

On Dropout and Nuclear Norm Regularization

Poorya Mianjy and Raman Arora

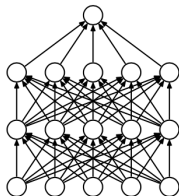
Johns Hopkins University

June 10, 2019

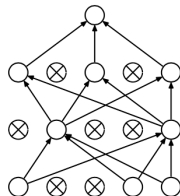
Motivation

- ▶ Algorithmic approaches endow deep learning systems with certain inductive biases that help generalization.
- ▶ In this paper we study dropout, one of the most popular algorithmic heuristics for training deep neural nets.

SRIVASTAVA, HINTON, KRIZHEVSKY, SUTSKEVER AND SALAKHUTDINOV



(a) Standard Neural Net



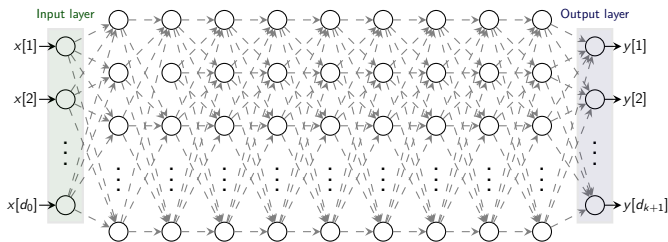
(b) After applying dropout.

Problem Setup

- ▶ Deep linear networks with k hidden layers

$$f_w : x \mapsto W_{k+1} \cdots W_1 x, \quad W_i \in \mathbb{R}^{d_i \times d_{i-1}}$$

where $w = \{W_i\}_{i=1}^{k+1}$ is the set of weight matrices.



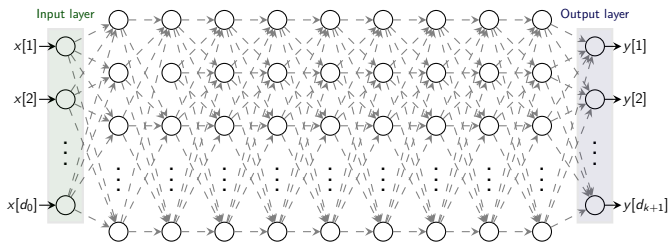
Problem Setup

- ▶ Deep linear networks with k hidden layers

$$f_w : x \mapsto W_{k+1} \cdots W_1 x, \quad W_i \in \mathbb{R}^{d_i \times d_{i-1}}$$

where $w = \{W_i\}_{i=1}^{k+1}$ is the set of weight matrices.

- ▶ $x \in \mathbb{R}^{d_0}$, $y \in \mathbb{R}^{d_{k+1}}$, $(x, y) \sim \mathcal{D}$. Assume $\mathbb{E}[xx^\top] = I$.



Problem Setup

- ▶ Deep linear networks with k hidden layers

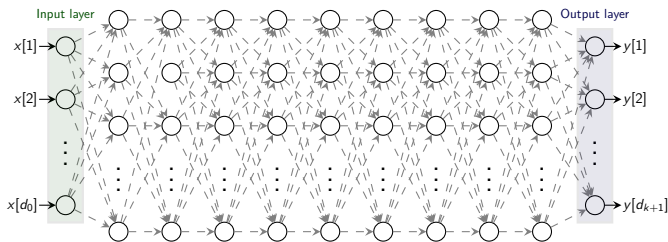
$$f_w : x \mapsto W_{k+1} \cdots W_1 x, \quad W_i \in \mathbb{R}^{d_i \times d_{i-1}}$$

where $w = \{W_i\}_{i=1}^{k+1}$ is the set of weight matrices.

- ▶ $x \in \mathbb{R}^{d_0}$, $y \in \mathbb{R}^{d_{k+1}}$, $(x, y) \sim \mathcal{D}$. Assume $\mathbb{E}[xx^\top] = I$.
- ▶ Learning problem: minimize the *population risk*

$$L(w) := \mathbb{E}_{(x,y) \sim \mathcal{D}} [\|y - f_w(x)\|^2]$$

based on iid samples from the distribution.

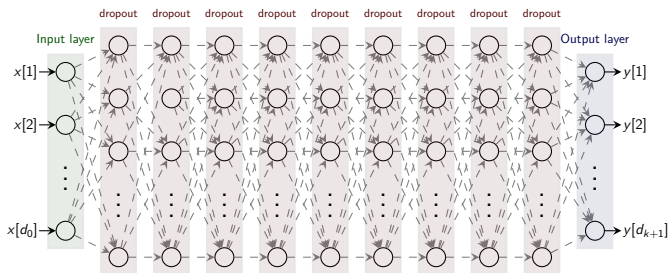


Problem Setup

- ▶ Network perturbed by dropping hidden nodes at random, computing

$$\bar{f}_w(\mathbf{x}) = W_{k+1}B_kW_k \cdots B_1W_1\mathbf{x},$$

where $B_i(j, j) = 0$ with probability $1 - \theta$, and $\frac{1}{\theta}$ with probability θ .



