

2022 WINTER ESC

Deep Learning Problem Solving Technique



고정민, 박상용



Week 3

CONTENTS

1. No Free Lunch Theorem

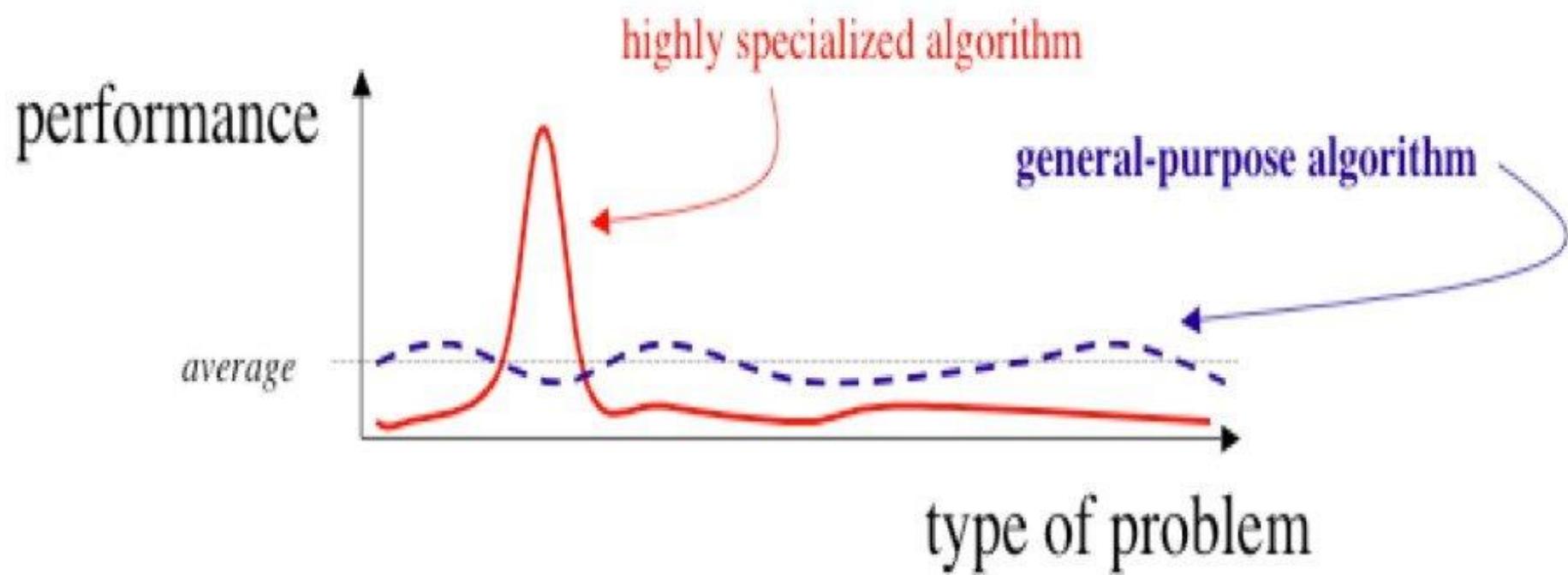
2. Weight Decay

3. Early Stopping

4. Ensemble method + Drop Out

5. Batch Normalization

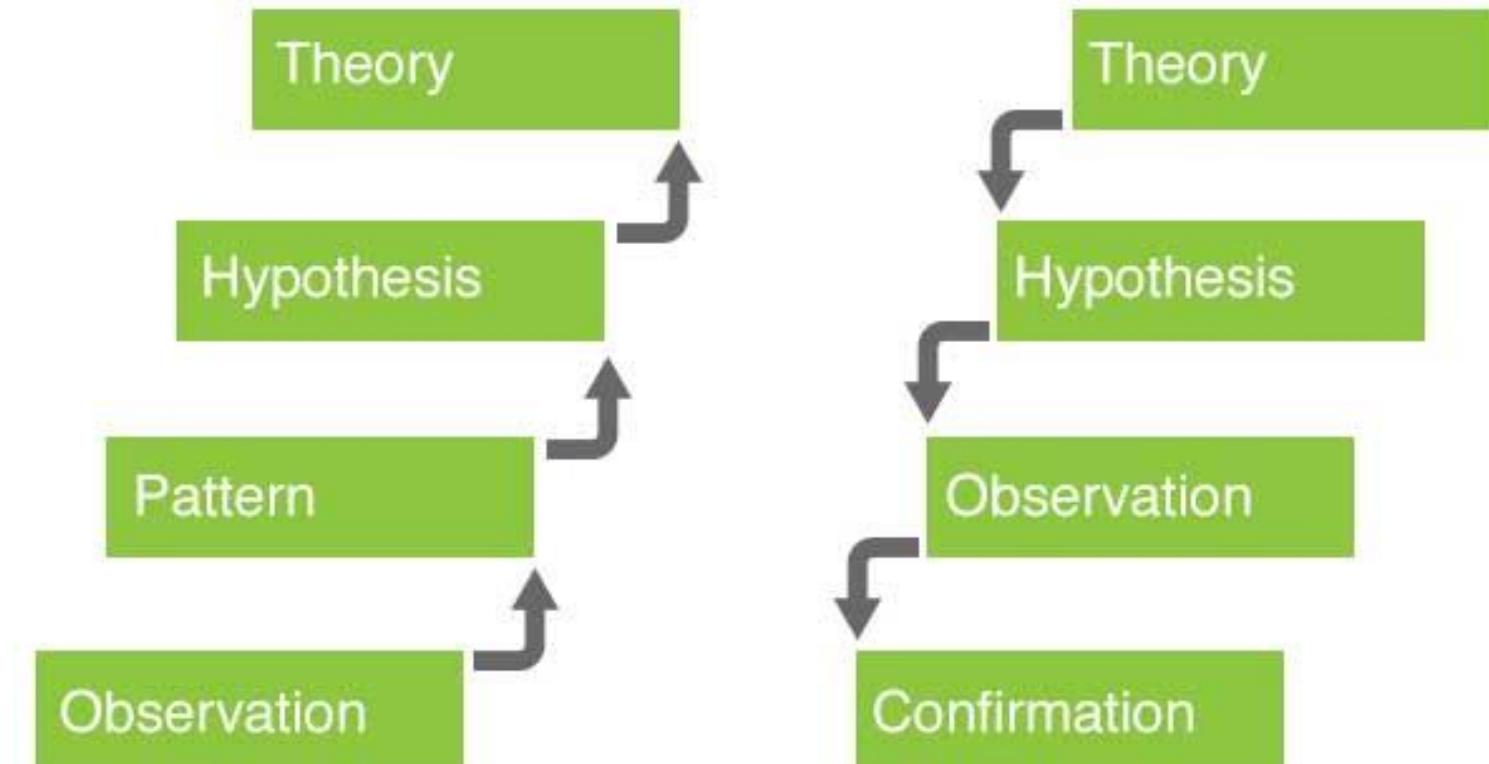
No Free Lunch Theorem for Optimization



We have dubbed the associated results "No Free Lunch" theorems because they demonstrate that if an algorithm performs well on a certain class of problems then it necessarily pays for that with degraded performance on the set of all remaining problems.

No Free Lunch Theorem for Optimization

Inductive Reasoning vs Deductive Reasoning



No Free Lunch Theorem for Optimization

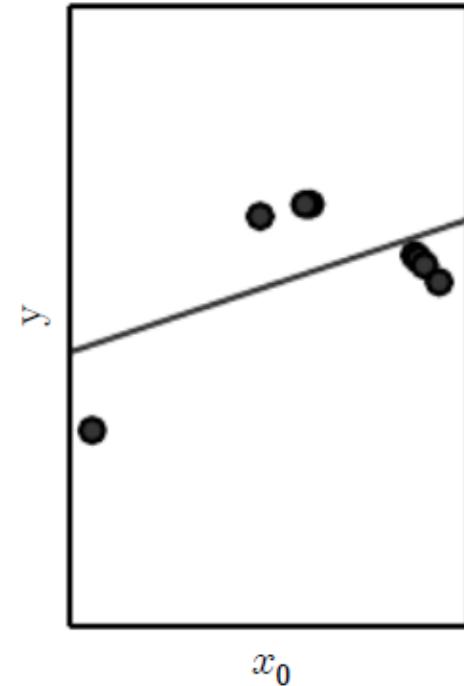
The factors determining how well a machine learning algorithm will perform

1. Make the training error small.
2. Make the gap between training and test error small.

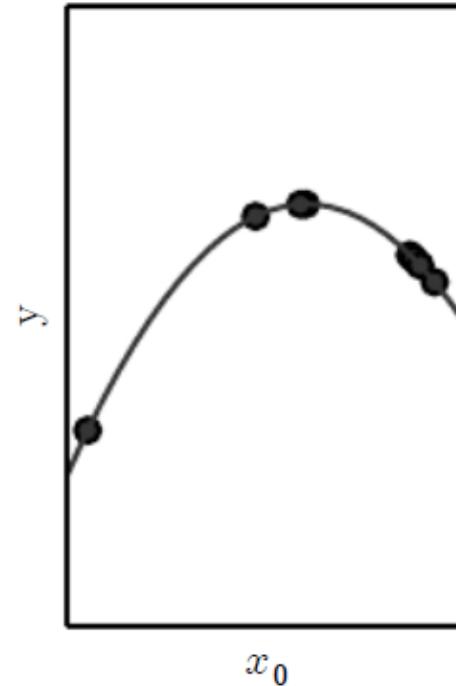
No Free Lunch Theorem for Optimization

Need to control the capacity of a learning algorithm

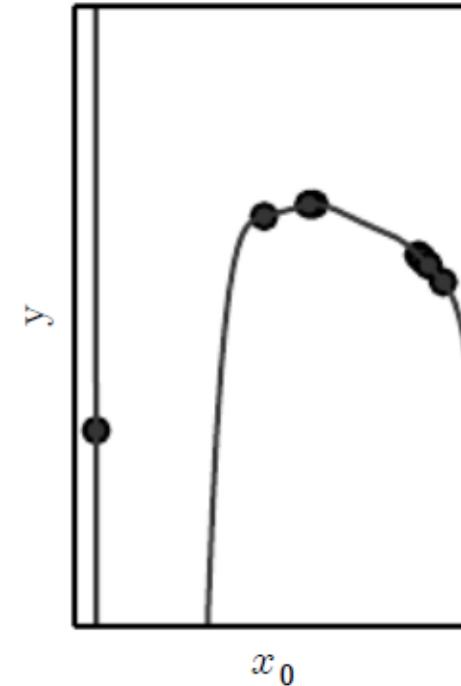
Underfitting



Appropriate capacity



Overfitting



No Free Lunch Theorem for Optimization

Bias - Variance Trade-Off

$\hat{f}(x_0) = \hat{f}$ evaluated at $X = x_0$

$$\mathbb{E} [\{\hat{f}(x_0) - y_0\}^2 \mid X = x_0]$$

$$= \mathbb{E} [\{\hat{f}(x_0) - (f(x_0) + \varepsilon)\}^2 \mid X = x_0] \quad (y = f(x) + \varepsilon, \mathbb{E}(\varepsilon) = 0, \text{Var}(\varepsilon) = \sigma_\varepsilon^2)$$

$$= \mathbb{E} [\{\hat{f}(x_0) - (f(x_0) + \varepsilon) + \mathbb{E}(\hat{f}(x_0)) - \mathbb{E}(\hat{f}(x_0))\}^2 \mid X = x_0]$$

No Free Lunch Theorem for Optimization

Bias - Variance Trade-Off

$$= \mathbb{E} \left[\left\{ \varepsilon - [\mathbb{E}(\hat{f}(x_0)) - f(x_0)] - [\hat{f}(x_0) - \mathbb{E}(\hat{f}(x_0))] \right\}^2 \mid X = x_0 \right]$$

$$= \mathbb{E}(\varepsilon^2) + \mathbb{E} \left[\left\{ \mathbb{E}(\hat{f}(x_0)) - f(x_0) \right\}^2 \right] + \mathbb{E} \left[\left\{ \hat{f}(x_0) - \mathbb{E}(\hat{f}(x_0)) \right\}^2 \mid X = x_0 \right]$$

