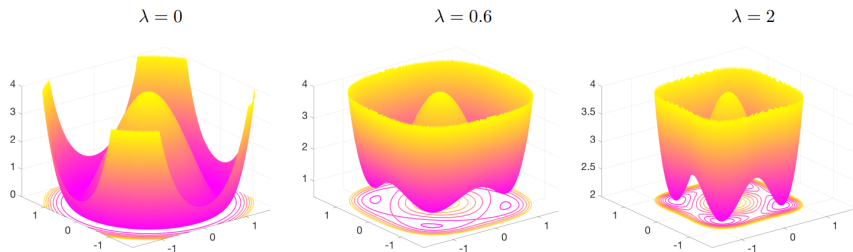
If  $\mathbf{U}$  is not equalized, then any  $\epsilon$ -neighborhood of  $\mathbf{U}$  contains a point with dropout objective strictly smaller than  $L_{\theta}(\mathbf{U})$ .

## Theorem 5 ([Mianjy et al., 2018])

*All local minima of Problem (4) are equalized, i.e. if  $\mathbf{U}$  is a local optimum, then  $\|u_i\| = \|u_j\| \forall i, j \in [r]$ .*

- Dropout tends to give **equal weights to all hidden nodes**.
- No matter how small the dropout rate – all local minima become equalized.

# Implicit Bias in Local Optima – Illustration



**Figure:** Optimization landscape for a single hidden-layer linear autoencoder network with dropout, for different regularization parameter  $\lambda$ .

- (Middle) All local minima are global, and are equalized, i.e. the weights are parallel to  $(\pm 1, \pm 1)$ .
- (Right) As  $\lambda$  increases, global optima shrink further.

# Strict Saddle Point/Property

## Definition 6 (Strict saddle point/property)

Let  $f : \mathcal{U} \rightarrow \mathbb{R}$  be a twice differentiable function and let  $U \in \mathcal{U}$  be a critical point of  $f$ .

Then,  $U$  is a **strict saddle point** of  $f$  if the Hessian of  $f$  at  $U$  has at least one negative eigenvalue, i.e.  $\lambda_{\min}(\nabla^2 f(U)) < 0$ .

Furthermore,  $f$  satisfies **strict saddle property** if all saddle points of  $f$  are strict saddle.

- Strict saddle property ensures that for any critical point  $U$  that is not a local optimum, **the Hessian has a significant negative eigenvalue**.
- SGD can escape saddle points and converge to a local minimum [Ge et al., 2015].

# Landscape Properties

- For the special case of no dropout (i.e.  $\lambda = 0$ ), Problem (4) has been shown to have **no spurious local minima** and **satisfy strict saddle property** ([Baldi and Hornik, 1989, Jin et al., 2017]).
- **Question:** Can the regularizer induced by dropout potentially introduce new spurious local minima and/or degenerate saddle points?
- The answer is **no**, at least when the dropout rate is sufficiently small.

# Landscape properties

## Theorem 7 ([Mianjy et al., 2018])

Let  $r := \text{Rank}(M)$ . Assume that  $d_1 \leq d_0$  and that the regularization parameter satisfies  $\lambda < \frac{r\lambda_r(M)}{(\sum_{i=1}^r \lambda_i(M)) - r\lambda_r(M)}$ . Then it holds for Problem (4) that

1. all local minima are global,
2. all saddle points are strict saddle points.

- The theorem guarantees that any critical point  $U$  that is not a global optimum is a strict saddle point.
- This property allows SGD to escape such saddle points.



# Proof sketch

## Lemma 8

*All critical points of Problem (4) that are not equalized, are strict saddle points.*

## Lemma 9

*Let  $r := \text{Rank}(M)$ . Assume that  $d_1 \leq d_0$  and  $\lambda < \frac{r\lambda_r}{\sum_{i=1}^p (\lambda_i - \lambda_r)}$ . Then all equalized local minima are global. All other equalized critical points are strict saddle points.*

- Theorem 5 and lemma 8 show that non-equalized critical points are not local optima, they are strict saddle points.
- If  $\lambda$  is chosen appropriately, then all critical points that are not global optimum, are strict saddle points.

# Summary

- **Dropout** is a popular regularization with limited understanding.
- **Instantiate explicit forms of regularizers** due to Dropout and how they provide **capacity control** in various machine learning problems:
  - Gaussian matrix sensing
  - Matrix completion
  - Deep learning
- Dropout effects on condensation and flatness, resulting good generalization.
  - All local minima are equalized
  - All local minima are global
  - All saddle points are non-degenerate

# References I



Arora, R., Bartlett, P., Mianjy, P., and Srebro, N. (2021).

Dropout: Explicit forms and capacity control.

In *International Conference on Machine Learning*, pages 351–361. PMLR.



Ba, J. L., Kiros, J. R., and Hinton, G. E. (2016).

Layer normalization.



Baldi, P. and Hornik, K. (1989).

Neural networks and principal component analysis: Learning from examples without local minima.

*Neural networks*, 2(1):53–58.



Caruana, R., Lawrence, S., and Giles, L. (2001).

Overfitting in neural nets: Backpropagation, conjugate gradient, and early stopping.

*Advances in neural information processing systems*, pages 402–408.



Foygel, R., Salakhutdinov, R., Shamir, O., and Srebro, N. (2011).

Learning with the weighted trace-norm under arbitrary sampling distributions.

*arXiv preprint arXiv:1106.4251*.



Gai, K. and Zhang, S. (2021).

A mathematical principle of deep learning: Learn the geodesic curve in the wasserstein space.

*arXiv preprint arXiv:2102.09235*.



Ge, R., Huang, F., Jin, C., and Yuan, Y. (2015).

Escaping from saddle points—online stochastic gradient for tensor decomposition.

In *Conference on learning theory*, pages 797–842. PMLR.

# References II



Ioffe, S. and Szegedy, C. (2015).

Batch normalization: Accelerating deep network training by reducing internal covariate shift.  
*In International conference on machine learning*, pages 448–456. PMLR.



Jin, C., Ge, R., Netrapalli, P., Kakade, S. M., and Jordan, M. I. (2017).

How to escape saddle points efficiently.  
*In International Conference on Machine Learning*, pages 1724–1732. PMLR.



Mianjy, P., Arora, R., and Vidal, R. (2018).

On the implicit bias of dropout.  
*In International Conference on Machine Learning*, pages 3540–3548. PMLR.



Neyshabur, B., Tomioka, R., and Srebro, N. (2015).

Norm-based capacity control in neural networks.  
*In Conference on Learning Theory*, pages 1376–1401. PMLR.



Salakhutdinov, R. and Srebro, N. (2010).

Collaborative filtering in a non-uniform world: Learning with the weighted trace norm.  
*arXiv preprint arXiv:1002.2780*.



Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014).

Dropout: A simple way to prevent neural networks from overfitting.  
*Journal of Machine Learning Research*, 15(56):1929–1958.



Xu, Z.-Q. J., Zhang, Y., and Luo, T. (2024).

Overview frequency principle/spectral bias in deep learning.  
*Communications on Applied Mathematics and Computation*, pages 1–38.

# References III



Zhou, H., Qixuan, Z., Luo, T., Zhang, Y., and Xu, Z.-Q. (2022).

Towards understanding the condensation of neural networks at initial training.

*Advances in Neural Information Processing Systems*, 35:2184–2196.