Problem Setup

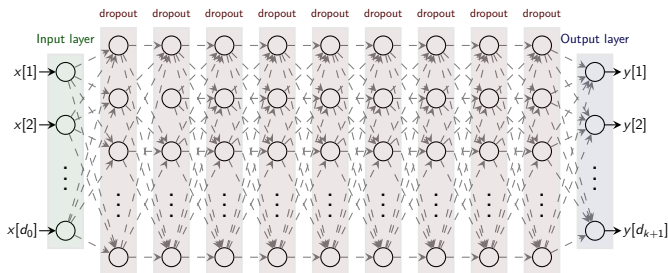
- ▶ Network perturbed by dropping hidden nodes at random, computing

$$\bar{f}_w(x) = W_{k+1}B_kW_k \cdots B_1W_1x,$$

where $B_i(j, j) = 0$ with probability $1 - \theta$, and $\frac{1}{\theta}$ with probability θ .

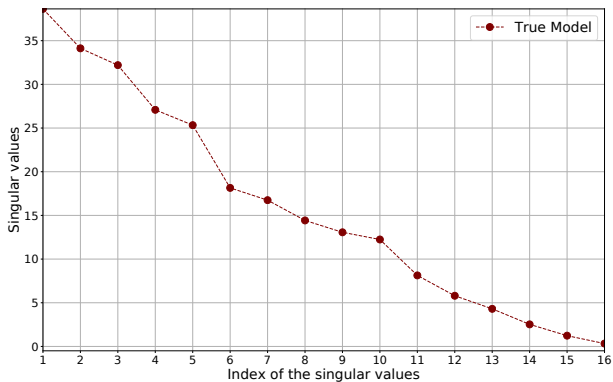
- ▶ dropout boils down to SGD on the *dropout objective*

$$L_\theta(w) := \mathbb{E}_{\{B_i\}, (x, y)} \|y - \bar{f}_w(x)\|^2$$



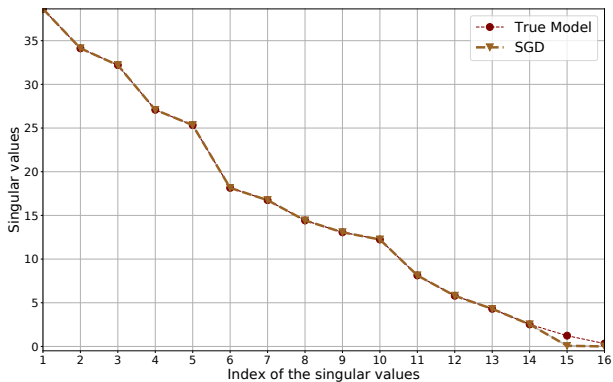
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



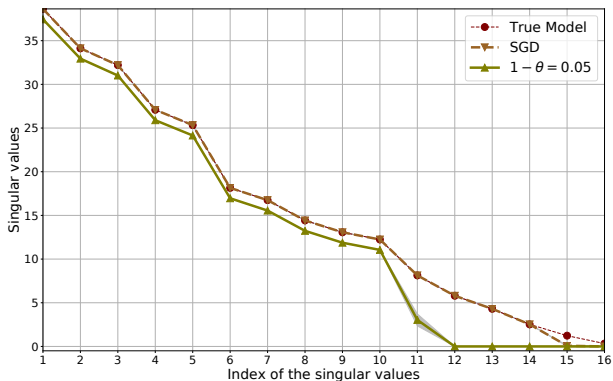
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



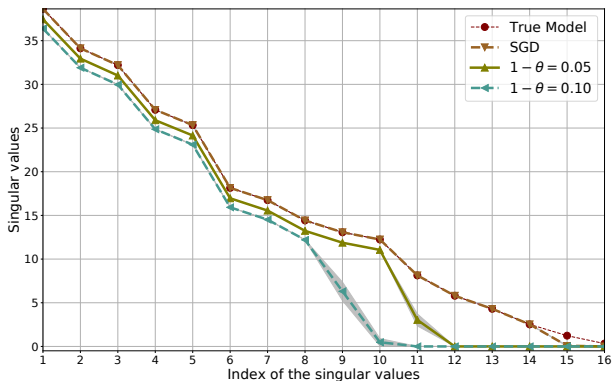
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