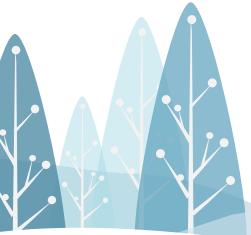
$$= \sigma_\varepsilon^2 + \left\{ \mathbb{E}(\hat{f}(x_0)) - f(x_0) \right\}^2 + Var(\hat{f}(x_0))$$

$$= \sigma_\varepsilon^2 + Bias^2(\hat{f}(x_0)) + Var(\hat{f}(x_0))$$

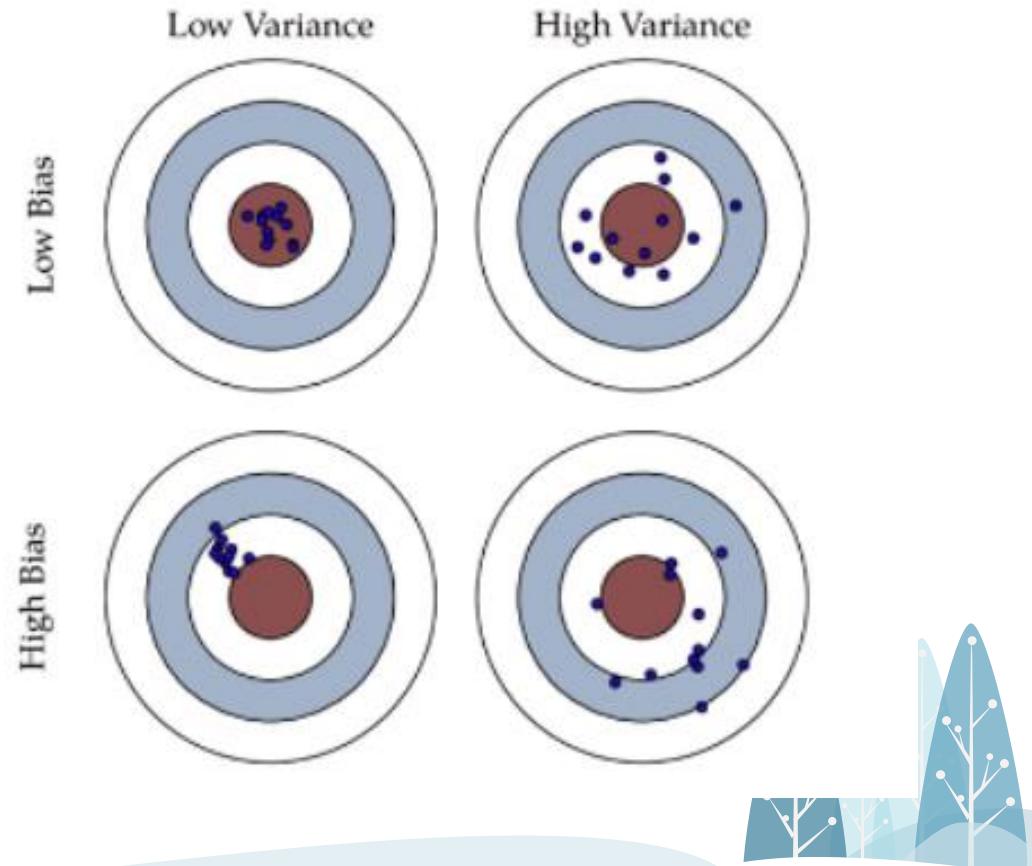
No Free Lunch Theorem for Optimization

Bias - Variance Trade-Off

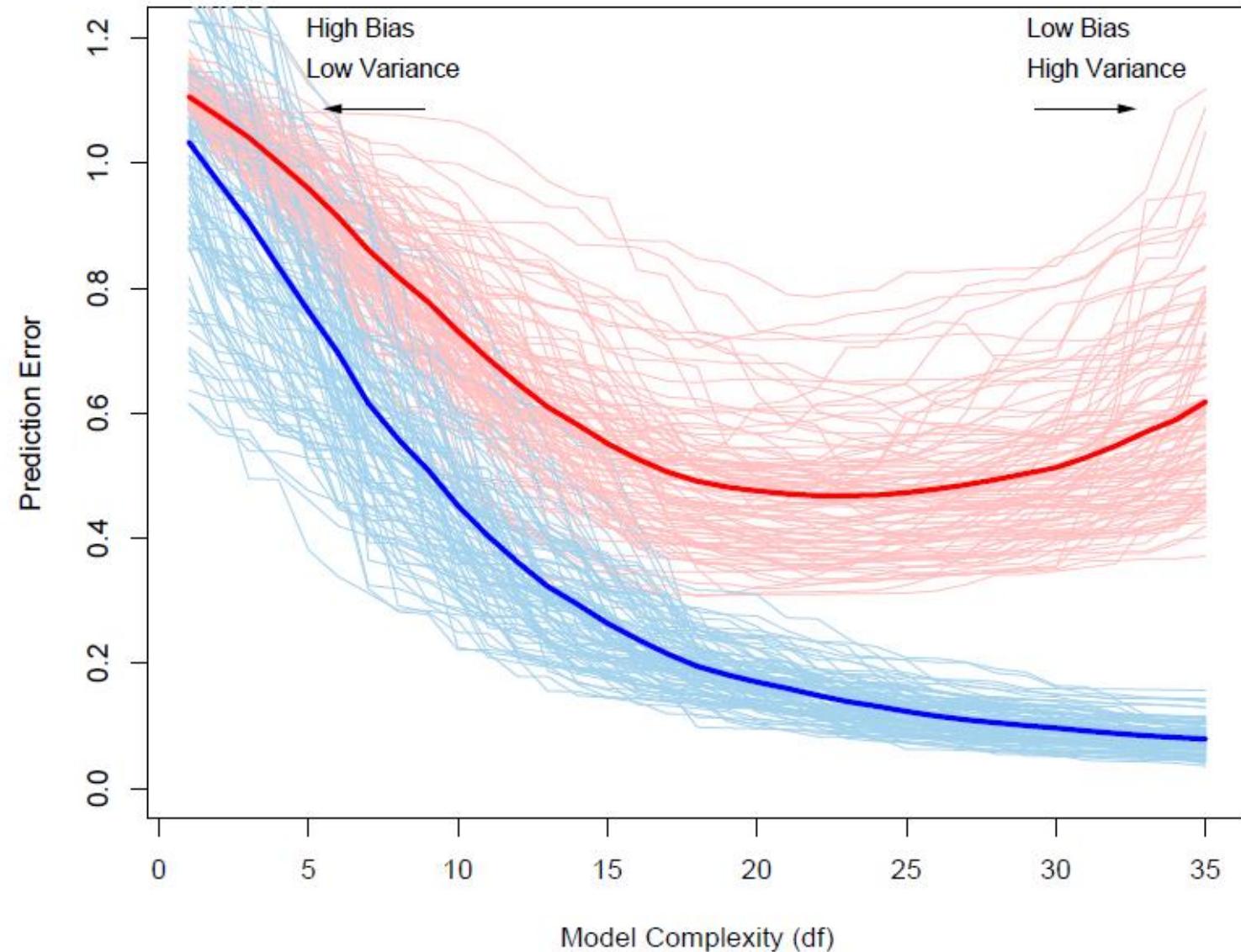
$$\sigma_{\varepsilon}^2 + \text{Bias}^2 (\hat{f}(x_0)) + \text{Var}(\hat{f}(x_0))$$



Irreducible Error + Bias² + Variance



No Free Lunch Theorem for Optimization



Regularization - Weight Decay



너무 좀 말이 많으세요. 말이

출처 : <https://youtu.be/edMByVIO8MM>

Regularization - Weight Decay

$$f(\mathbf{x} ; \boldsymbol{\theta}) \leftarrow \Omega(\mathbf{w}) = \mathbf{w}^\top \mathbf{w} \quad (\text{Regularizer})$$

$$J(\mathbf{w}) = MSE_{train} + \lambda \mathbf{w}^\top \mathbf{w}$$

Minimizing $J(w)$ results in a choice of weights that make a trade-off between fitting the training data and being small.

Regularization - Weight Decay

Norm

$$\|x\|_p = \left(\sum_i |x_i|^p \right)^{\frac{1}{p}}, \quad p \in \mathbb{R}, \quad p \geq 1$$

Regularization - Weight Decay

L^2 Norm (Euclidean Norm)

$$\|x\|_2 = \left(\sum_i |x_i|^2 \right)^{\frac{1}{2}} \rightarrow \|x\|_2^2 = \sum_i |x_i|^2$$

Squared L^2 Norm

L^1 Norm

$$\|x\|_1 = \sum_i |x_i|$$

Regularization - Weight Decay

Parameter Norm Penalties

$$\tilde{J}(\theta; X, y) = J(\theta; X, y) + \alpha \Omega(\theta)$$

Different choices for the parameter norm Ω can result in different solutions being preferred.

Regularization - Weight Decay

L^2 Parameter Regularization
(a.k.a. Ridge Regression / Tikhonov Regularization)

$$\Omega(\theta) = \frac{1}{2} \|w\|_2^2 = \frac{1}{2} w^\top w$$

$$\tilde{J}(w; X, y) = \frac{\alpha}{2} w^\top w + J(w; X, y)$$

Regularization - Weight Decay

L^2 Parameter Regularization

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \frac{\alpha}{2} \mathbf{w}^\top \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

with the corresponding parameter gradient,

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

Regularization - Weight Decay

L^2 Parameter Regularization

To take a single gradient step to update the weights, we perform this update:

$$\mathbf{w} \leftarrow \mathbf{w} - \epsilon (\alpha \mathbf{w} + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y}))$$

Written another way, the update is:

$$\mathbf{w} \leftarrow (1 - \epsilon \alpha) \mathbf{w} - \epsilon \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

Regularization - Weight Decay

L^1 Parameter Regularization (a.k.a. Lasso Regression)

$$\Omega(\boldsymbol{\theta}) = \|\boldsymbol{w}\|_1 = \sum_i |w_i|$$

$$\tilde{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = \alpha \|\boldsymbol{w}\|_1 + J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y})$$

Regularization - Weight Decay

L^1 Parameter Regularization

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \|\mathbf{w}\|_1 + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

with the corresponding gradient (actually, **sub-gradient**),

$$\nabla_{\mathbf{w}} \tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \cdot sign(\mathbf{w}) + \nabla_{\mathbf{w}} J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

Regularization - Weight Decay

Non-differentiable functions and subgradients

Given a convex function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,
a vector $z \in \mathbb{R}^n$ is said to be a *subgradient* of f at x_0 if

$$f(x) \geq f(x_0) + \langle z, x - x_0 \rangle \quad \text{for all } x \in \mathbb{R}^n$$

The set of all subgradients of f at x_0 is called the *subdifferential* and denoted by $\partial f(x_0)$.

When f is differentiable at x_0 , then the sub-differential reduces to a single vector - that is $\partial f(x_0) = \{\nabla f(x_0)\}$.

