



CVPR 2021 Tutorial

Normalization Techniques in Deep Learning: Methods, Analyses and Applications



Lei Huang

Beihang University, Beijing, China





Outline

01. Motivations of Normalization Techniques

02. Introduction of Normalization Methods

03. Analyses of Normalization

04. Applications of Normalization

Analyses of Normalization

- Topics of deep learning theory

Representation

Optimization

Generalization

BN property:

Control distribution

Stabilize training

Accelerate convergence

Improve generalization

Analyses of Normalization

- Topics of deep learning theory

Representation

Optimization

Generalization

BN property:

Control distribution

Stabilize training

Accelerate convergence

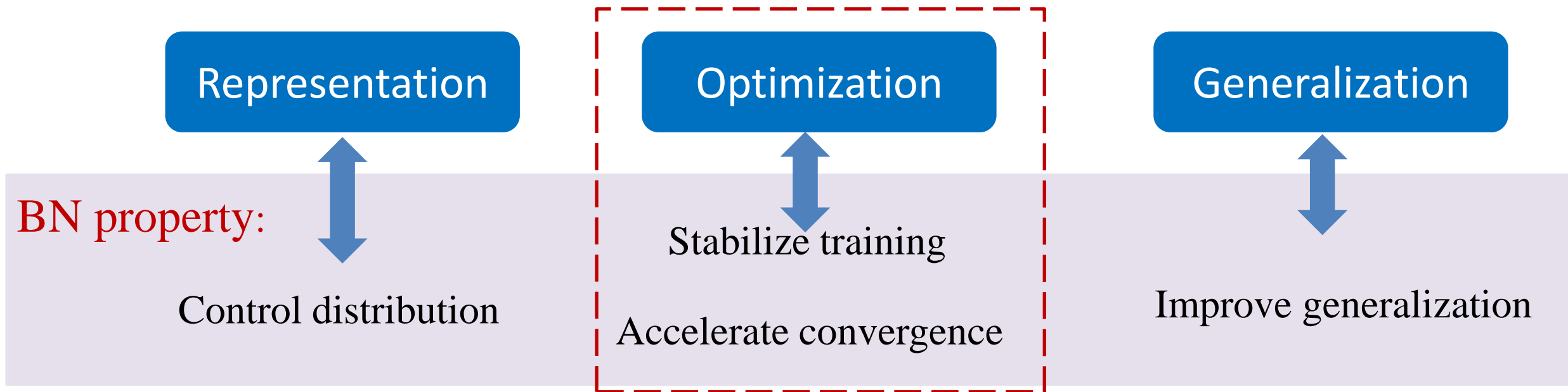
Improve generalization

- Disclaims in this tutorials:

- Not theoretic analysis in generalization bound or optimization theory
- But empirical analysis with high-level theoretic insights.

Analyses of Normalization

- Topics of deep learning theory



- Disclaims in this tutorials:
 - Not theoretic analysis in generalization bound or optimization theory
 - But empirical analysis with high-level theoretic insights.

Scale Invariant Analysis

- Scale invariant for weight in a normalization layer

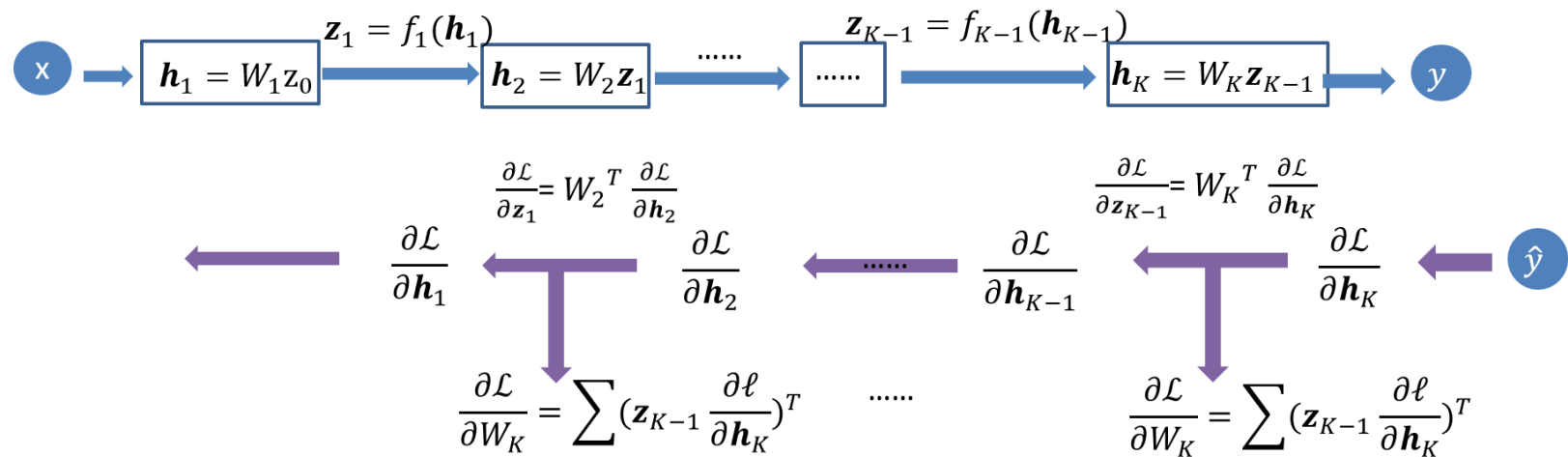


$$BN(\mathbf{x}; a\mathbf{W}) = BN(\mathbf{x}; \mathbf{W})$$
$$\frac{\partial BN(\mathbf{x}; a\mathbf{W})}{\partial(a\mathbf{W})} = \frac{1}{a} \frac{\partial BN(\mathbf{x}; \mathbf{W})}{\partial \mathbf{W}},$$

Applies also for LN or weight normalization

Scale Invariant Analysis

- Stabilize training in a network
 - Scale dynamics for a network



$$\widehat{W}_k = \alpha_k W_k$$

	Layer input	Output-gradient	Weight-gradient
Un-normalized	$\hat{x}_k = (\prod_{i=1}^k \alpha_i) x_k$	$\frac{\partial L}{\partial \hat{h}_k} = (\prod_{i=k+1}^K \alpha_i) \frac{\partial L}{\partial h_k}$	$\frac{\partial L}{\partial \widehat{W}_k} = (\prod_{i=1, i \neq k}^K \alpha_i) \frac{\partial L}{\partial W_k}$
normalized	$\hat{x}_k = x_k$	$\frac{\partial L}{\partial \hat{h}_k} = \frac{1}{\alpha_k} \frac{\partial L}{\partial h_k}$	$\frac{\partial L}{\partial \widehat{W}_k} = \frac{1}{\alpha_k} \frac{\partial L}{\partial W_k}$

