

Demystifying Batch Normalization in ReLU Networks: Equivalent Convex Optimization Models and Implicit Regularization

Tolga Ergen*, Arda Sahiner*, Batu Ozturkler, John Pauly, Morteza Mardani, Mert Pilanci

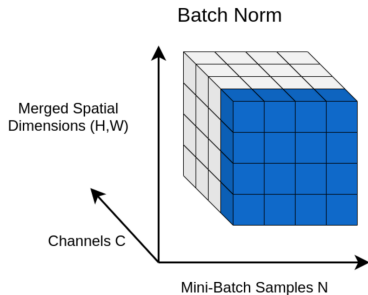
🐦 @tolgaergen_ 🏠 stanford.edu/~ergen/ ✉ ergen@stanford.edu



Why do we need Batch Normalization (BN)?

BN transforms a batch of data to zero mean and standard deviation one, and has two trainable parameters α , γ :

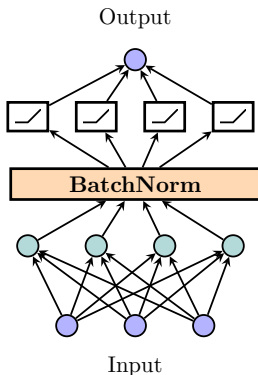
$$\text{BN}_{\gamma, \alpha}(\mathbf{x}) = \frac{(\mathbf{I}_d - \frac{1}{d} \mathbf{1} \mathbf{1}^T) \mathbf{x}}{\|(\mathbf{I}_d - \frac{1}{d} \mathbf{1} \mathbf{1}^T) \mathbf{x}\|_2} \gamma + \alpha$$



BN **stabilizes** and **accelerates** training of deep neural networks

Two-layer Neural Networks with ReLU Activation

Model:



Notation:

$X \in \mathbb{R}^{n \times d}$: Data matrix

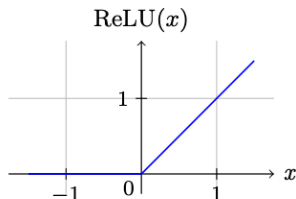
$Y \in \mathbb{R}^{n \times C}$: Label matrix

$\mathcal{L}(\cdot, \cdot)$: Arbitrary convex loss function

$\beta > 0$: Regularization coefficient

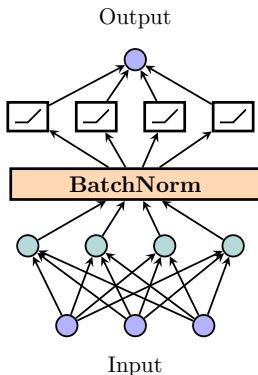
$W^{(1)} \in \mathbb{R}^{d \times m}, W^{(2)} \in \mathbb{R}^{m \times C}$: Layer weights

θ : Represents all trainable parameters



Two-layer Neural Networks with ReLU Activation

Model:



Notation:

$X \in \mathbb{R}^{n \times d}$: Data matrix

$Y \in \mathbb{R}^{n \times c}$: Label matrix

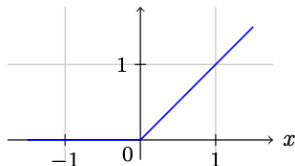
$\mathcal{L}(\cdot, \cdot)$: Arbitrary convex loss function

$\beta > 0$: Regularization coefficient

$W^{(1)} \in \mathbb{R}^{d \times m}, W^{(2)} \in \mathbb{R}^{m \times c}$: Layer weights

θ : Represents all trainable parameters

$\text{ReLU}(x)$



Optimization problem:

$$p_{\text{non-convex}} := \min_{\theta} \mathcal{L} \left(\phi \left(\text{BN}_{\gamma, \alpha} \left(XW^{(1)} \right) \right) W^{(2)}, Y \right) + \frac{\beta}{2} \|\theta\|_2^2$$

where $\phi(x) = \text{ReLU}(x) = (x)_+$ and $\mathcal{L}(\cdot, \cdot)$ is any convex loss

Two-layer Neural Networks with ReLU Activation

► High-dimensional regime ($n \leq d$)

Theorem

Let $n \leq d$ and \mathbf{X} is full row-rank, then an optimal solution is

$$\begin{aligned} (\mathbf{w}_j^{(1)*}, \mathbf{w}_j^{(2)*}) &= (\mathbf{x}^\dagger \mathbf{y}_j, (\|\mathbf{y}_j\|_2 - \beta)_+ \mathbf{e}_j) \\ (\gamma_j^*, \alpha_j^*) &= \left(\frac{\|\mathbf{y}_j - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{y}_j\|_2}{\|\mathbf{y}_j\|_2}, \frac{\mathbf{1}^T \mathbf{y}_j}{\sqrt{n} \|\mathbf{y}_j\|_2} \right), \forall j \in [C] \end{aligned}$$

where \mathbf{e}_j is the j^{th} ordinary basis vector.