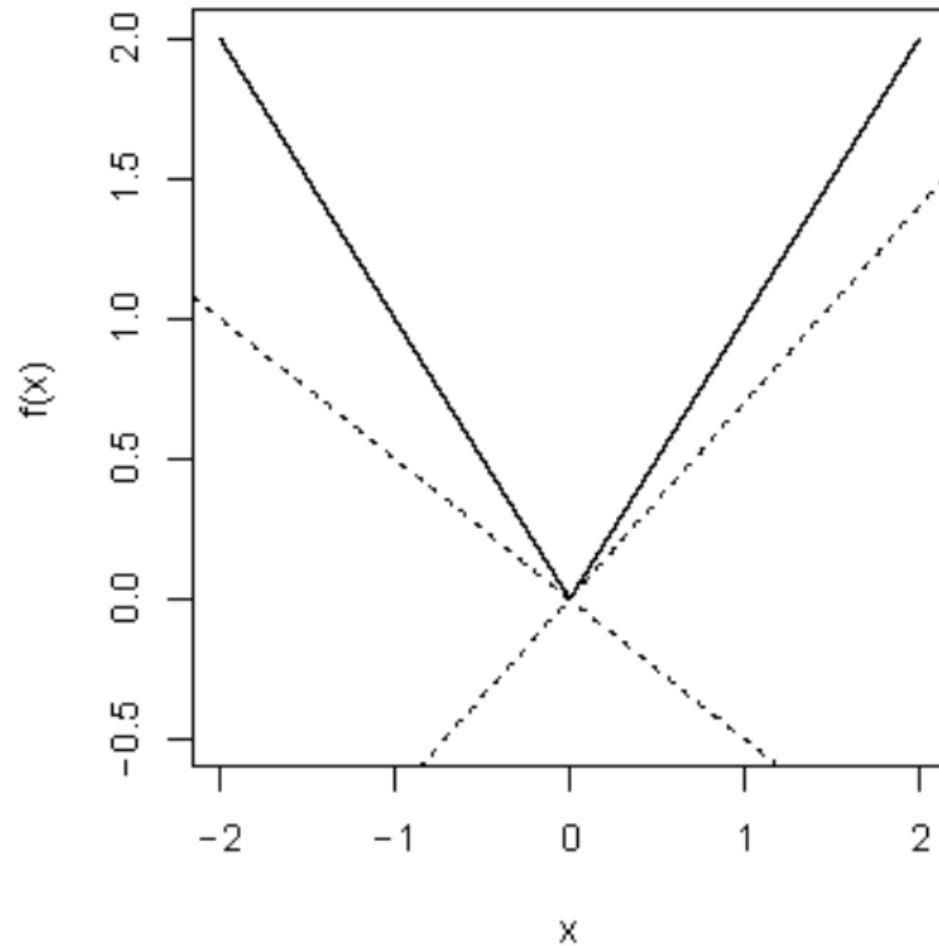
Regularization - Weight Decay

Non-differentiable functions and subgradients

Example 1

$$f : \mathbb{R} \rightarrow \mathbb{R},$$

$$f(x) = \|x\|_1 = |x|$$



Regularization - Weight Decay

Non-differentiable functions and subgradients

Example 1

$$f : \mathbb{R} \rightarrow \mathbb{R}, \quad f(x) = \|x\|_1 = |x|$$

$$\partial f(x) = \begin{cases} \{+1\}, & \text{if } x > 0 \\ \{-1\}, & \text{if } x < 0 \\ [-1, 1], & \text{if } x = 0 \end{cases}$$

Which is frequently denoted by $\text{sign}(x)$. $z \in \text{sign}(x) : z$ belongs to the sub-differential of the absolute value function at x .

Regularization - Weight Decay

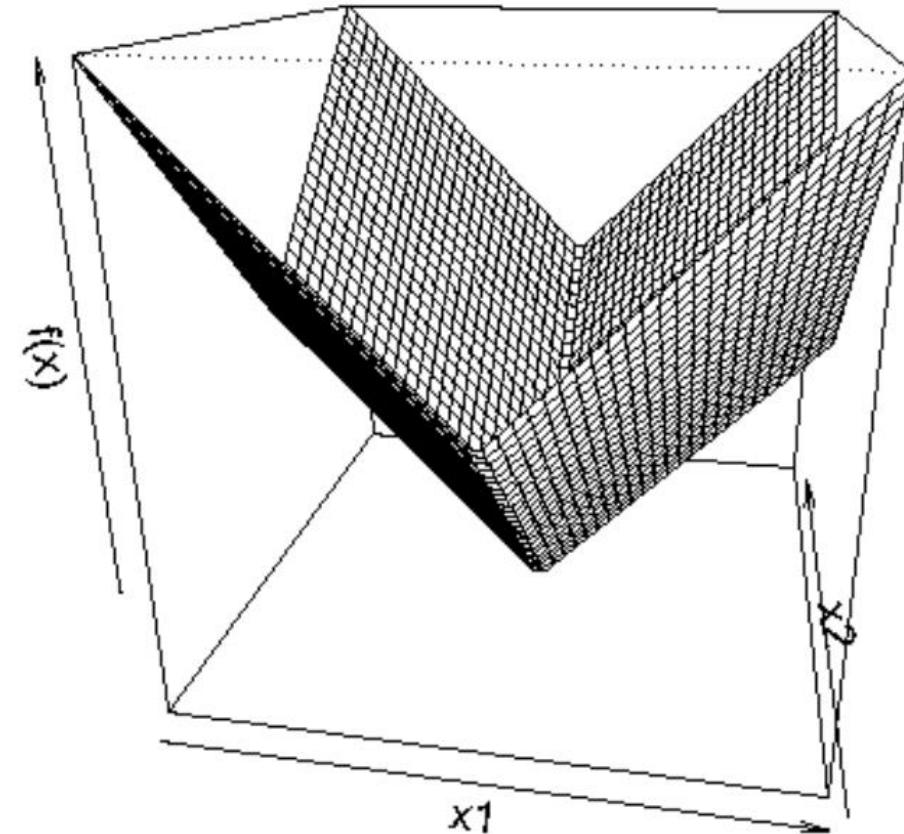
Non-differentiable functions and subgradients

Example 2

$$f : \mathbb{R}^n \rightarrow \mathbb{R},$$

$$f(x) = \|x\|_1 = |x|$$

$$x = (x_1, x_2, x_3, \dots, x_n)$$



Regularization - Weight Decay

Non-differentiable functions and subgradients

Example 2

$$f : \mathbb{R}^n \rightarrow \mathbb{R}, \quad f(x) = \|x\|_1 = |x|$$

$$\partial f(x_i) = \begin{cases} \{+1\}, & \text{if } x_i > 0 \\ \{-1\}, & \text{if } x_i < 0 \quad \text{for } i \in \{1, 2, \dots, n\} \\ [-1, 1], & \text{if } x_i = 0 \end{cases}$$

Which is frequently denoted by $\text{sign}(x_i)$. $z \in \text{sign}(x_i) : z$ belongs to the sub-differential of the absolute value function at x_i .

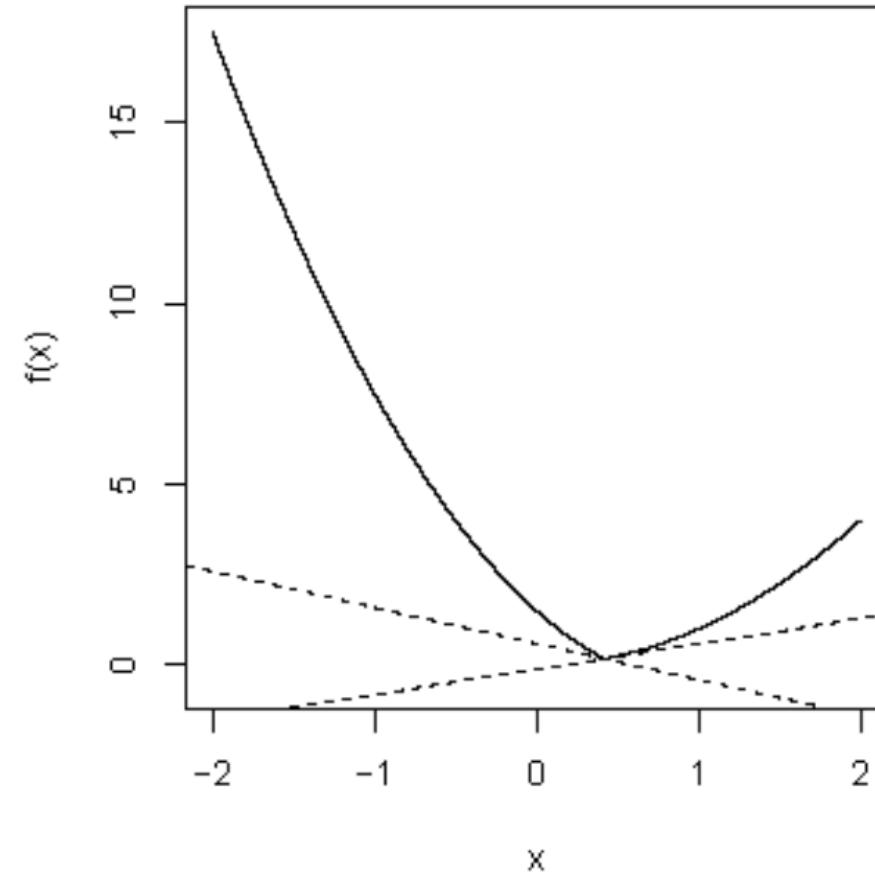
Regularization - Weight Decay

Non-differentiable functions and subgradients

Example 3

$$f(x) = \max\{f_1(x_1), f_2(x_2)\}$$

$$f_1, f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$$



Regularization - Weight Decay

Non-differentiable functions and subgradients

Example 3

$$f(x) = \max\{f_1(x), f_2(x)\}, \quad f_1, f_2 : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$f_1(x) > f_2(x) \rightarrow \partial f(x) = \nabla f_1(x)$$

$$f_1(x) < f_2(x) \rightarrow \partial f(x) = \nabla f_2(x)$$

$$f_1(x) = f_2(x) \rightarrow \partial f(x) \in \{\theta_1 \nabla f_1(x) + \theta_2 \nabla f_2(x) : \theta_1 + \theta_2 = 1, \theta_1 \geq 0, \theta_2 \geq 0\}$$

Regularization - Weight Decay

Elastic Net

$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = r\alpha \|\mathbf{w}\|_1 + \frac{1-r}{2} \alpha \mathbf{w}^\top \mathbf{w} + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

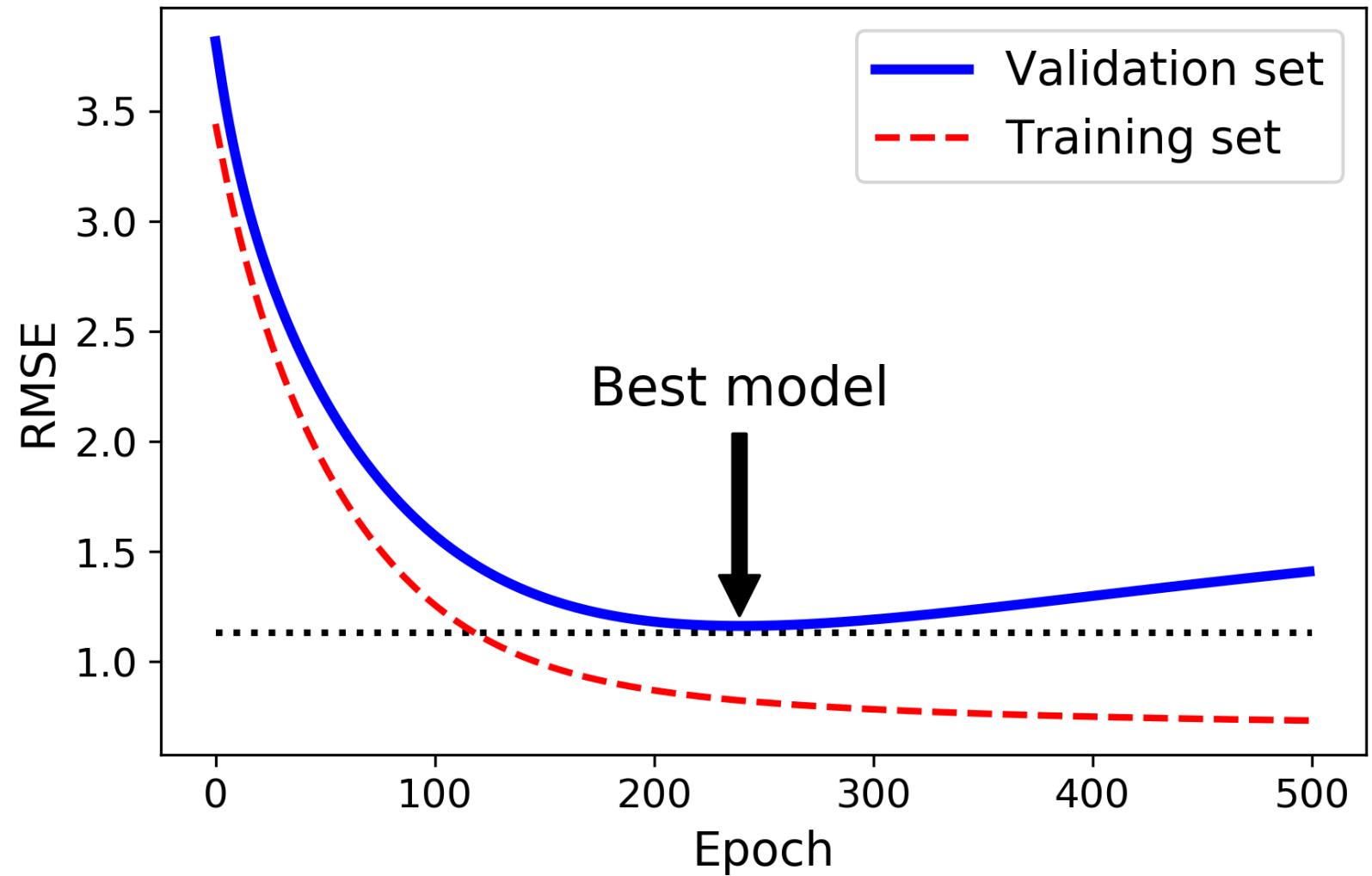
$r = 0 \rightarrow L^2$ regularization

$r = 1 \rightarrow L^1$ regularization

Regularization - Weight Decay

Ridge → Lasso → Elastic Net

Regularization - Early Stopping



Regularization - Early Stopping

Algorithm 7.1 The early stopping meta-algorithm for determining the best amount of time to train. This meta-algorithm is a general strategy that works well with a variety of training algorithms and ways of quantifying error on the validation set.