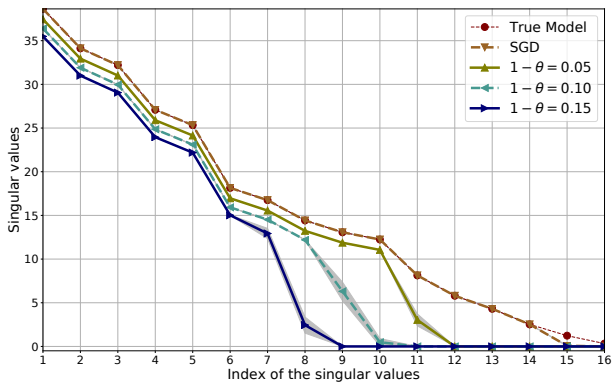
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



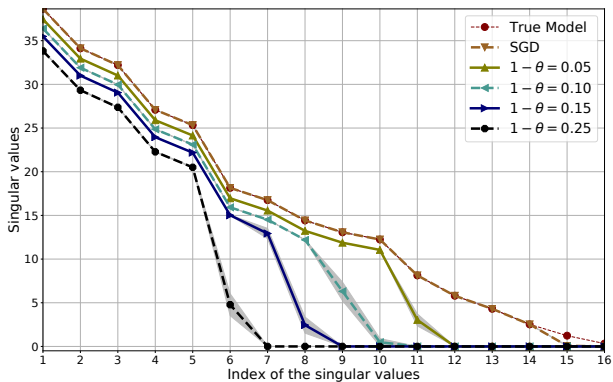
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



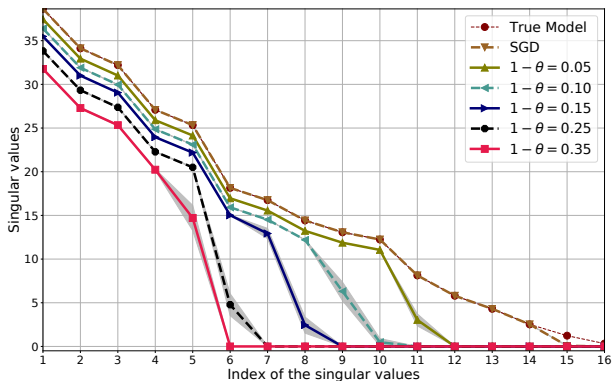
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



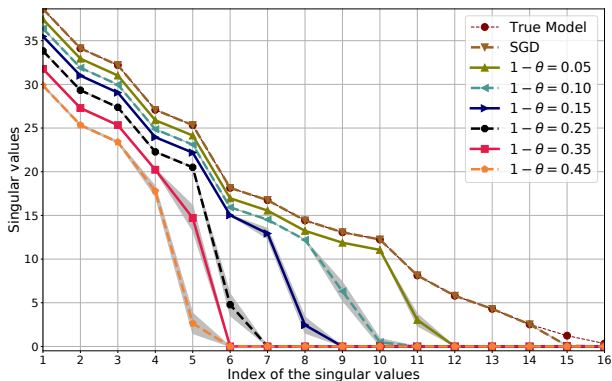
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



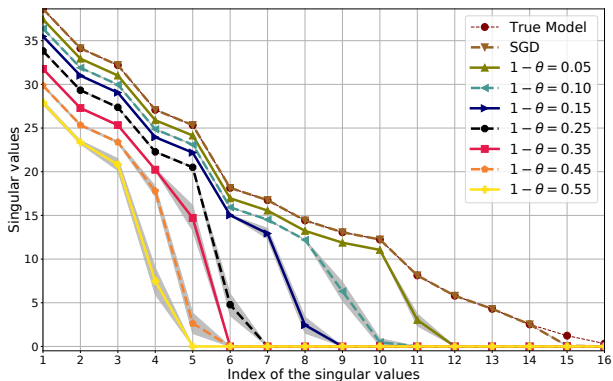
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



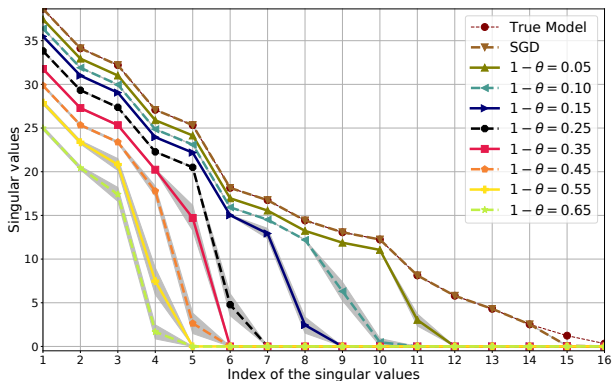
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