Scale Invariant Analysis

- The auto-tuning effects of learning rate

Scale invariant

$$f(\lambda \mathbf{w}) = f(\mathbf{w})$$

$$\frac{\partial f(\lambda \mathbf{w})}{\partial \lambda \mathbf{w}} = \frac{1}{\lambda} \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}}$$

$$\left\langle \mathbf{w}, \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right\rangle = 0$$

Weight update:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \eta \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}}$$

$$\|\mathbf{w}_{t+1}\|_2^2 = \|\mathbf{w}_t\|_2^2 + \eta^2 \left\| \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right\|_2^2$$

$\mathbf{w} \uparrow$

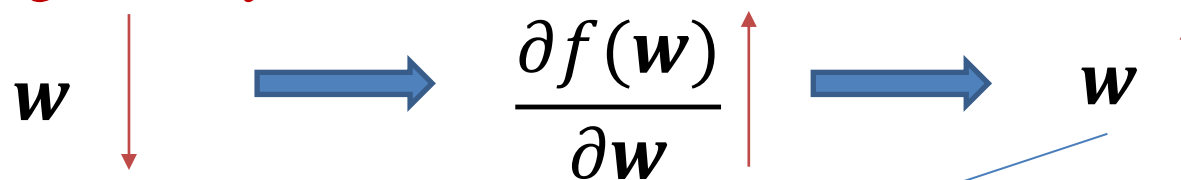
$$\frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \downarrow$$

Scale Invariant Analysis

- The auto-tuning effects + weight decay

$$\left\{ \begin{array}{l} f(\lambda \mathbf{w}) = f(\mathbf{w}) \\ \frac{\partial f(\lambda \mathbf{w})}{\partial \lambda \mathbf{w}} = \frac{1}{\lambda} \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \end{array} \right. \quad \left\{ \begin{array}{l} \langle \mathbf{w}, \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \rangle = 0 \\ \|\mathbf{w}_{t+1}\|_2^2 = \|\mathbf{w}_t\|_2^2 + \eta^2 \left\| \frac{\partial f(\mathbf{w})}{\partial \mathbf{w}} \right\|_2^2 \end{array} \right.$$

weight decay

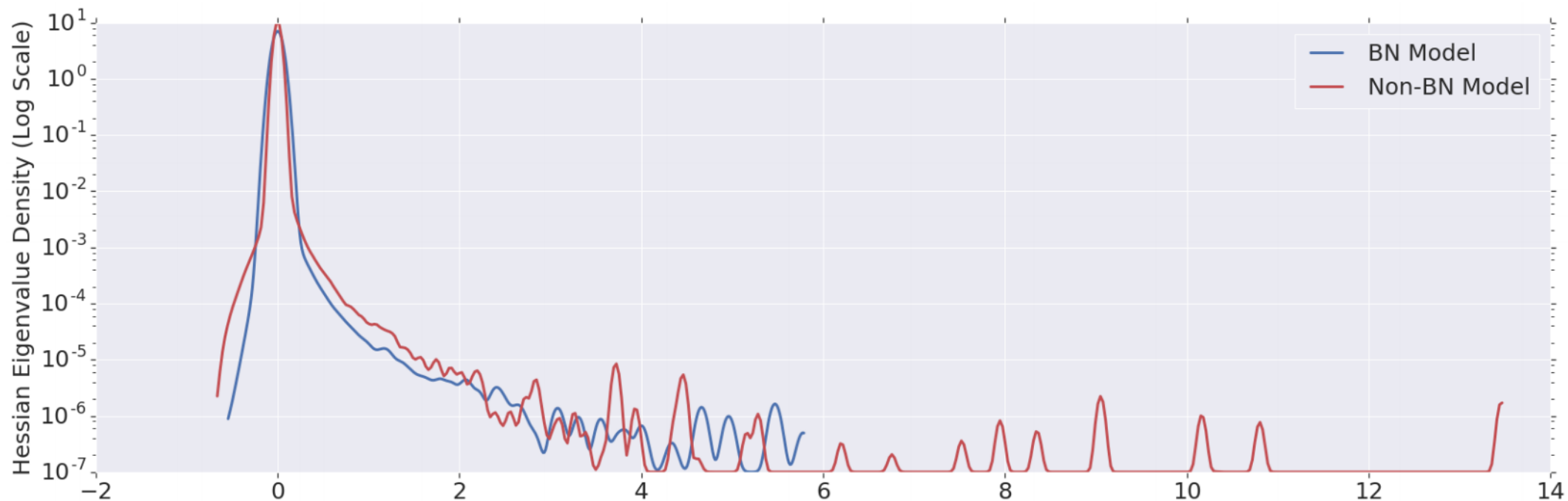


Equilibrium

Improved Conditioning in Optimization



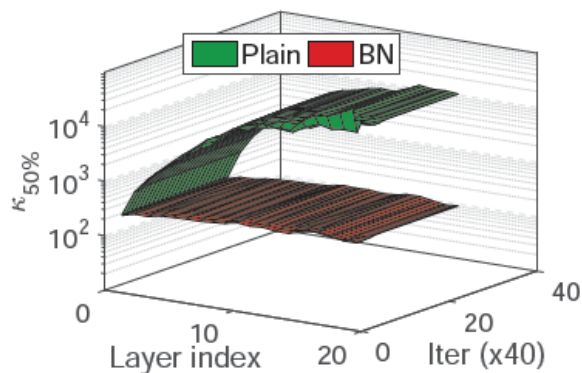
- The full Hessian spectrum



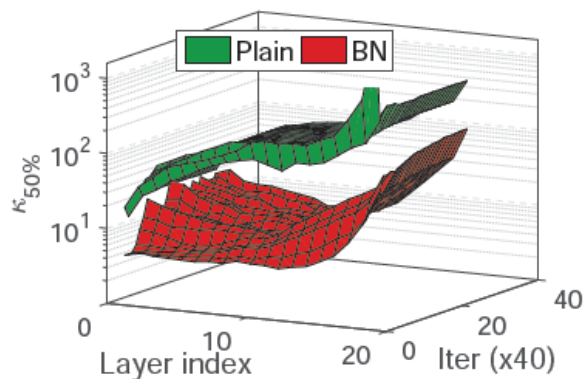
The eigenvalue comparison of the Hessian of the VGG network with BN (blue) and without BN (red)

Improved Conditioning in Optimization

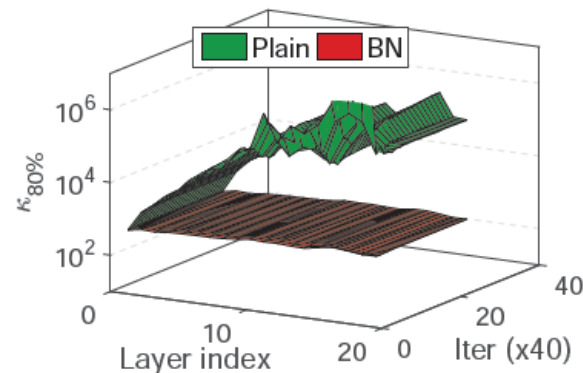
- The Layer-wise spectrum
 - Covariance of layer input Σ_x
 - Covariance of output-gradient $\Sigma_{\nabla h}$
 - Criteria 2: Σ_x and $\Sigma_{\nabla h}$ are well conditioned (**in layer**)