Let n be the number of steps between evaluations.

Let p be the “patience,” the number of times to observe worsening validation set error before giving up.

Let θ_o be the initial parameters.

Regularization - Early Stopping

Let n be the number of steps between evaluations.

Let p be the “patience,” the number of times to observe worsening validation set error before giving up.

Let θ_o be the initial parameters.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

while $j < p$ **do**

 Update θ by running the training algorithm for n steps.

$i \leftarrow i + n$

$v' \leftarrow \text{ValidationSetError}(\theta)$

if $v' < v$ **then**

$j \leftarrow 0$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

$v \leftarrow v'$

else

$j \leftarrow j + 1$

end if

end while

Best parameters are θ^* , best number of training steps is i^*

“Unobtrusive”



Regularization - Early Stopping

Algorithm 7.2 A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

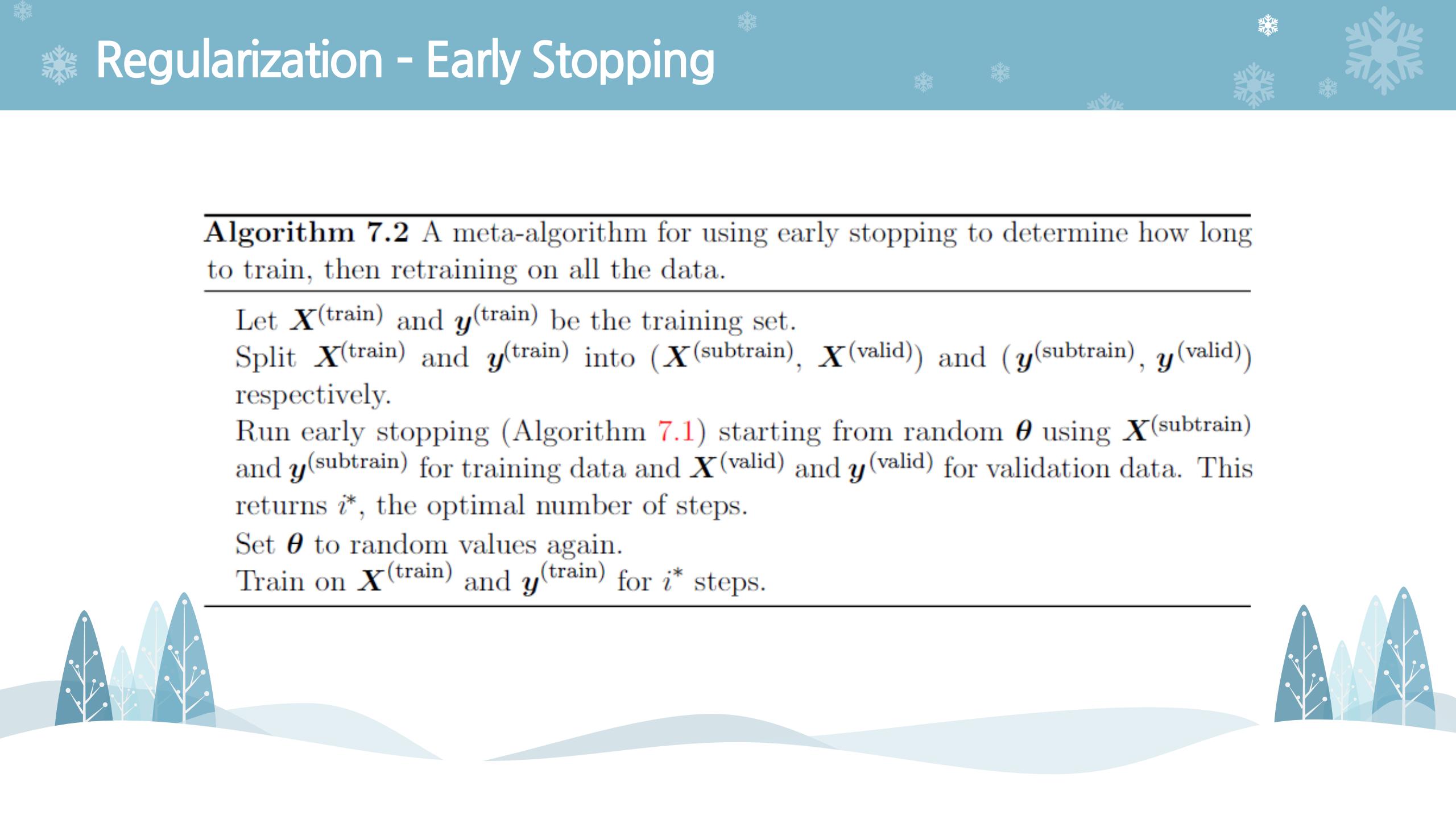
Let $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ be the training set.

Split $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ into $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$ and $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$ respectively.

Run early stopping (Algorithm 7.1) starting from random $\boldsymbol{\theta}$ using $\mathbf{X}^{(\text{subtrain})}$ and $\mathbf{y}^{(\text{subtrain})}$ for training data and $\mathbf{X}^{(\text{valid})}$ and $\mathbf{y}^{(\text{valid})}$ for validation data. This returns i^* , the optimal number of steps.

Set $\boldsymbol{\theta}$ to random values again.

Train on $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ for i^* steps.



Regularization - Early Stopping

Algorithm 7.2 A meta-algorithm for using early stopping to determine how long to train, then retraining on all the data.

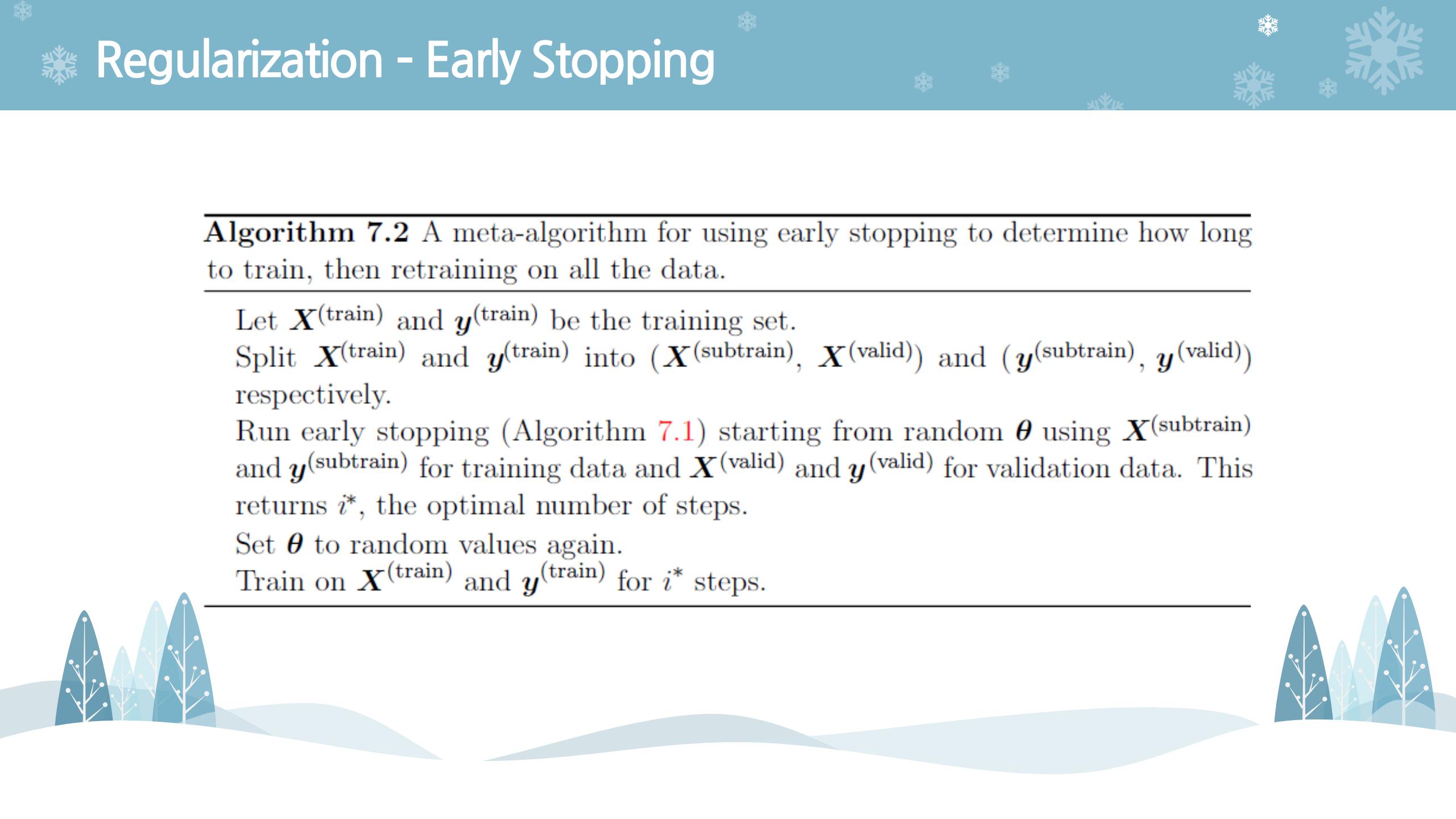
Let $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ be the training set.

Split $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ into $(\mathbf{X}^{(\text{subtrain})}, \mathbf{X}^{(\text{valid})})$ and $(\mathbf{y}^{(\text{subtrain})}, \mathbf{y}^{(\text{valid})})$ respectively.

Run early stopping (Algorithm 7.1) starting from random $\boldsymbol{\theta}$ using $\mathbf{X}^{(\text{subtrain})}$ and $\mathbf{y}^{(\text{subtrain})}$ for training data and $\mathbf{X}^{(\text{valid})}$ and $\mathbf{y}^{(\text{valid})}$ for validation data. This returns i^* , the optimal number of steps.

Set $\boldsymbol{\theta}$ to random values again.

Train on $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ for i^* steps.



