

Improved Dropout for Shallow and Deep Learning

¹Zhe Li, ²Boqing Gong, ¹Tianbao Yang

¹The University of Iowa, ²University of Central Florida

Main contribution

- Proposed a multinomial dropout for shallow learning.
- Demonstrated that this proposed distribution-dependent dropout leads to a faster convergence and a smaller generalization error through the risk bound analysis.
- Proposed an efficient evolutionary dropout for deep learning.
- Justified the proposed dropouts for both shallow and deep learning empirically.

Problem Setup

- Let (\mathbf{x}, y) denote a feature vector and a label, where $\mathbf{x} \in \mathbb{R}^d$ and $y \in \mathcal{Y}$.
- Denote by \mathcal{P} the joint distribution of (\mathbf{x}, y) and by \mathcal{D} the marginal distribution of \mathbf{x} .
- The goal is to learn a linear prediction function ($f(x) = \mathbf{w}^\top \mathbf{x}$) that minimizes the expected risk (considering loss function $\ell(\cdot, y)$):

$$\min_{\mathbf{w} \in \mathbb{R}^d} \mathcal{L}(\mathbf{w}) \triangleq \mathbb{E}_{\mathcal{P}}[\ell(\mathbf{w}^\top \mathbf{x}, y)] \quad (1)$$

- Denote by $\epsilon \sim \mathcal{M}$ a dropout noise vector of dimension d .
- The corrupted feature vector is given by $\hat{\mathbf{x}} = \mathbf{x} \circ \epsilon$, where the operator \circ represents the element-wise multiplication.
- Denote by $\hat{\mathcal{P}}$ the joint distribution of the new data $(\hat{\mathbf{x}}, y)$ and by $\hat{\mathcal{D}}$ the marginal distribution of $\hat{\mathbf{x}}$.
- With the corrupted data, the risk minimization becomes

$$\min_{\mathbf{w} \in \mathbb{R}^d} \hat{\mathcal{L}}(\mathbf{w}) \triangleq \mathbb{E}_{\hat{\mathcal{P}}}[\ell(\mathbf{w}^\top (\mathbf{x} \circ \epsilon), y)] \quad (2)$$

Learning with Multinomial Dropout

Definition 1. A **multinomial dropout** is defined as $\hat{\mathbf{x}} = \mathbf{x} \circ \epsilon$, where $\epsilon_i = \frac{m_i}{kp_i}$, $i \in [d]$ and $\{m_1, \dots, m_d\}$ follow a multinomial distribution $Mult(p_1, \dots, p_d; k)$ with $\sum_{i=1}^d p_i = 1$ and $p_i \geq 0$.

Proposition 1. If $\ell(z, y) = \log(1 + \exp(-yz))$, then

$$\mathbb{E}_{\hat{\mathcal{P}}}[\ell(\mathbf{w}^\top \hat{\mathbf{x}}, y)] = \mathbb{E}_{\mathcal{P}}[\ell(\mathbf{w}^\top \mathbf{x}, y)] + R_{\mathcal{D}, \mathcal{M}}(\mathbf{w})$$

where \mathcal{M} denotes the distribution of ϵ and $R_{\mathcal{D}, \mathcal{M}}(\mathbf{w}) = \mathbb{E}_{\mathcal{D}, \mathcal{M}} \left[\log \frac{\exp(\mathbf{w}^\top \mathbf{x} \circ \epsilon) + \exp(-\mathbf{w}^\top \mathbf{x} \circ \epsilon)}{\exp(\mathbf{w}^\top \mathbf{x}/2) + \exp(-\mathbf{w}^\top \mathbf{x}/2)} \right]$. Dropout is a data-dependent regularizer

Learning with Multinomial Dropout:

- Give the initial solution \mathbf{w}_1 .
- Update the model at t^{th} iteration: $\mathbf{w}_{t+1} = \mathbf{w}_t - \eta_t \nabla \ell(\mathbf{w}_t^\top (\mathbf{x}_t \circ \epsilon_t), y_t)$
- Output the final solution: $\hat{\mathbf{w}}_n = \frac{1}{n} \sum_{t=1}^n \mathbf{w}_t$