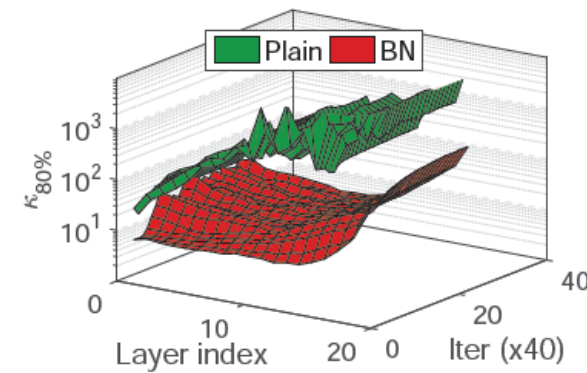
(a) $\kappa_{50\%}(\Sigma_x)$



(b) $\kappa_{50\%}(\Sigma_{\nabla h})$



(c) $\kappa_{80\%}(\Sigma_x)$



(d) $\kappa_{80\%}(\Sigma_{\nabla h})$



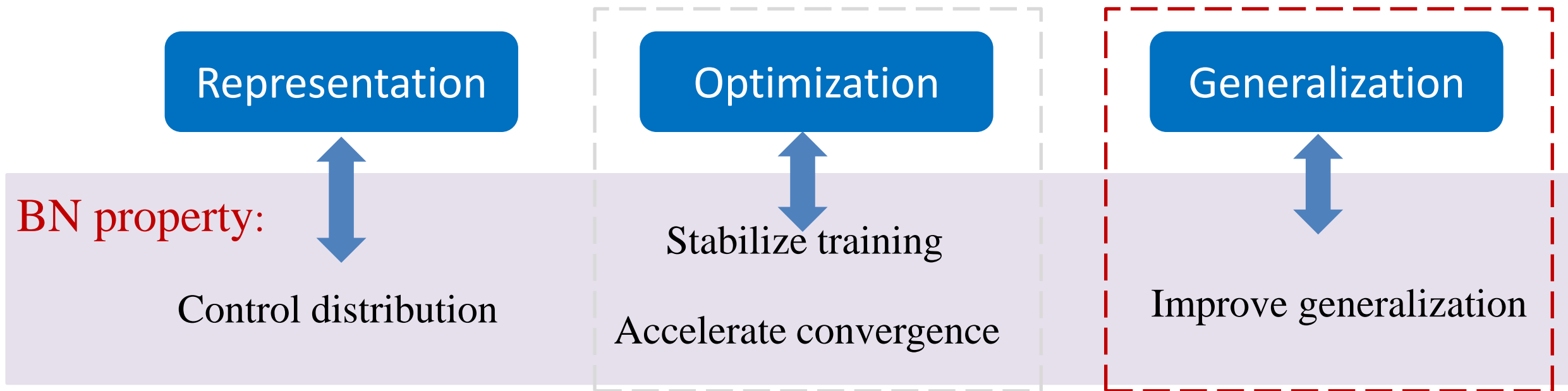
Theoretical Analysis under Assumptions



- Mean field theory of batch normalization
 - Yang et al, ICLR 2019
 - Wei et al, 2019
- Theoretical accelerated converge of BN on learning half-space with Gaussian inputs [Kohler et al, AISTATS 2019]
- Global convergence results for two-layer NN with ReLU and weight normalization [Dukler et al, ICML 2020]

Analyses of Normalization

- Topics of deep learning theory



- Disclaims in this tutorials:
 - Not theoretic analysis in generalization bound or optimization theory
 - But empirical analysis with high-level theoretic insights.

The stochasticity

- Introduced Stochasticity

- During training

$$\hat{X}_m = \frac{X_m - \mu_{X_m}}{\sigma_{X_m}}$$

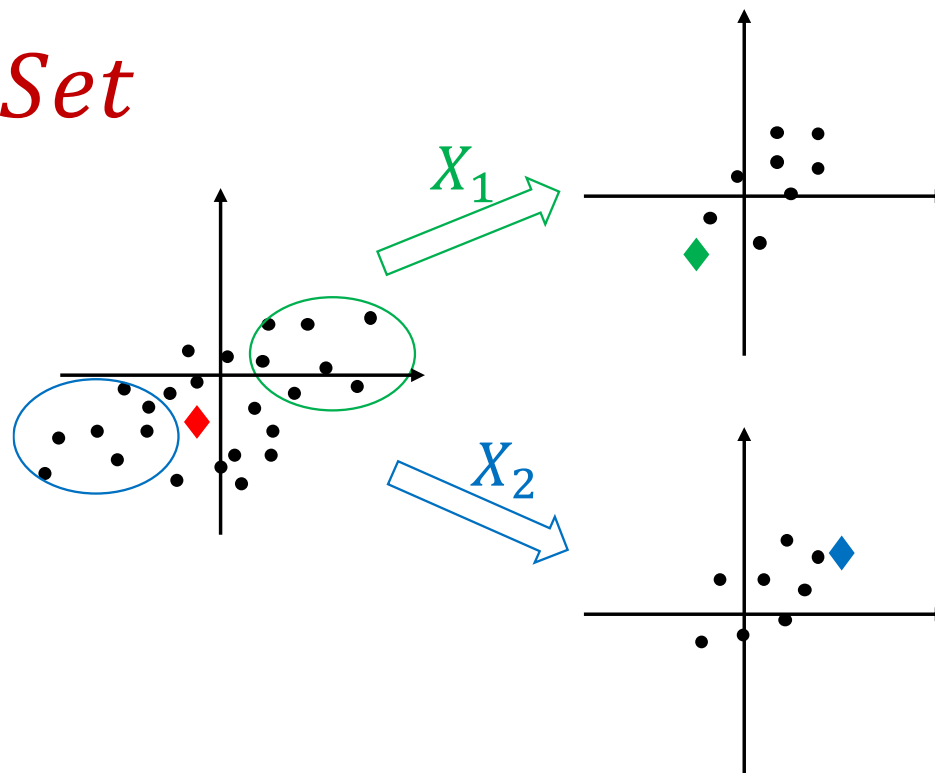
$X_m \sim \text{DataSet}$

- During inference

Population statistics $\{\hat{\mu}, \hat{\sigma}\}$:

$$\hat{\mu} = (1 - \lambda)\hat{\mu} + \lambda\mu_{X_m}$$

$$\hat{\sigma} = (1 - \lambda)\hat{\sigma} + \lambda\sigma_{X_m}$$

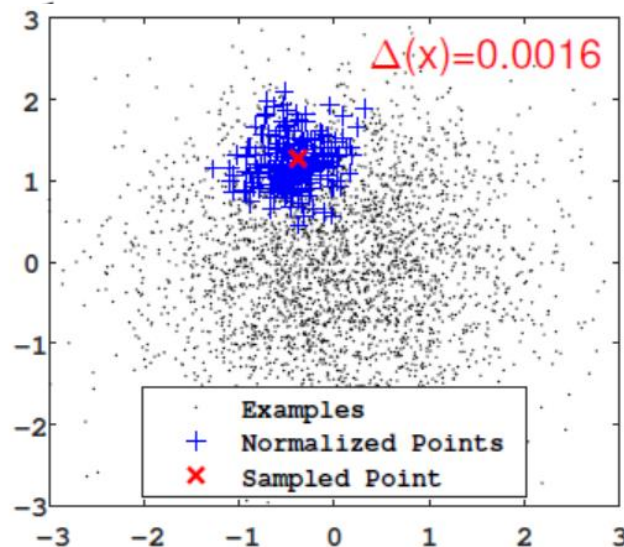


Quantitatively Evaluation

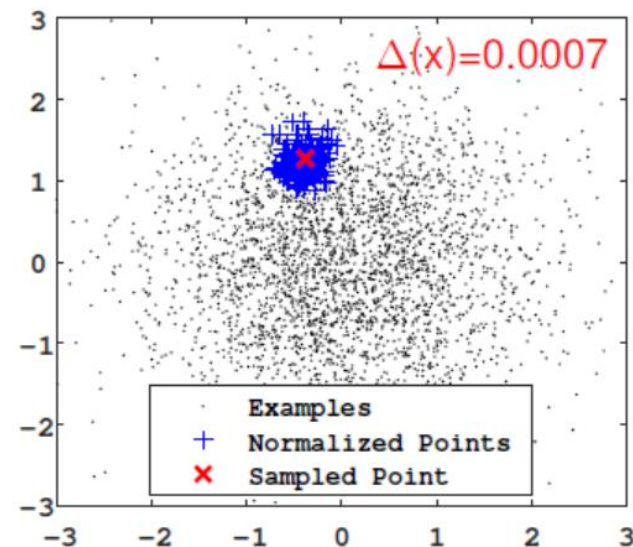
- Stochastic Normalization Disturbance (SND)

$$\text{Empirical SND: } \hat{\Delta}_{\mathbf{G}}(\mathbf{x}) = \frac{1}{s} \sum_{i=1}^s \|\mathbf{G}(\mathbf{X}_i^B; \mathbf{x}) - \frac{1}{s} \sum_{j=1}^s \mathbf{G}(\mathbf{X}_j^B; \mathbf{x})\|$$

$\{\mathbf{X}_j^B\}_{j=1}^s \sim \text{Dataset}$



(a) batch size of 16

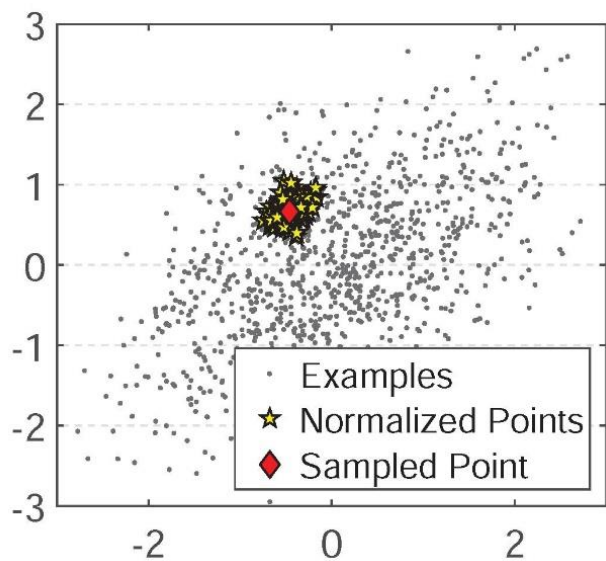


(b) batch size of 64

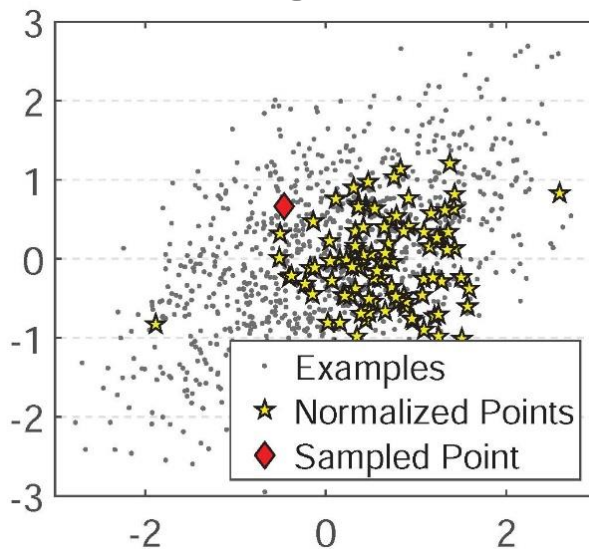
Qualitatively Visualization

- Stochastic outputs combining different mini-batches for different batch whitening