Ensemble Method

Wisdom of Crowds



Ensemble Method

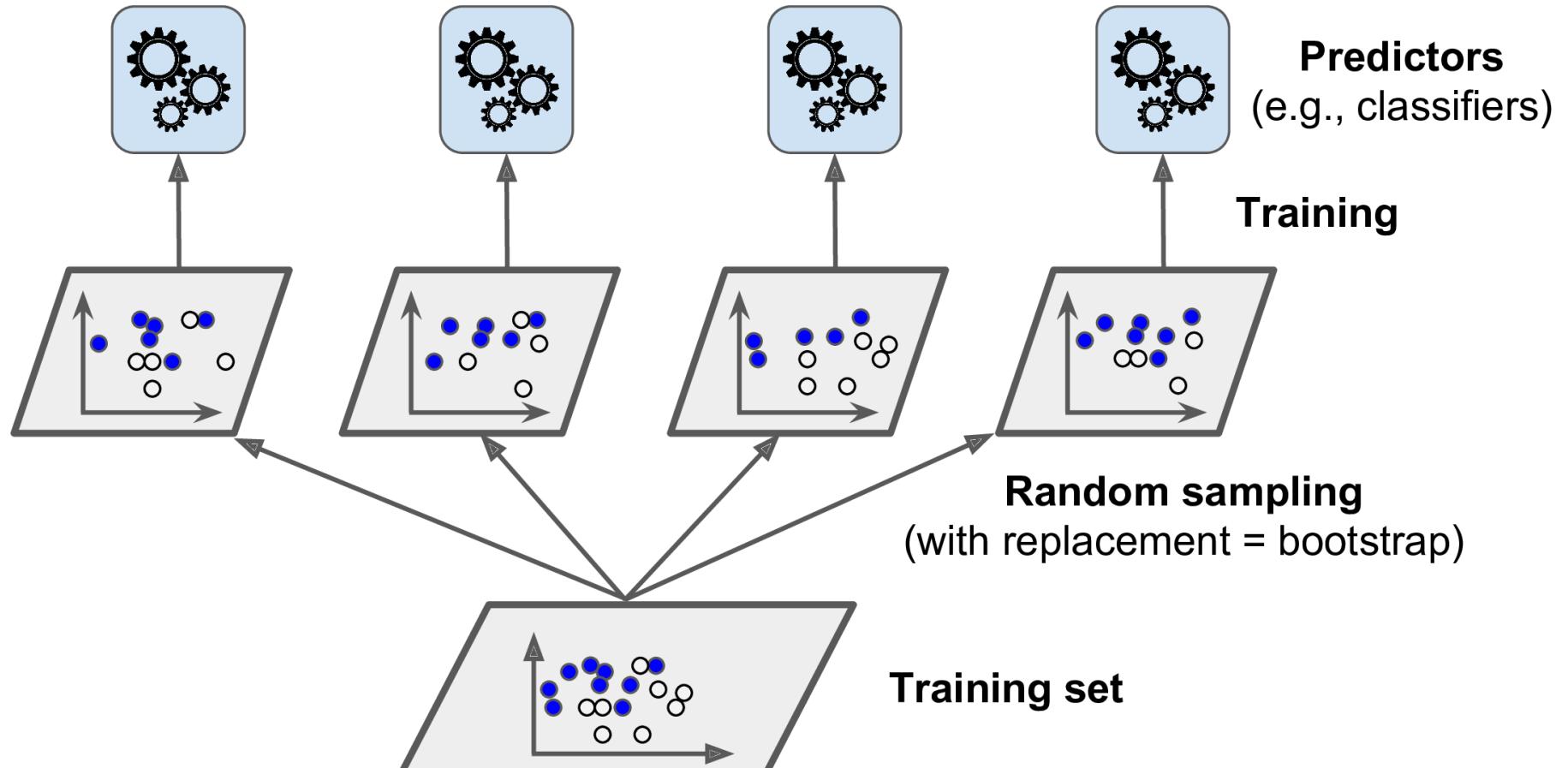
Ensemble Method - Model Averaging

The idea of ensemble learning is to build a prediction model by combining the strengths of a collection of simpler base models.

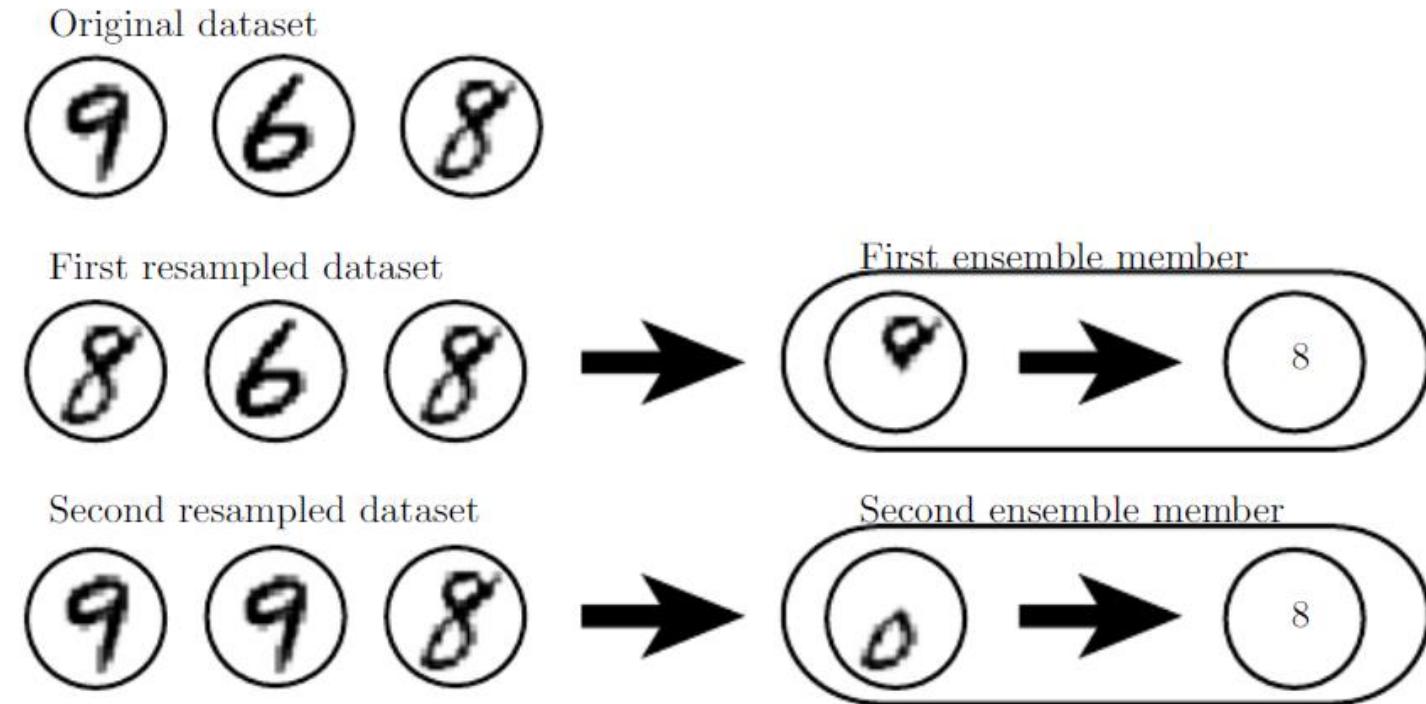
Bagging (Bootstrap aggregation)

The idea is to train several different models separately, then have all of the models vote on the output for test examples.

Ensemble Method - Bagging



Ensemble Method - Bagging



The first dataset omits the '9' and repeats the '8'. On this dataset, the detector learns that a loop on top of the digit corresponds to an '8'.

On the second dataset, we repeat the '9' and omit the '6'. In this case, the detector learns that a loop on the bottom of the digit corresponds to an '8'.

Ensemble Method - Bagging

Variance of the Average of Prediction

Consider for example a set of k regression models.

Suppose that each model makes an error ϵ_i on each example.

$$\epsilon_i \sim N(m, \sigma^2), \quad \text{Corr}(\epsilon_i, \epsilon_j) = \rho$$

$$\frac{1}{\sigma^2} \mathbb{E}[(\epsilon_i - m)(\epsilon_j - m)] = \rho > 0$$

$$\mathbb{E}[\epsilon_i \epsilon_j] = \rho \sigma^2 + m^2$$

Ensemble Method - Bagging

Variance of the Average of Prediction

$$\mathbb{E}[\epsilon_i] = m \quad \mathbb{E}[\epsilon_i \epsilon_j] = \begin{cases} \rho\sigma^2 + m^2 & i \neq j \\ \sigma^2 + m^2 & i = j \end{cases}$$

$$Var\left(\frac{1}{k}\sum_i \epsilon_i\right) = \frac{1}{k^2} Var\left(\sum_i \epsilon_i\right) = \frac{1}{k^2} \left[\mathbb{E}\left[\left(\sum_i \epsilon_i\right)^2\right] - \mathbb{E}\left[\sum_i \epsilon_i\right]^2 \right]$$

Ensemble Method - Bagging

Variance of the Average of Prediction

$$\mathbb{E} \left[\sum_i \epsilon_i \right] = \sum_i \mathbb{E}[\epsilon_i] = km$$

$$\left(\sum_i \epsilon_i \right)^2 = \sum_{i,j} \epsilon_i \epsilon_j \quad \mathbb{E} \left[\left(\sum_i \epsilon_i \right)^2 \right] = \sum_{i,j} \mathbb{E}[\epsilon_i \epsilon_j] = k\mathbb{E}[\epsilon_i^2] + (k^2 - k)\mathbb{E}[\epsilon_i \epsilon_j]$$
$$= k(\sigma^2 + m^2) + k(k-1)(\rho\sigma^2 + m^2)$$

Ensemble Method - Bagging

Variance of the Average of Prediction

$$k(\sigma^2 + m^2) + k(k-1)(\rho\sigma^2 + m^2)$$

$$= k\sigma^2 + k^2\rho\sigma^2 - k\rho\sigma^2 + k^2m^2$$

$$Var\left(\frac{1}{k}\sum_i \epsilon_i\right) = \frac{1}{k^2}(k\sigma^2 + k^2\rho\sigma^2 - k\rho\sigma^2) = \frac{1}{k}(\sigma^2 + k\rho\sigma^2 - \rho\sigma^2)$$

$$= \rho\sigma^2 + \frac{1-\rho}{k}\sigma^2$$



Break Time

Session Flow



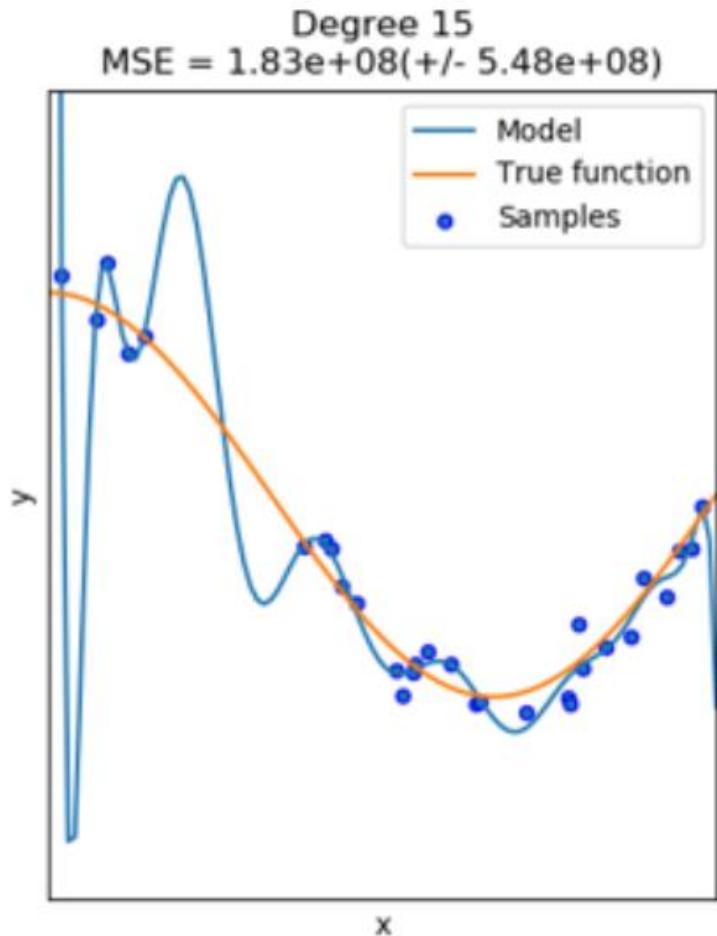
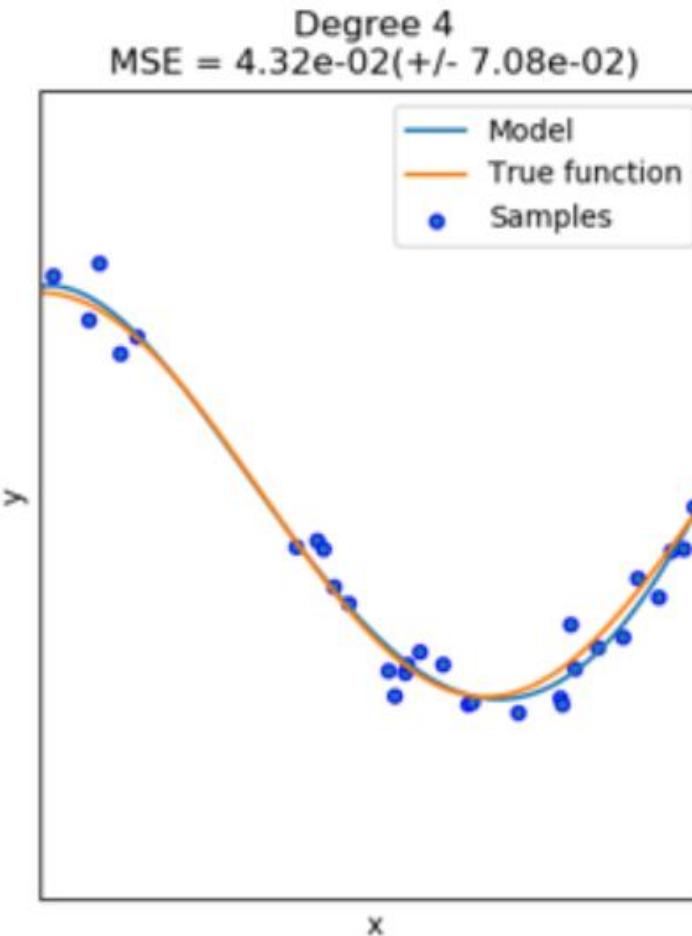
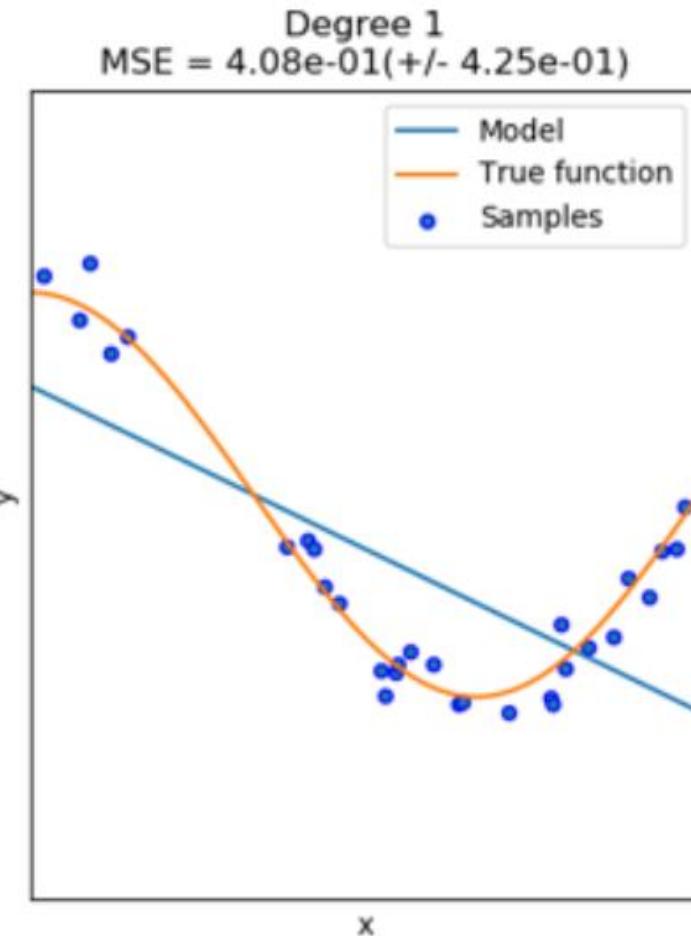
Curse of
Dimensionality