Improved Dropout for Shallow Learning

Risk Bound of $\hat{\mathbf{w}}_n$ in Expectation

Theorem 1: Let $\mathcal{L}(\mathbf{w})$ be the expected risk of \mathbf{w} defined in (1). Assume $\mathbb{E}_{\hat{\mathcal{D}}}[\|\mathbf{x} \circ \epsilon\|_2^2] \leq B^2$ and $\ell(z, y)$ is convex and G -Lipschitz continuous. For any $\|\mathbf{w}_*\|_2 \leq r$, by appropriately choosing η , we can have

$$\mathbb{E}[\mathcal{L}(\hat{\mathbf{w}}_n) + R_{\mathcal{D}, \mathcal{M}}(\hat{\mathbf{w}}_n)] \leq \mathcal{L}(\mathbf{w}_*) + R_{\mathcal{D}, \mathcal{M}}(\mathbf{w}_*) + \frac{GBr}{\sqrt{n}}$$

Theoretically, minimizing the term $\mathbb{E}_{\hat{\mathcal{D}}}[\|\mathbf{x} \circ \epsilon\|_2^2]$ and the relaxed upper bound of term $R_{\mathcal{D}, \mathcal{M}}(\mathbf{w}_*)$ yields the optimal sampling probabilities:

$$p_i^* = \frac{\sqrt{\mathbb{E}_{\mathcal{D}}[x_i^2]}}{\sum_{j=1}^d \sqrt{\mathbb{E}_{\mathcal{D}}[x_j^2]}}, i = 1, \dots, d \quad (3)$$

Practically, we use the empirical second-order statistics to compute the probabilities:

$$p_i = \frac{\sqrt{\frac{1}{n} \sum_{j=1}^n [[\mathbf{x}_j]_i^2]}}{\sum_{i'=1}^d \sqrt{\frac{1}{n} \sum_{j=1}^n [[\mathbf{x}_j]_{i'}^2]}}, i = 1, \dots, d \quad (4)$$

Improved Dropout for Deep Learning

Let $X^l = (\mathbf{x}_1^l, \dots, \mathbf{x}_m^l)$ denote the outputs of the l^{th} layer for a mini-batch of m examples, calculate the probabilities for dropout by

$$p_i^l = \frac{\sqrt{\frac{1}{m} \sum_{j=1}^m [[\mathbf{x}_j^l]_i^2]}}{\sum_{i'=1}^d \sqrt{\frac{1}{m} \sum_{j=1}^m [[\mathbf{x}_j^l]_{i'}^2]}}, i = 1, \dots, d \quad (5)$$

Evolutional Dropout for Deep Learning

Input: a batch of outputs of a layer: $X^l = (\mathbf{x}_1^l, \dots, \mathbf{x}_m^l)$ and dropout level parameter $k \in [0, d]$

Output: $\hat{X}^l = X^l \circ \Sigma^l$
Compute sampling probabilities by (5)

For $j = 1, \dots, m$

Sample $\mathbf{m}_j^l \sim Mult(p_1^l, \dots, p_d^l; k)$

Construct $\epsilon_j^l = \frac{\mathbf{m}_j^l}{k \mathbf{p}^l} \in \mathbb{R}^d$, where $\mathbf{p}^l = (p_1^l, \dots, p_d^l)^\top$

Let $\Sigma^l = (\epsilon_1^l, \dots, \epsilon_m^l)$ and compute $\hat{X}^l = X^l \circ \Sigma^l$

Figure 1: Evolutional Dropout applied to a layer over a mini-batch

Remark 1. Similar to Batch Normalization, evolutionary dropout can also address the internal covariate shift issue by adapting the sampling probabilities to the evolving distribution of layers' output.

