# A Random Matrix Analysis of Learning with $\alpha$-Dropout
## Artemiss Workshop ICML 2020

**MEA. Seddik[12*], R. Couillet[23], M. Tamaazousti[1]**

[1]CEA List, France
[2]CentraleSupélec, L2S, France
[3]GIPSA Lab Grenoble-Alpes University, France

[*]`http://melaseddik.github.io/`

July 15, 2020

# Abstract

**Context:**

- ▶ Study of a one-hidden-layer network with $\alpha$-**Dropout**.

**Motivation:**

- ▶ Classical Dropout[1] corresponds to *zero-imputation*.
- ▶ *Zero-imputation* alter neural networks performances[2].

**Results:**

- ▶ **Asymptotic generalization performances** on a binary classification problem.
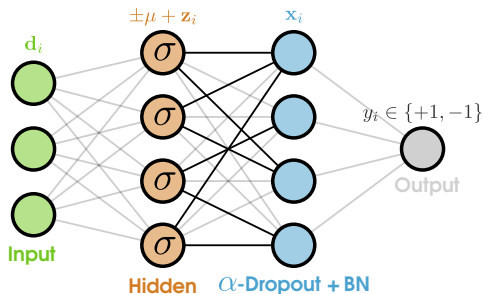- ▶ An aftermath analysis exhibits $\alpha \neq 0$ which improves generalization.

---

[1]Srivastava et al., *Dropout: a simple way to prevent neural networks from overfitting*. JMLR 2014.

[2]Yi et al., *Why not to use zero imputation? correcting sparsity bias in training neural networks*. ICLR 2019.

## Model and Problem Statement

Let $\boldsymbol{d}_1, \ldots, \boldsymbol{d}_n \in \mathbb{R}^q$ in two classes $\mathcal{C}_1$ and $\mathcal{C}_2$, and $\sigma : \mathbb{R}^q \to \mathbb{R}^p$ s.t. for $\boldsymbol{d}_i \in \mathcal{C}_a$

$$\mathbb{E}[\sigma(\boldsymbol{d}_i)] = (-1)^a \boldsymbol{\mu} \qquad \mathbb{E}[\sigma(\boldsymbol{d}_i)\sigma(\boldsymbol{d}_i)^\mathsf{T}] = \boldsymbol{I}_p + \boldsymbol{\mu}\boldsymbol{\mu}^\mathsf{T}$$



After the $\alpha$-**Dropout layer and BN**, the features matrix $\boldsymbol{X}_{\alpha,\varepsilon} = [\boldsymbol{x}_1, \ldots, \boldsymbol{x}_n] \in \mathcal{M}_{p,n}$ is

$$\boldsymbol{X}_{\alpha,\varepsilon} = \frac{(\boldsymbol{B}_\varepsilon \odot (\boldsymbol{Z} + \boldsymbol{\mu}\boldsymbol{y}^\mathsf{T})) \boldsymbol{P}_n + \alpha \boldsymbol{B}_\varepsilon \boldsymbol{P}_n}{\sqrt{\varepsilon + \alpha^2 \varepsilon (1 - \varepsilon)}}$$

with $[\boldsymbol{B}_\varepsilon]_{ij} \sim \mathrm{Ber}(\varepsilon)$, $Z_{ij} \sim \mathcal{N}(0, 1)$ and $\boldsymbol{P}_n = \boldsymbol{I}_n - \frac{1}{n}\boldsymbol{1}_n\boldsymbol{1}_n^\mathsf{T}$.

# Learning with $\alpha$-Dropout

We consider the Ridge-classifier with $\ell_2$-loss

$$\mathcal{E}(\boldsymbol{w}) = \frac{1}{n}\|\boldsymbol{y} - \boldsymbol{X}_{\alpha,\varepsilon}^{\mathsf{T}}\boldsymbol{w}\|^2 + \gamma\|\boldsymbol{w}\|^2$$

The solution of which is explicitly given by, for $z \in \mathbb{C} \setminus \mathbb{R}^-$

$$\boldsymbol{w} = \frac{1}{n}\boldsymbol{Q}(\gamma)\boldsymbol{X}_{\alpha,\varepsilon}\boldsymbol{y}, \qquad \boldsymbol{Q}(z) \equiv \left(\frac{1}{n}\boldsymbol{X}_{\alpha,\varepsilon}\boldsymbol{X}_{\alpha,\varepsilon}^{\mathsf{T}} + z\boldsymbol{I}_p\right)^{-1}$$