Dropout &
Drop-Connect

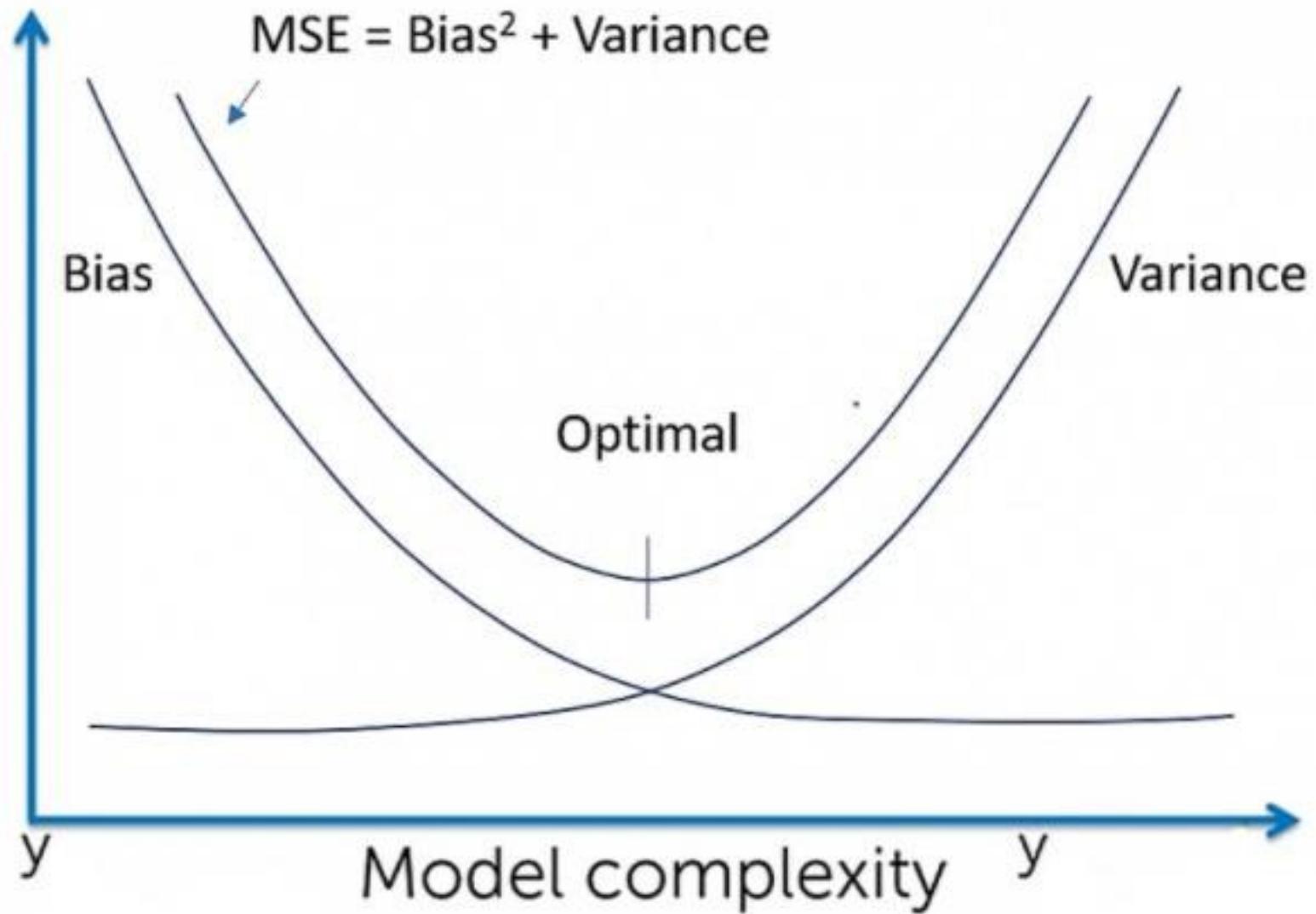
Batch-Normalization

Q & A

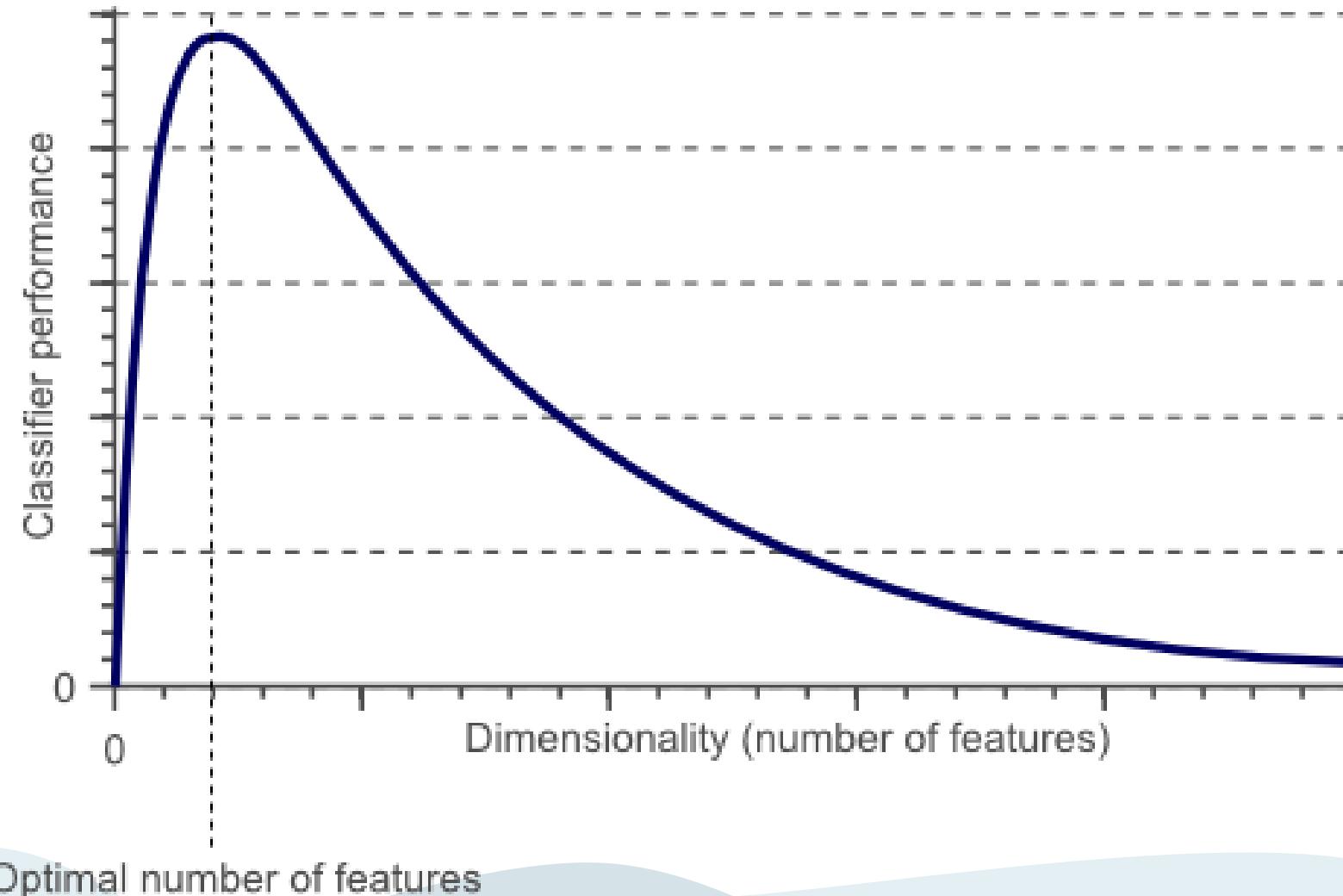
1. Over-fitting



Bias-Variance Trade-off



Curse of Dimensionality



Curse of Dimensionality

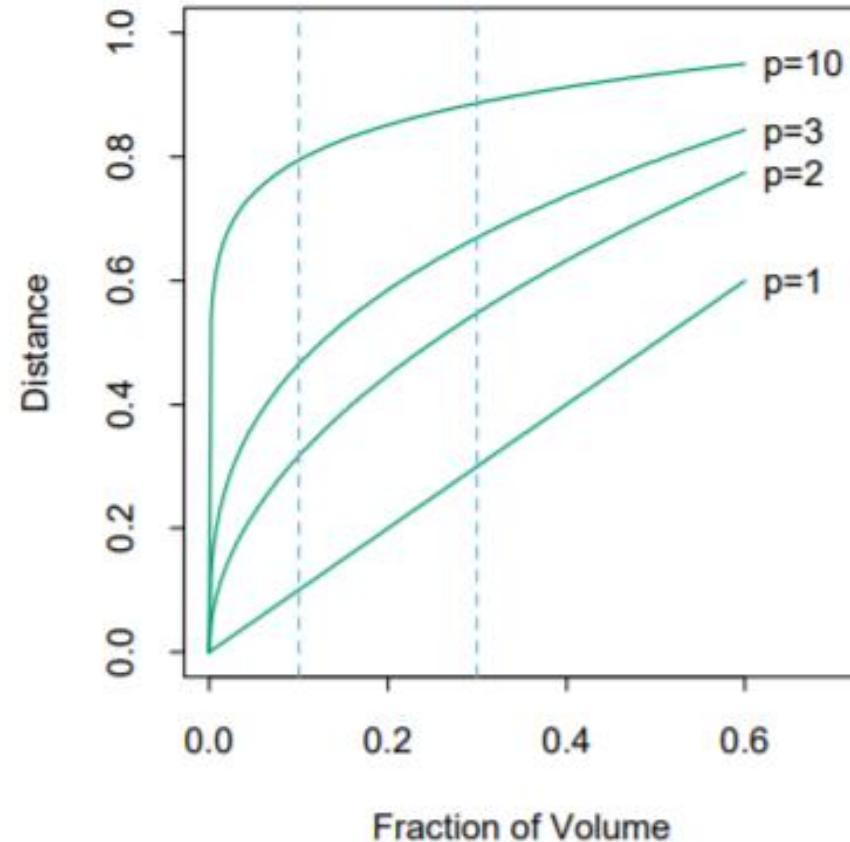
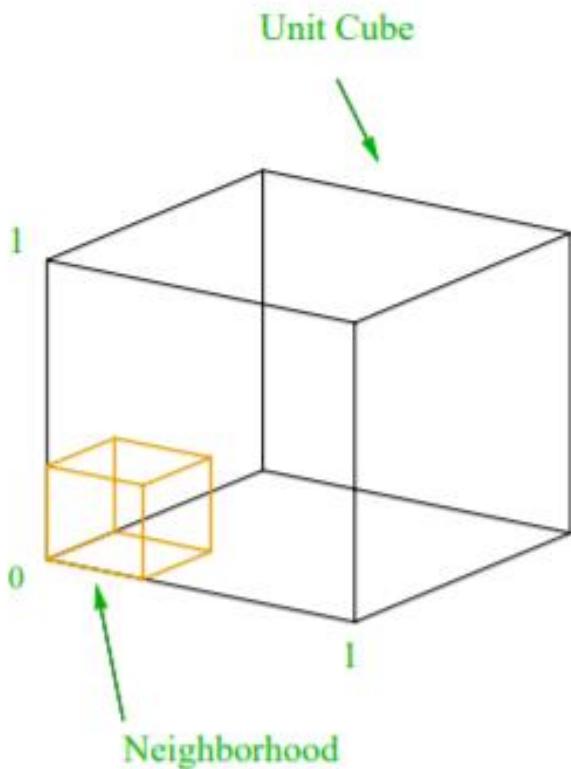
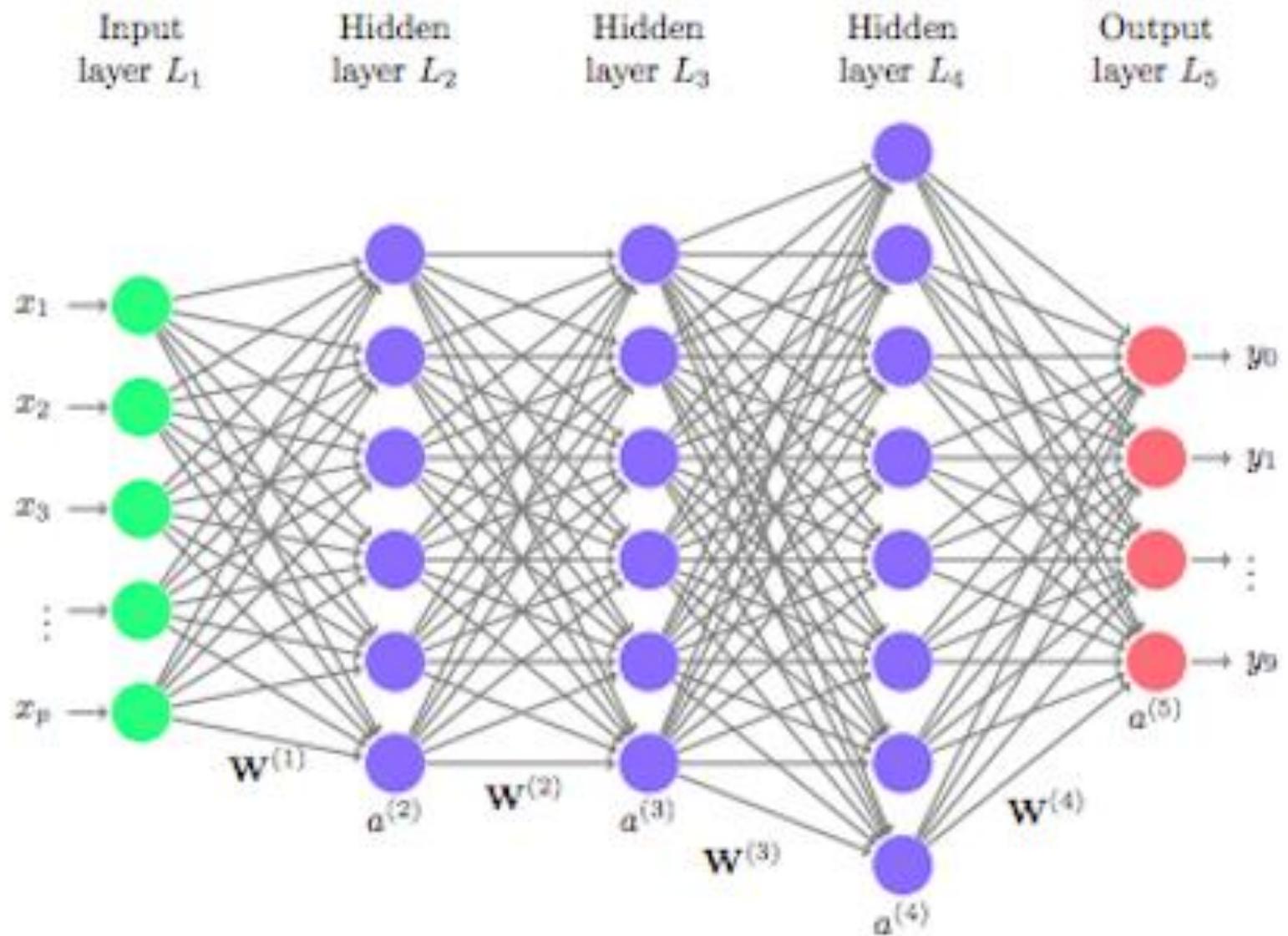
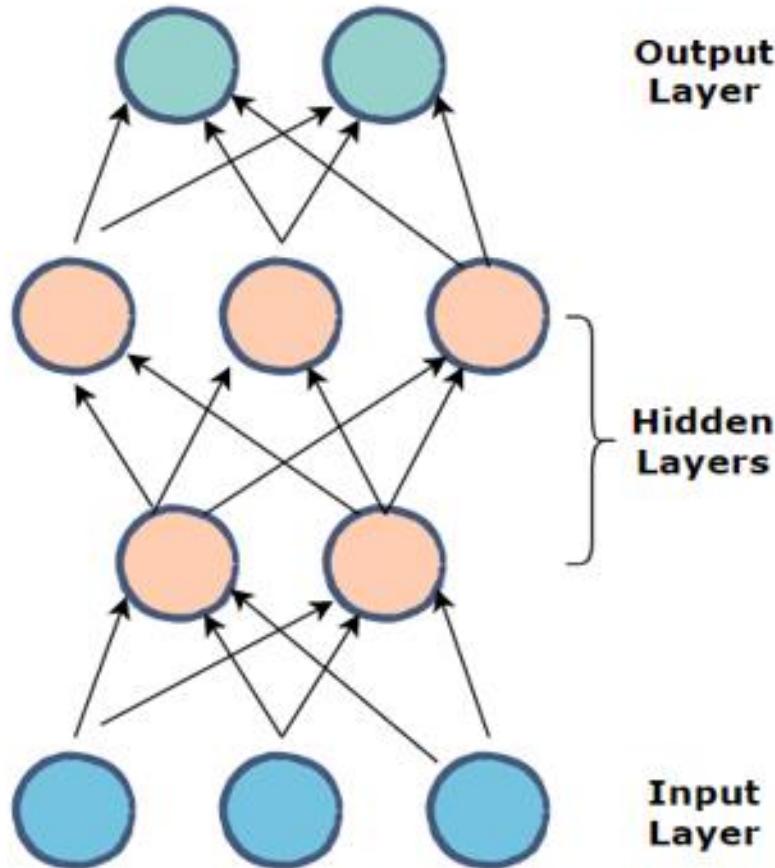


FIGURE 2.6. The curse of dimensionality is well illustrated by a subcubical