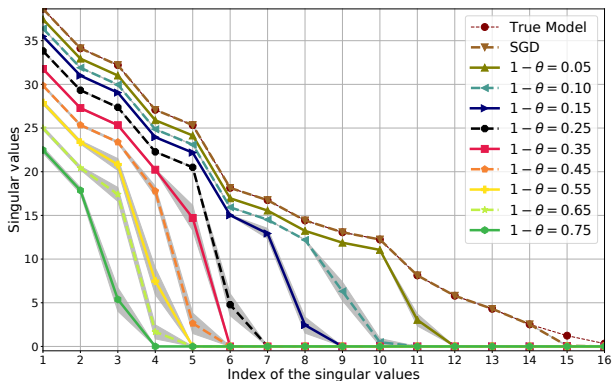
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



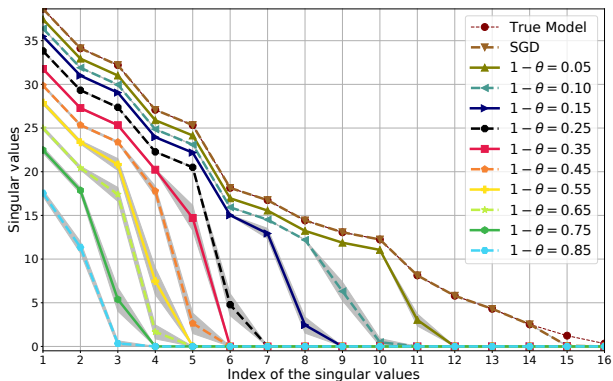
Empirical Observation

- ▶ 3-layer network with width/input/output dimensionality = 20.



Empirical Observation

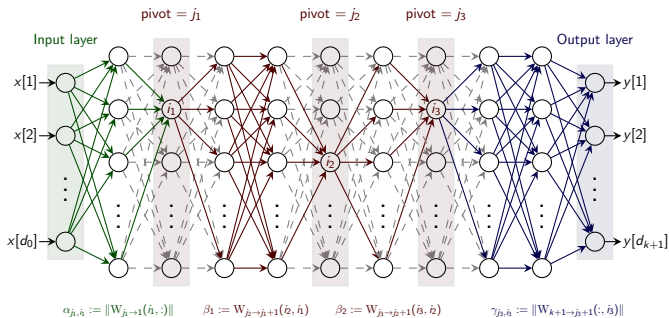
- ▶ 3-layer network with width/input/output dimensionality = 20.



Main Results

Explicit Regularizer

Give full characterization of $R(w) := L_{\theta}(w) - L(w)$



Main Results

Explicit Regularizer

Give full characterization of $R(\mathbf{w}) := L_\theta(\mathbf{w}) - L(\mathbf{w})$

Induced Regularizer

$$\Theta(M) := \min_{f_{\mathbf{w}}=M} R(\mathbf{w})$$

Main Results

Explicit Regularizer

Give full characterization of $R(w) := L_\theta(w) - L(w)$

Induced Regularizer

$$\Theta(M) := \min_{f_w=M} R(w)$$

- ▶ Multi-dimensional output

$$\Theta^{**}(f_w) = \nu_{\{d_i\}} \|f_w\|_*^2$$

Main Results

Explicit Regularizer

Give full characterization of $R(\mathbf{w}) := L_\theta(\mathbf{w}) - L(\mathbf{w})$

Induced Regularizer

$$\Theta(M) := \min_{f_{\mathbf{w}}=M} R(\mathbf{w})$$

- ▶ Multi-dimensional output

$$\Theta^{**}(f_{\mathbf{w}}) = \nu_{\{d_i\}} \|f_{\mathbf{w}}\|_*^2$$

- ▶ One-dimensional output

$$\Theta(f_{\mathbf{w}}) = \Theta^{**}(f_{\mathbf{w}}) = \nu_{\{d_i\}} \|f_{\mathbf{w}}\|^2$$

Main Results

Explicit Regularizer

Give full characterization of $R(\mathbf{w}) := L_\theta(\mathbf{w}) - L(\mathbf{w})$

Induced Regularizer

$$\Theta(M) := \min_{f_{\mathbf{w}}=M} R(\mathbf{w})$$

- ▶ Multi-dimensional output

$$\Theta^{**}(f_{\mathbf{w}}) = \nu_{\{d_i\}} \|f_{\mathbf{w}}\|_*^2$$

- ▶ One-dimensional output

$$\Theta(f_{\mathbf{w}}) = \Theta^{**}(f_{\mathbf{w}}) = \nu_{\{d_i\}} \|f_{\mathbf{w}}\|^2$$

Effective Regularization Parameter

$\nu_{\{d_i\}}$ increases with depth and decreases with width

deeper and narrower networks are more biased towards low-rank solutions

Thanks for your attention!

Stop by [Poster 79](#) for more information.