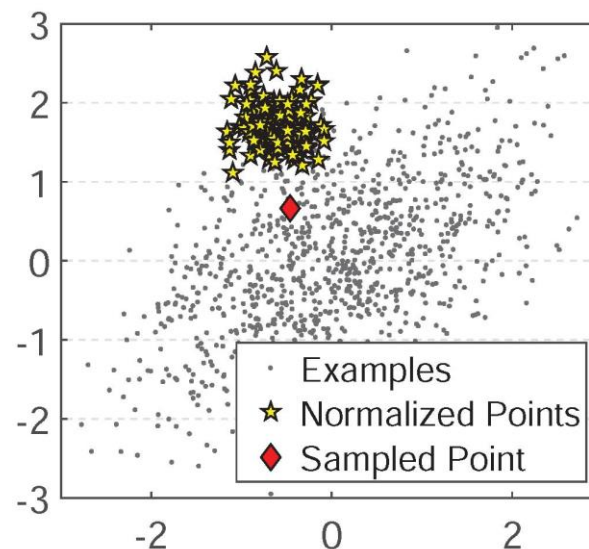
BN



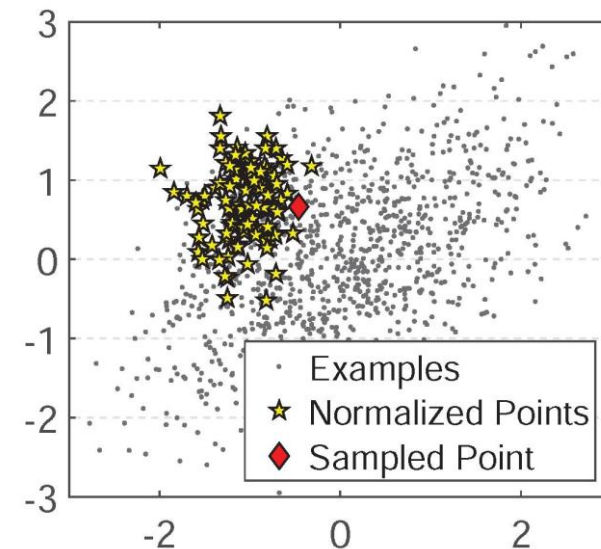
PCA



ZCA

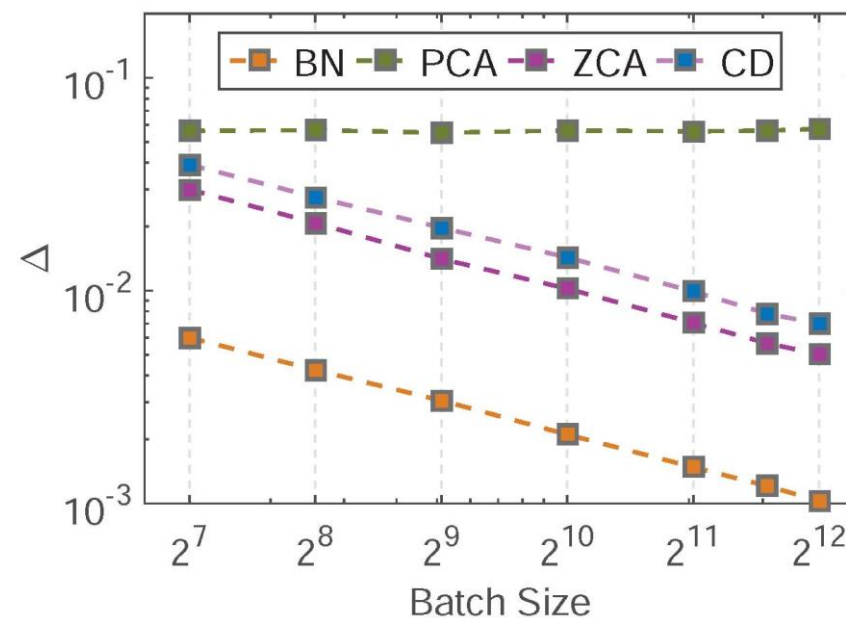
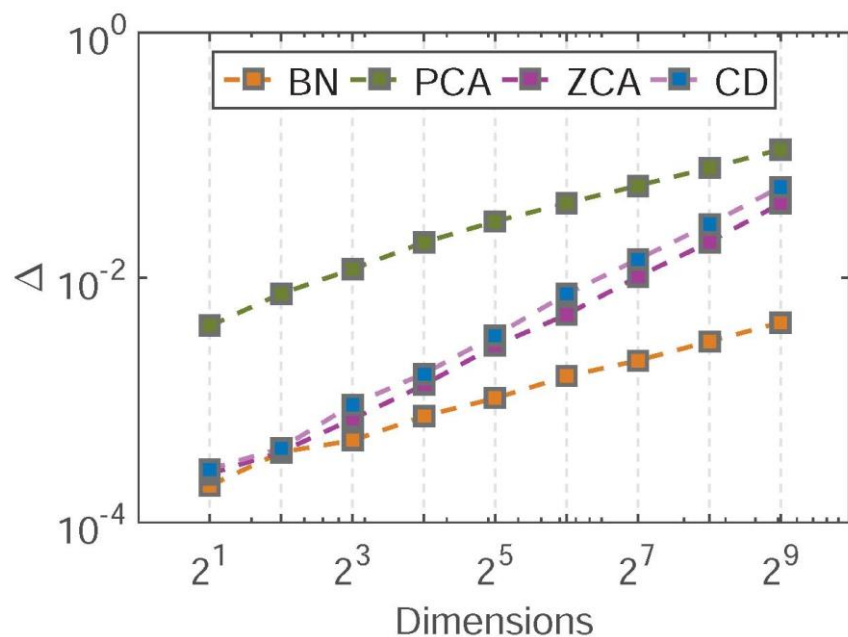


CD



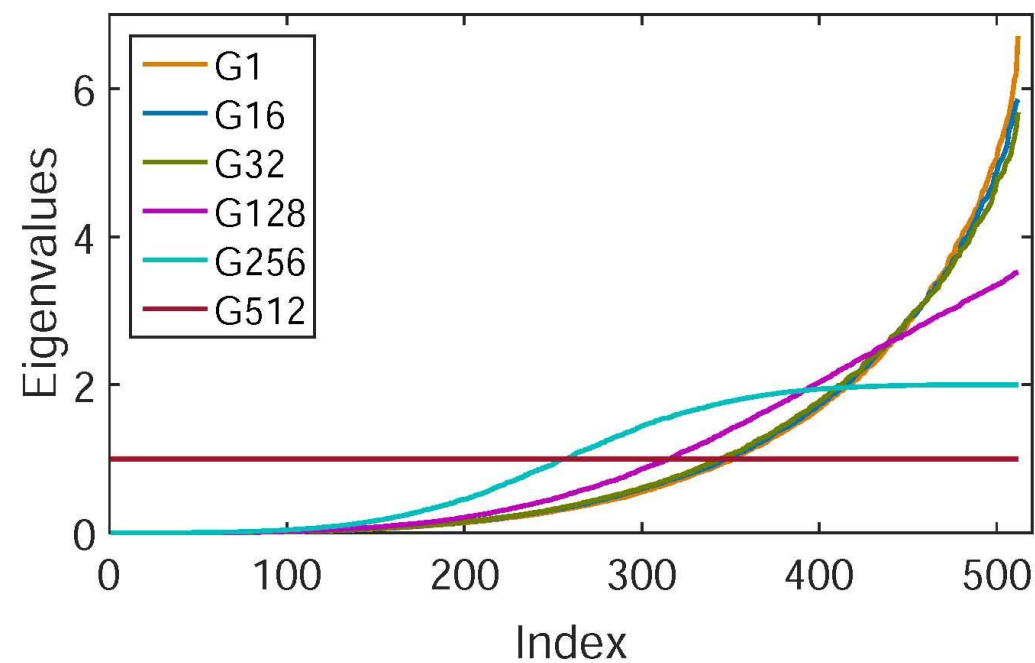
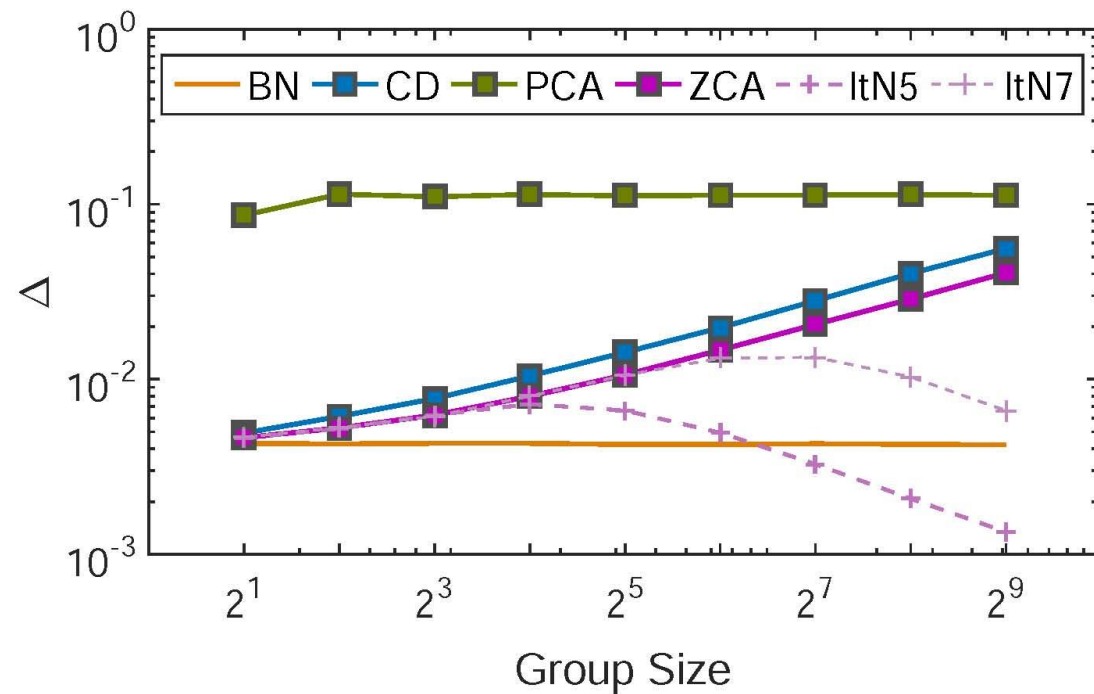
Investigation into the Stochasticity