Two-layer Neural Networks with ReLU Activation

- ▶ High-dimensional regime ($n \leq d$)

Theorem

Let $n \leq d$ and \mathbf{X} is full row-rank, then an optimal solution is

$$\begin{aligned} (\mathbf{w}_j^{(1)*}, \mathbf{w}_j^{(2)*}) &= (\mathbf{X}^\dagger \mathbf{y}_j, (\|\mathbf{y}_j\|_2 - \beta)_+ \mathbf{e}_j) \\ (\gamma_j^*, \alpha_j^*) &= \left(\frac{\|\mathbf{y}_j - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{y}_j\|_2}{\|\mathbf{y}_j\|_2}, \frac{\mathbf{1}^T \mathbf{y}_j}{\sqrt{n} \|\mathbf{y}_j\|_2} \right), \forall j \in [C] \end{aligned}$$

where \mathbf{e}_j is the j^{th} ordinary basis vector.

- ▶ Generic case (arbitrary n, d):

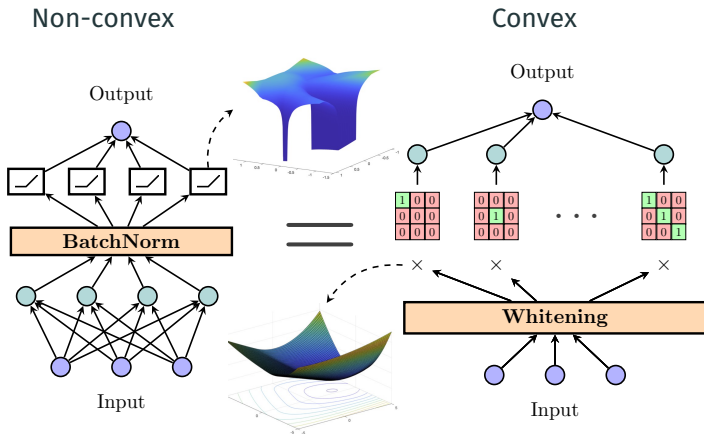
Theorem

Given the SVD of \mathbf{X} as $\mathbf{X} := \mathbf{U} \mathbf{\Sigma} \mathbf{V}^T$, the non-convex training problem is equivalent to

$$p_{\text{convex}} := \min_{\mathbf{S}_i} \mathcal{L} \left(\sum_{i=1}^P \mathbf{D}_i \mathbf{U} \mathbf{S}_i, \mathbf{Y} \right) + \beta \sum_{i=1}^P \|\mathbf{S}_i\|_{C_i, *}$$

where $\mathbf{D}_1, \dots, \mathbf{D}_P$ are fixed diagonal matrices, and $\|\cdot\|_{C_i, *}$ is a constrained version of the standard nuclear norm.

ReLU+BN \equiv Convex+Low-rank+Whitening



$$X \in \mathbb{R}^{n \times d} \xrightarrow{\text{ReLU+BN}} \tilde{X} = [D_1 U \quad \dots \quad D_p U] \in \mathbb{R}^{n \times dp}$$

ReLU+BN \equiv Low-rank convex model applied to whitened data \tilde{X}

Deep ReLU Networks

Model: $f_{\theta,L}(X) := \mathbf{A}^{(L-1)}\mathbf{W}^{(L)}$, where $\mathbf{A}^{(l)} := \left(\text{BN}_{\gamma,\alpha} \left(\mathbf{A}^{(l-1)}\mathbf{W}^{(l)} \right) \right)_+$

Theorem

Assume the network is overparameterized s.t. $\text{Range}(\mathbf{A}^{(L-2)}) = \mathbb{R}^n$, then optimal solution in closed-form is as follows

$$\begin{aligned} \left(\mathbf{w}_j^{(L-1)*}, \mathbf{w}_j^{(L)*} \right) &= \left(\mathbf{A}^{(L-2)\dagger} \mathbf{y}_j, (\|\mathbf{y}_j\|_2 - \beta)_+ \mathbf{e}_j \right) \\ \left(\gamma_j^{(L-1)*}, \alpha_j^{(L-1)*} \right) &= \left(\frac{\|\mathbf{y}_j - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{y}_j\|_2}{\|\mathbf{y}_j\|_2}, \frac{\mathbf{1}^T \mathbf{y}_j}{\sqrt{n} \|\mathbf{y}_j\|_2} \right), \forall j \in [C] \end{aligned}$$

where \mathbf{e}_j is the j^{th} ordinary basis vector.