neighborhood for uniform data in a unit cube. The figure on the right shows the side-length of the subcube needed to capture a fraction r of the volume of the data, for different dimensions p . In ten dimensions we need to cover 80% of the range of each coordinate to capture 10% of the data.

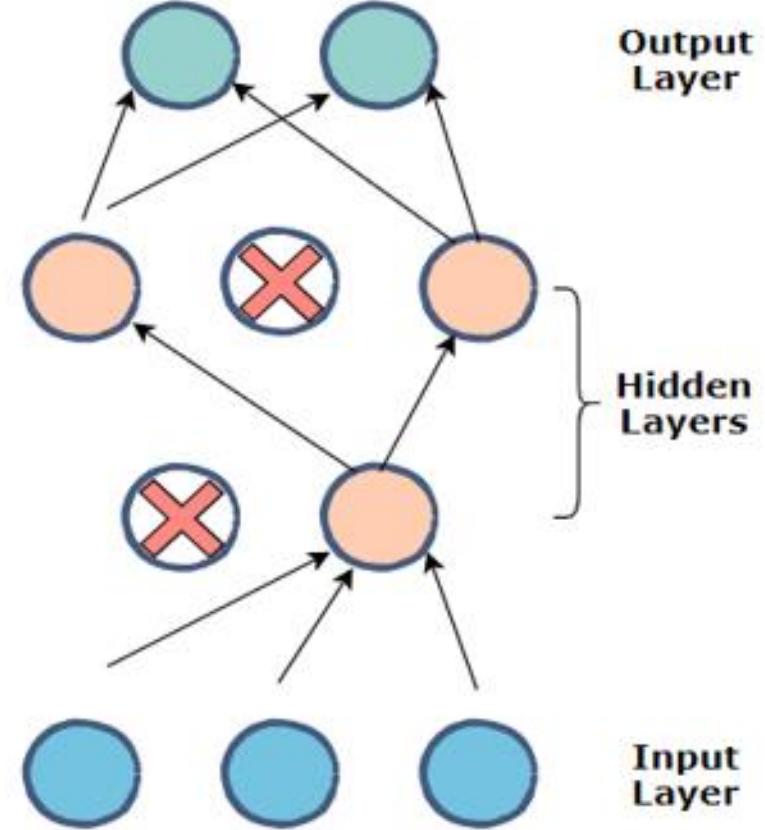
Dimensionality in DL



2. What is Dropout?

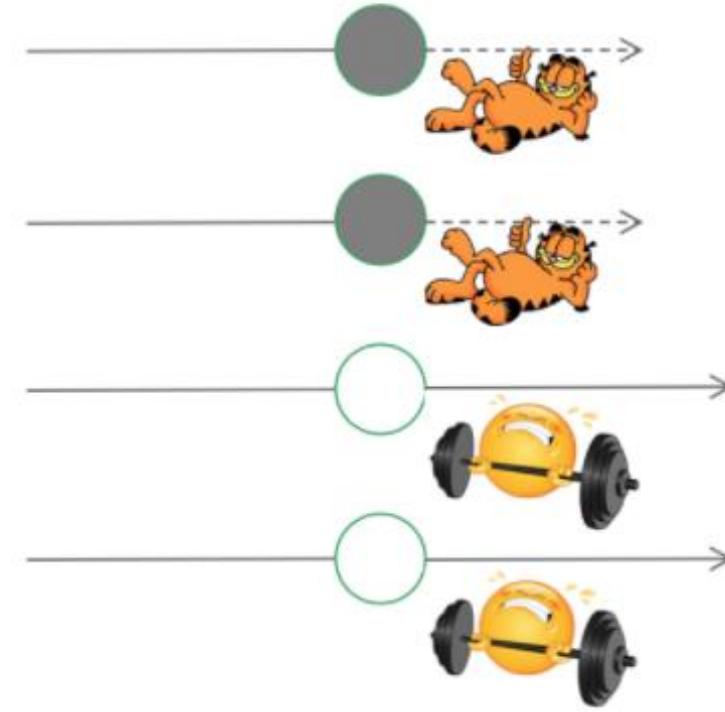
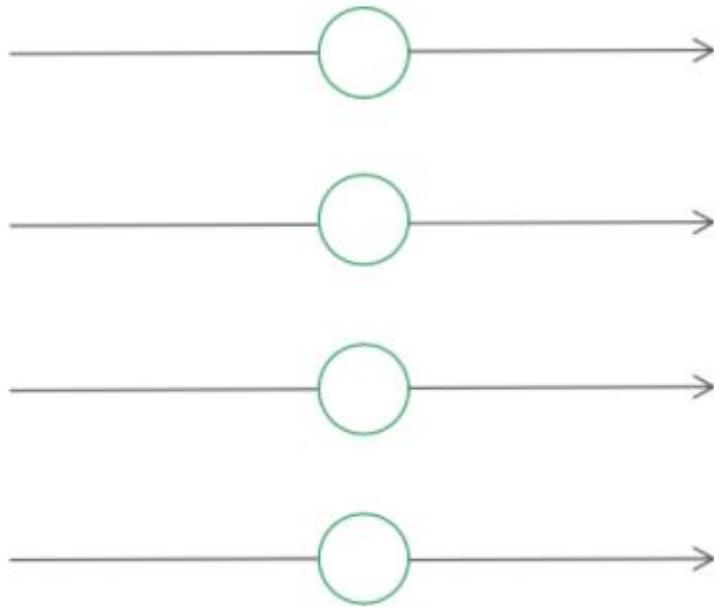