- Stochasticity related to batch size and dimension
 - Increased SND with increasing dimension and decreasing batch size
 - SND: $PCA > CD > ZCA > BN$
 - PCA has nearly same large SND among batches



Investigation into the Stochasticity

- Controlling the stochasticity of BW by groups



Theoretical Model

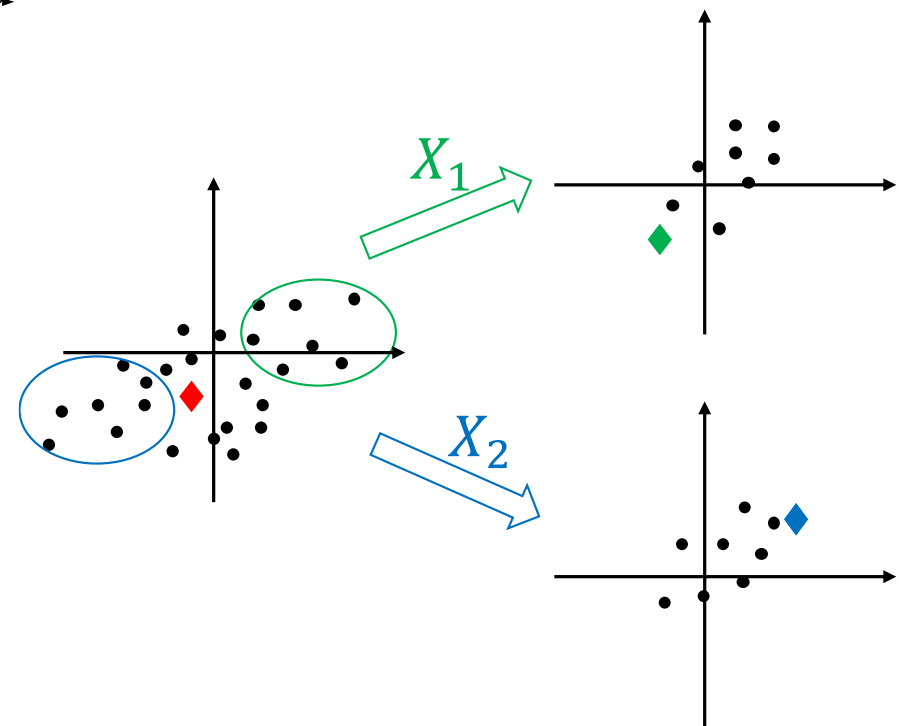
- Assumption:** X is approximately normal distribution with population statistics $(\hat{\mu}, \hat{\sigma})$, X_m is i.i.d sampled

Consider an example x ,
mini-batch output:

$$\frac{x - \mu_{X_m}}{\sigma_{X_m}} = \left(\frac{x - \hat{\mu}}{\hat{\sigma}} + \frac{\hat{\mu} - \mu_{X_m}}{\hat{\sigma}} \right) \frac{\hat{\sigma}}{\sigma_{X_m}}$$

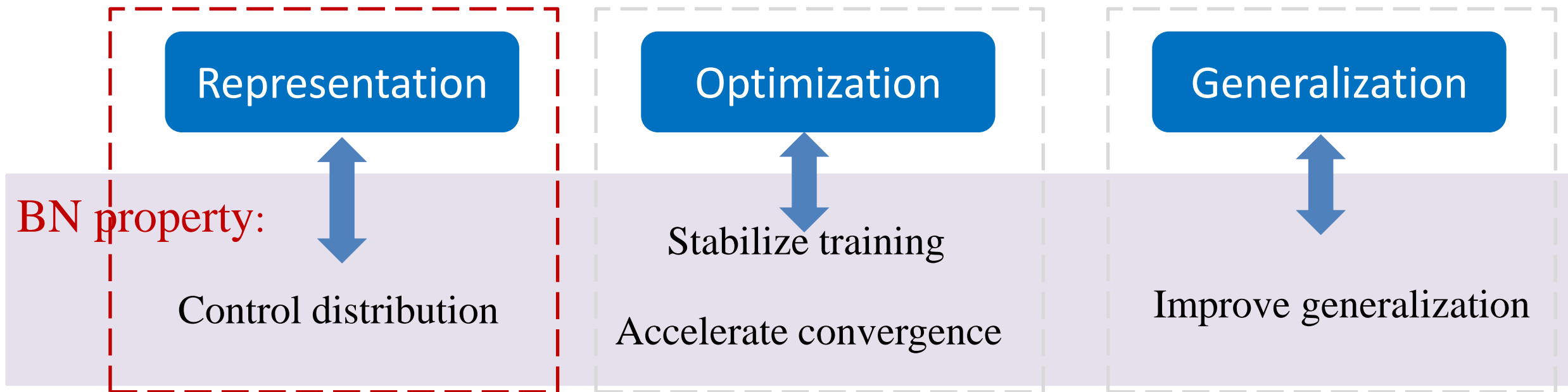
$$\frac{\hat{\mu} - \mu_{X_m}}{\hat{\sigma}} \sim \frac{1}{\sqrt{m}} \mathcal{N}(0,1)$$