▶ The corresponding (hard) decision function is

$$g(\boldsymbol{x}) \equiv \boldsymbol{x}^{\mathsf{T}}\boldsymbol{w} = \frac{1}{n}\boldsymbol{x}^{\mathsf{T}}\boldsymbol{Q}(\gamma)\boldsymbol{X}_{\alpha,\varepsilon}\boldsymbol{y} \underset{\mathcal{C}_2}{\overset{\mathcal{C}_1}{\lessgtr}} 0$$

## Assumptions (Growth rate)

As $n \to \infty$,

1. $\frac{q}{n} \to r \in (0,\infty)$ and $\frac{p}{n} \to c \in (0,\infty)$;

2. For $a \in \{1,2\}$, $\frac{n_a}{n} \to c_a \in (0,1)$;

3. $\|\boldsymbol{\mu}\| = \mathcal{O}(1)$.

## Main Results

### Deterministic equivalent of $Q(z)$

Under the previous Assumptions,

$$Q(z) \leftrightarrow \bar{Q}(z) \equiv \mathcal{D}_z - \frac{\frac{\varepsilon}{1+\alpha^2(1-\varepsilon)}\mathcal{D}_z\mu\mu^\intercal\mathcal{D}_z}{1 + cq(z) + \frac{\varepsilon}{1+\alpha^2(1-\varepsilon)}\mu^\intercal\mathcal{D}_z\mu},$$

where $\mathcal{D}_z \equiv q(z)\mathrm{diag}\left\{\frac{1+cq(z)}{1+cq(z)+\frac{(1-\varepsilon)q(z)}{1+\alpha^2(1-\varepsilon)}\mu_i^2}\right\}_{i=1}^p$ with $q(z) \equiv \frac{c-z-1+\sqrt{(c-z-1)^2+4zc}}{2zc}$.

### Gaussian Approximation of $g(x)$

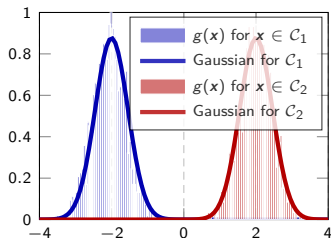Under the previous Assumptions, for $x \in \mathcal{C}_a$ with $a \in \{1,2\}$,

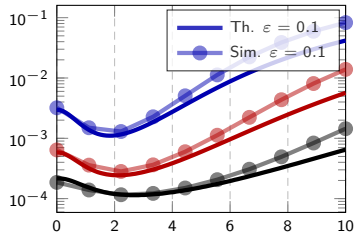$$\nu^{-\frac{1}{2}}\left(g(x) - m_a\right) \xrightarrow{\mathcal{D}} \mathcal{N}(0,1)$$

where

$$m_a \equiv (-1)^a \sqrt{\frac{\varepsilon}{1 + \alpha^2(1-\varepsilon)}} \frac{\mu^\intercal \bar{Q}(\gamma)\mu}{1 + \delta(\gamma)}$$

$$\nu \equiv \frac{1}{(1+\delta(\gamma))^2}\left(\eta(C_1) + \frac{\varepsilon}{1+\alpha^2(1-\varepsilon)} \times \left[\mu^\intercal\left(\Delta(C_1) - \bar{Q}(\gamma)\right)\mu - \frac{2\,\eta(C_1)\mu^\intercal\bar{Q}(\gamma)\mu}{1+\delta(\gamma)}\right]\right)$$

## Take Away Messages



Test (generalization) scores $g(x)$

Test Error in terms of $\alpha$

**Highlights:**

- ▶ Existence of $\alpha \neq 0$ which minimizes **the test error**.
- ▶ In our setting, such $\alpha$ satisfies $\frac{1}{m_a} \frac{\partial m_a}{\partial \alpha} = \frac{1}{\sqrt{\nu}} \frac{\partial \sqrt{\nu}}{\partial \alpha}$.

**Perspectives:**

- ▶ Extend the analysis to a $k$-**class model** with $\alpha_\ell$'s for each class.
- ▶ Validation of the $\alpha$-Dropout approach with **real data**.
- ▶ Extend to **multi-layers** networks.