Experimental Results

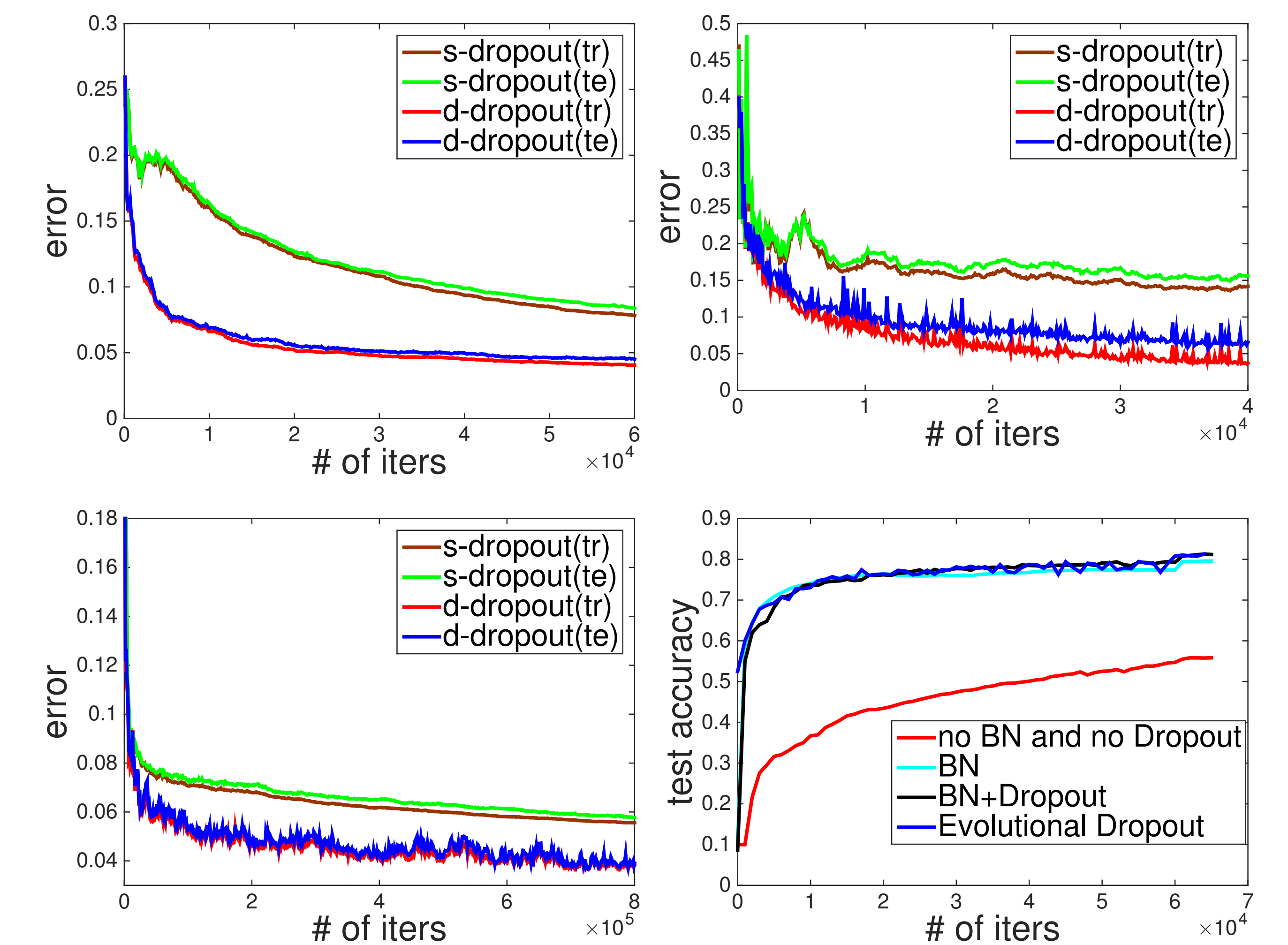


Figure 2: data-dependent dropout vs. standard dropout on three datasets (real-sim, news20 and RCV1) for logistic regression; Lower Right Corner: Evolutional dropout vs BN on CIFAR-10.