$$\frac{\hat{\sigma}}{\sigma_{X_m}} \sim \sqrt{m} \chi_{m-1}^{-1}$$



Analyses of Normalization

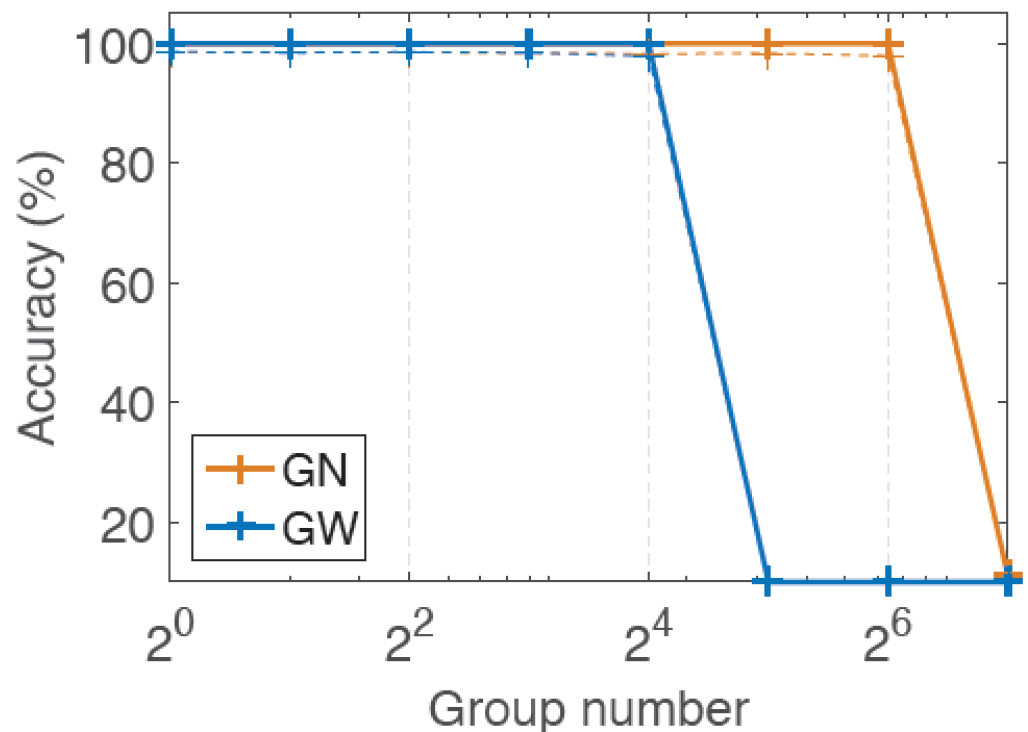
- Topics of deep learning theory



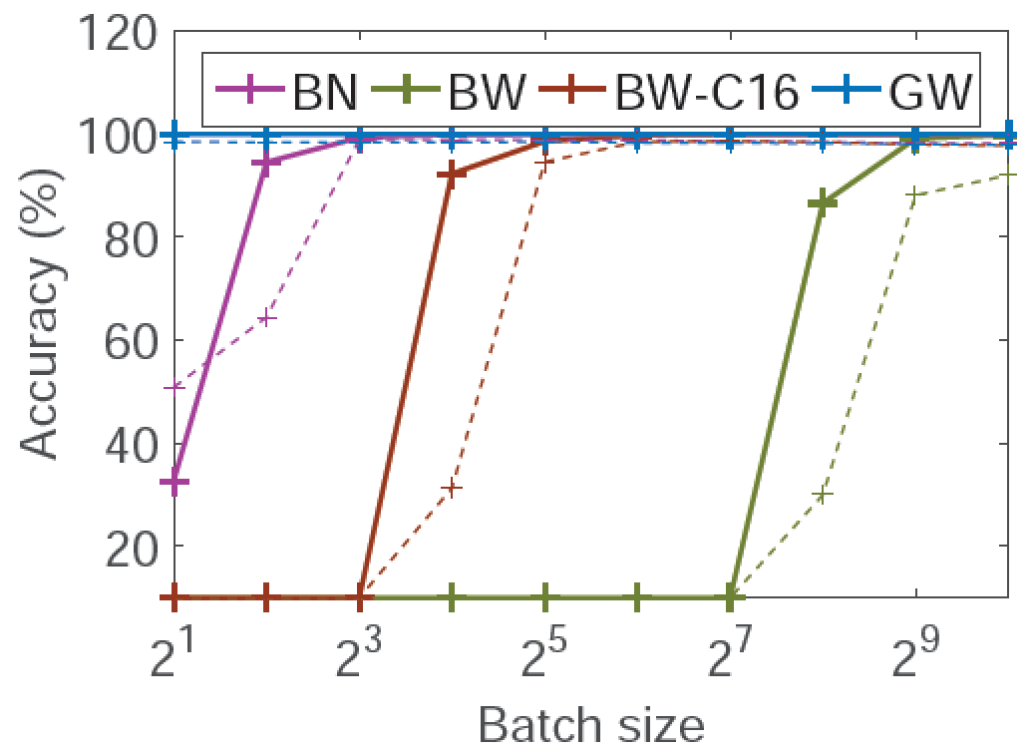
- Disclaims in this tutorials:
 - Not theoretic analysis in generalization bound or optimization theory
 - But empirical analysis with high-level theoretic insights.