Deep ReLU Networks

Model: $f_{\theta,L}(X) := \mathbf{A}^{(L-1)}\mathbf{W}^{(L)}$, where $\mathbf{A}^{(l)} := \left(\text{BN}_{\gamma,\alpha} \left(\mathbf{A}^{(l-1)}\mathbf{W}^{(l)} \right) \right)_+$

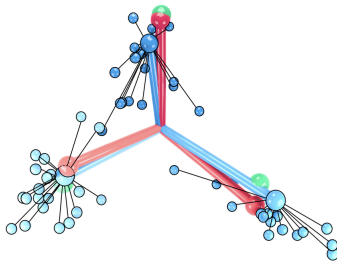
Theorem

Assume the network is overparameterized s.t. $\text{Range}(\mathbf{A}^{(L-2)}) = \mathbb{R}^n$, then optimal solution in closed-form is as follows

$$\begin{aligned} (\mathbf{w}_j^{(L-1)*}, \mathbf{w}_j^{(L)*}) &= (\mathbf{A}^{(L-2)\dagger} \mathbf{y}_j, (\|\mathbf{y}_j\|_2 - \beta)_+ \mathbf{e}_j) \\ (\gamma_j^{(L-1)*}, \alpha_j^{(L-1)*}) &= \left(\frac{\|\mathbf{y}_j - \frac{1}{n} \mathbf{1} \mathbf{1}^T \mathbf{y}_j\|_2}{\|\mathbf{y}_j\|_2}, \frac{\mathbf{1}^T \mathbf{y}_j}{\sqrt{n} \|\mathbf{y}_j\|_2} \right), \forall j \in [C] \end{aligned}$$

where \mathbf{e}_j is the j^{th} ordinary basis vector.

This also explains **Neural Collapse** in (Papayan et al., 2020)

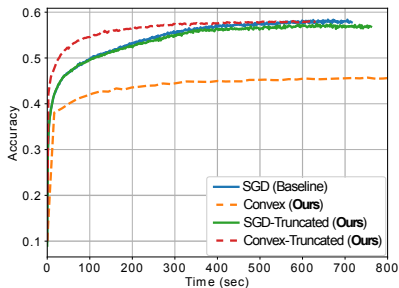


Numerical Experiments on Image Datasets

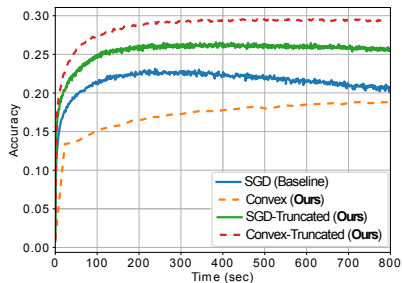
Truncation: Given the SVD $\mathbf{X} = \mathbf{U}\mathbf{\Sigma}\mathbf{V}^T$, the truncated data matrix is defined as

$$\hat{\mathbf{X}} := \mathbf{U}\mathcal{T}(\mathbf{\Sigma})\mathbf{V}^T,$$

with $\mathcal{T}_k(\mathbf{\Sigma})_{ii} := \sigma_i \mathbb{1}\{i \geq k\}$, which takes the k top singular values.



(a) CIFAR-10



(b) CIFAR-100

Figure 1: Image classification on CIFAR datasets

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied
 - don't need heuristics or hyperparameter search, e.g., learning rate and initialization

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied
 - don't need heuristics or hyperparameter search, e.g., learning rate and initialization
- ▶ BN networks seek **low-rank** structure via nuclear norm regularization

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied
 - don't need heuristics or hyperparameter search, e.g., learning rate and initialization
- ▶ BN networks seek **low-rank** structure via nuclear norm regularization
- ▶ BN induces **whitening** effect on the data

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied
 - don't need heuristics or hyperparameter search, e.g., learning rate and initialization
- ▶ BN networks seek **low-rank** structure via nuclear norm regularization
- ▶ BN induces **whitening** effect on the data

Future research directions:

- ▶ convex programs for deeper networks (more than 2 layers)

Takeaways and Open Problems

- ▶ ReLU network training with BN is a convex optimization problem in high dimension
 - convex optimization theory and solvers can be applied
 - don't need heuristics or hyperparameter search, e.g., learning rate and initialization
- ▶ BN networks seek **low-rank** structure via nuclear norm regularization
- ▶ BN induces **whitening** effect on the data

Future research directions:

- ▶ convex programs for deeper networks (more than 2 layers)
- ▶ extensions to various architectures such as GAN, RNN, ResNets ...