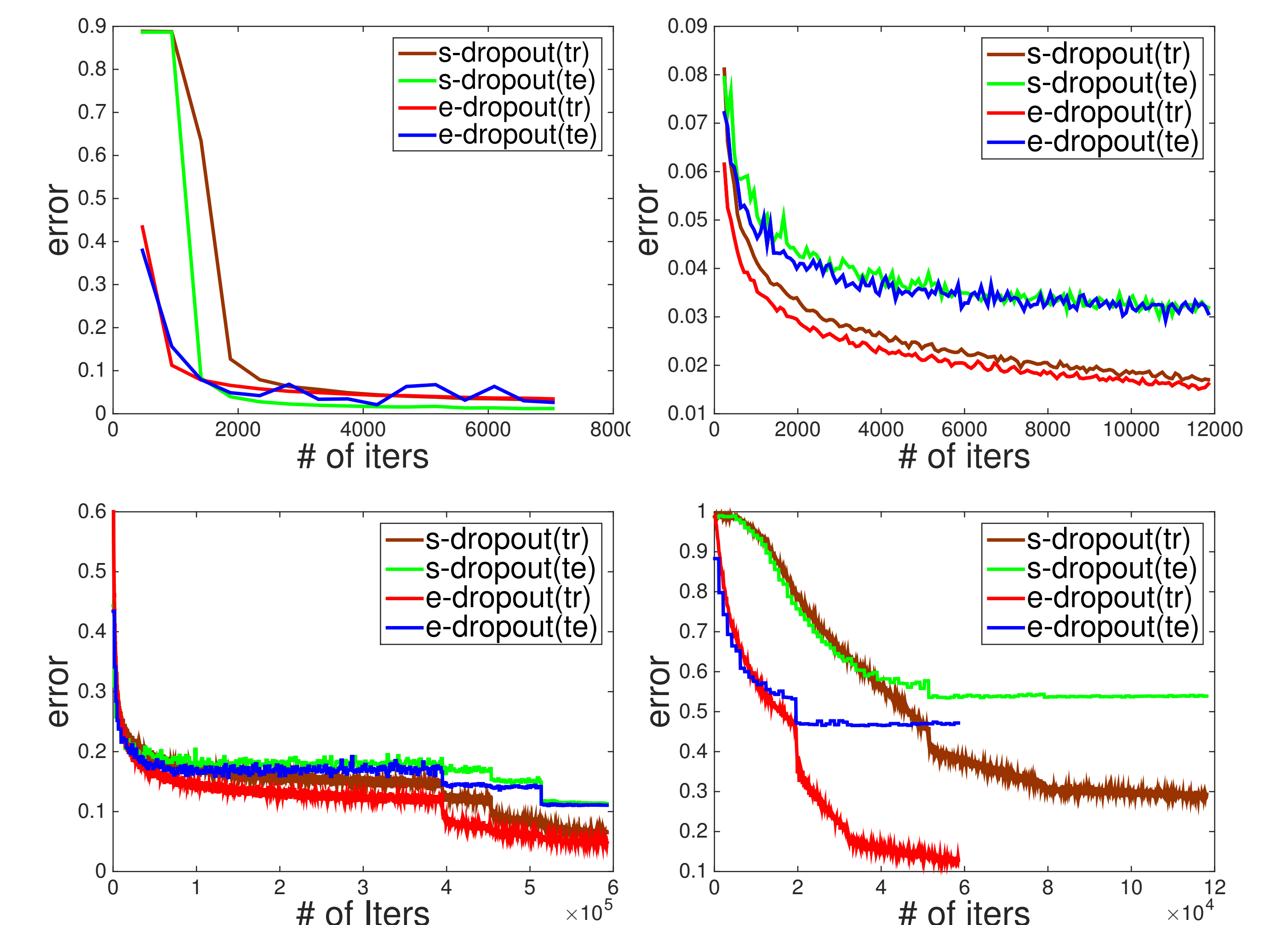


Figure 3: Evolutional dropout vs. standard dropout on four benchmark datasets (MNIST, SVHN, CIFAR-10 and CIFAR-100) for deep learning