Revisiting the Constraint of Normalization

- The large group problem of GN

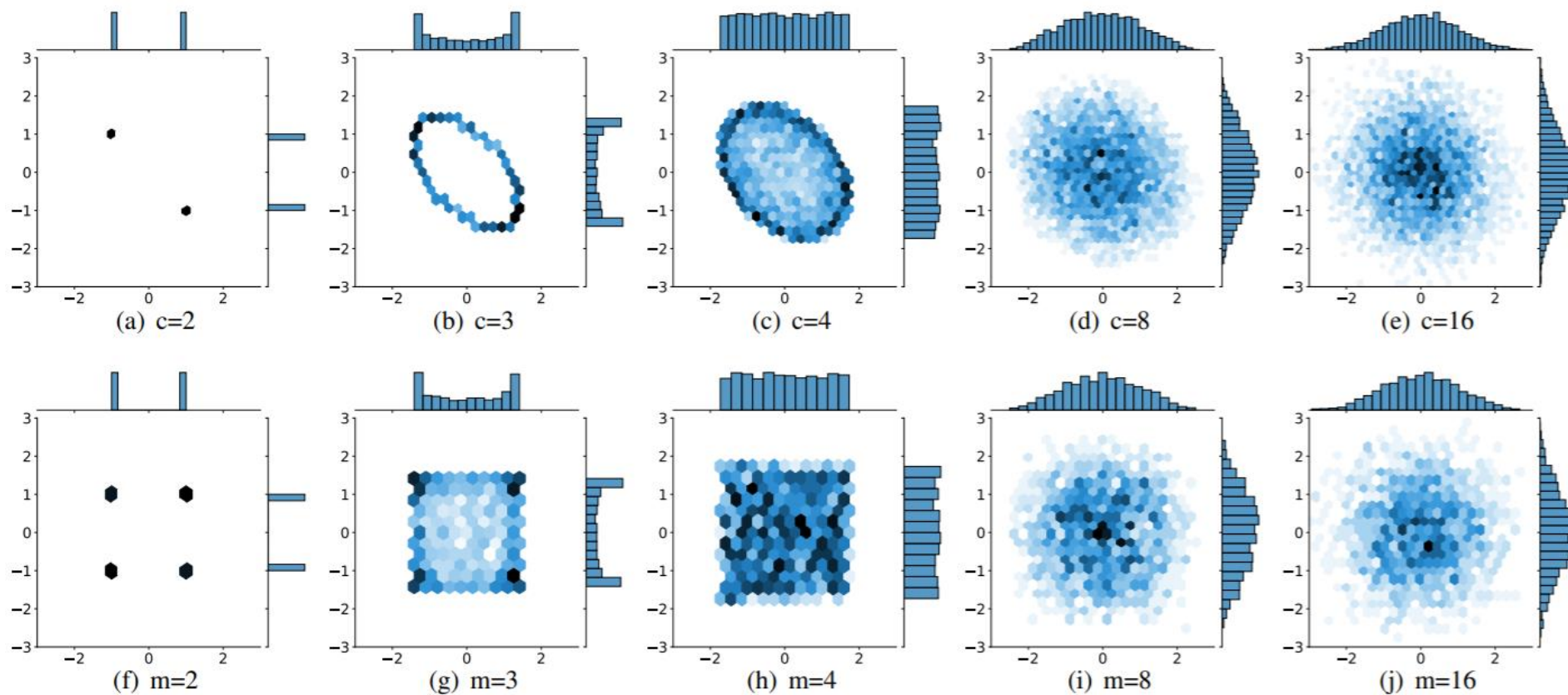


- The small-batch problem of BN/BW



Revisiting the Constraint of Normalization

- Feature representation in a layer



Revisiting the Constraint of Normalization

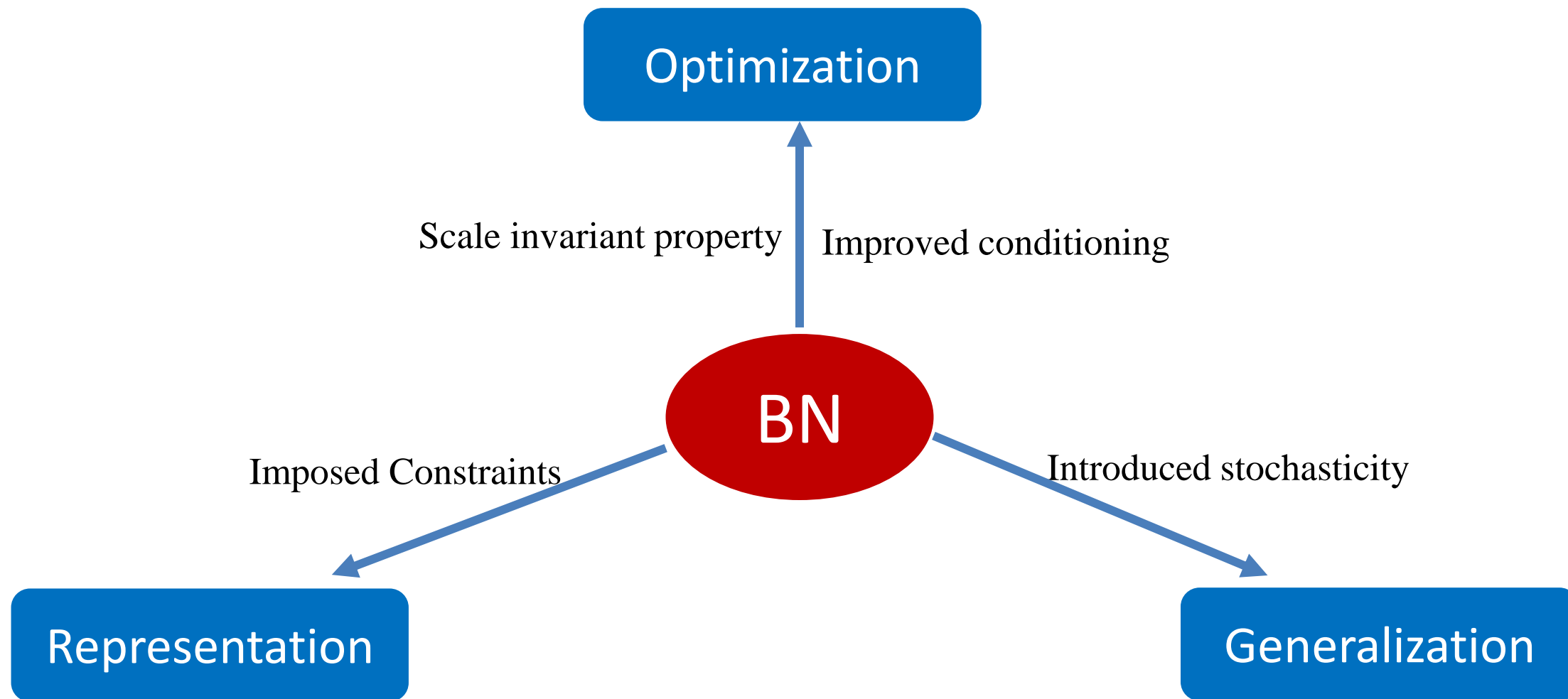
- Constraints as a system of equations: $\Upsilon_{\phi}(\hat{X})$

$$\Upsilon_{\phi_{BN}}(\hat{X}): \sum_{j=1}^m \hat{X}_{ij} = 0 \text{ and } \sum_{j=1}^m \hat{X}_{ij}^2 - m = 0$$

- Constrain number $\zeta(\phi; X)$: the number of equations in $\Upsilon_{\phi}(\hat{X})$

	Normalization along a batch		Normalization along a group of neurons	
	BN	BW	GN	GW
$\zeta(\phi; \mathbf{X})$	$2d$	$\frac{d(d+3)}{2}$	$2gm$	$\frac{mg(g+3)}{2}$
$\zeta(\phi; \mathbb{D})$	$\frac{2Nd}{m}$	$\frac{Nd(d+3)}{2m}$	$2gN$	$\frac{Ng(g+3)}{2}$
Ranges of m/g	$m \geq 2$	$m \geq \frac{d+3}{2}$	$g \leq \frac{d}{2}$	$g \leq \frac{\sqrt{8d+9}-3}{2}$

Summary





Challenges and Opportunities

- Theoretic analysis of DNN itself is difficult
- Normalizing activations make theoretic analysis of DNN challenging
 - Ruin the assumption of representation analysis tools
 - Data dependent
- Normalizing weight attribute to theoretic analysis of DNN
 - Pros: Lipschitz constant can be controlled/bounded during training by normalizing the weight
 - Certificated defense against adversarial attacks [Tsuzuku et al, NeurIPS 2018]
 - Theoretically analyzing DNNs generalization [Bartlett et al, NeurIPS 2017]]
 - Challenges: normalizing weights is still not as effective as normalizing activations in practice [Huang et al, CVPR 2020; Brock et al, ICLR 2021]