Original network

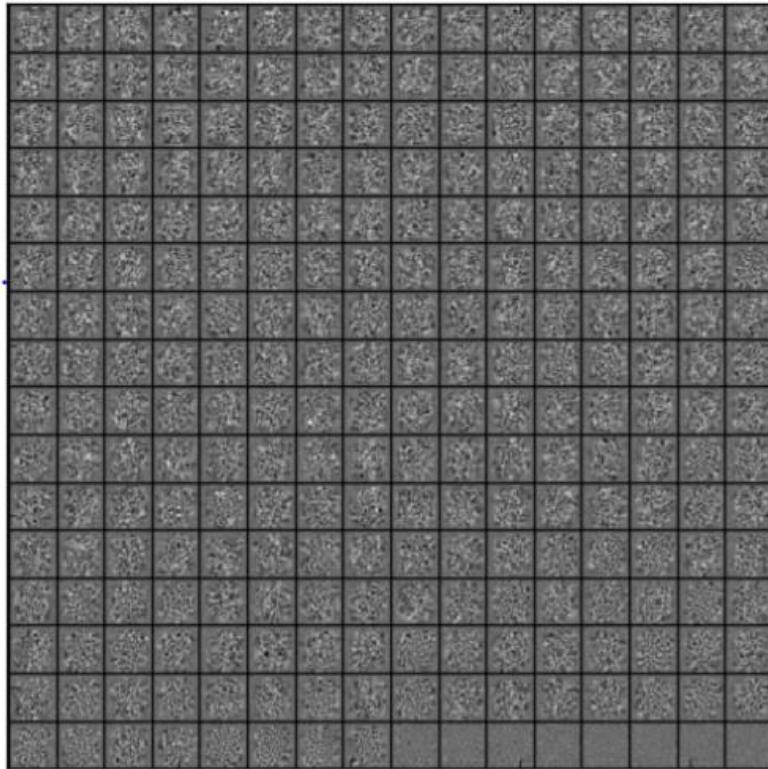


drop out networks

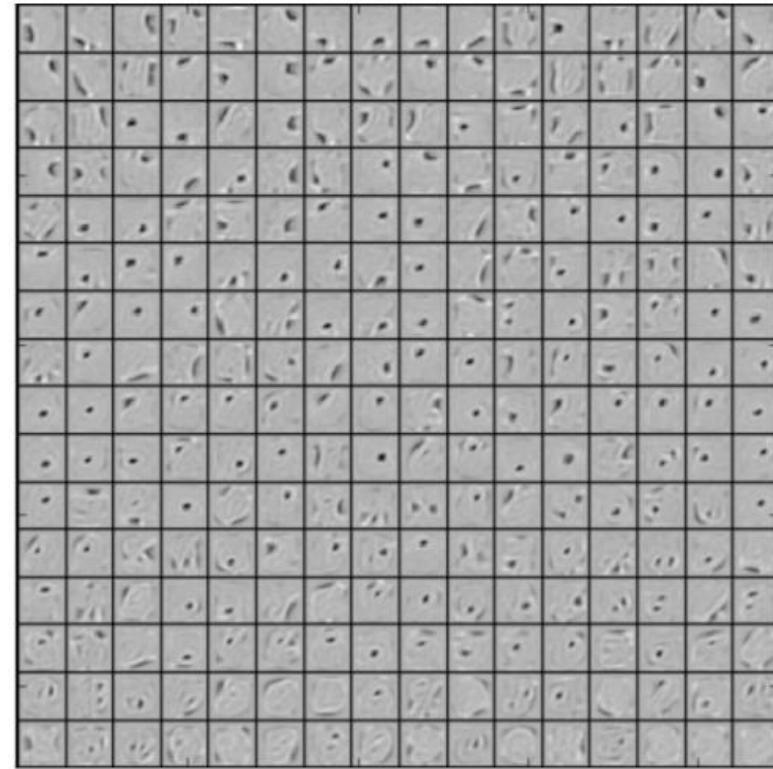
Co-adaptation



Co-adaptation

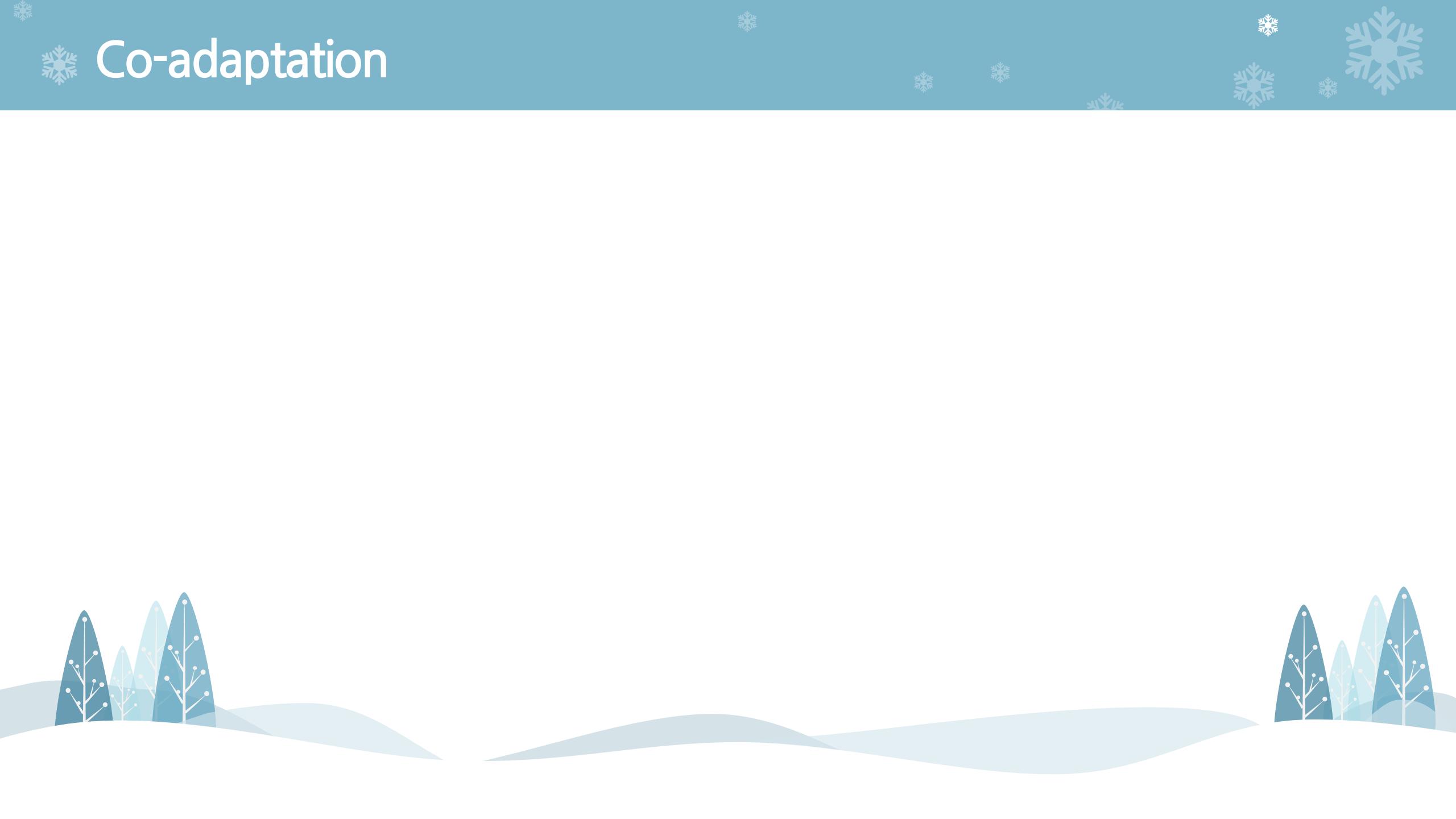


(a) Without dropout

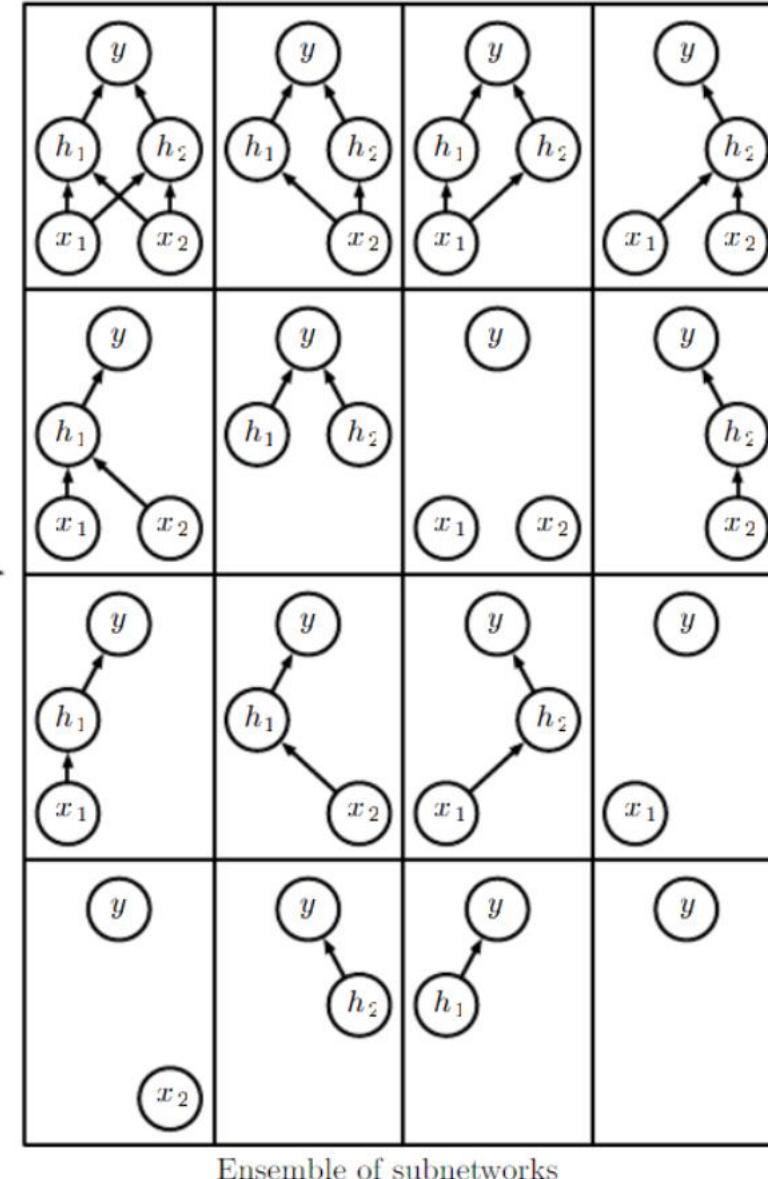
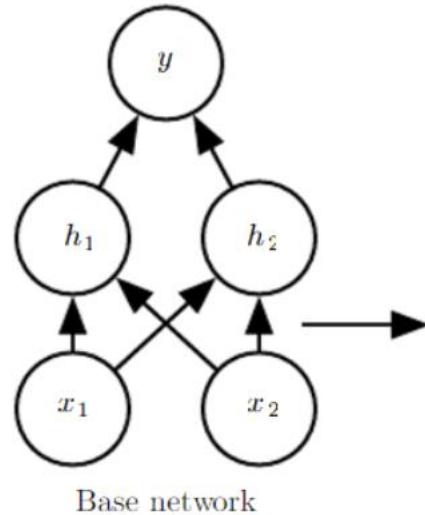


(b) Dropout with $p = 0.5$.

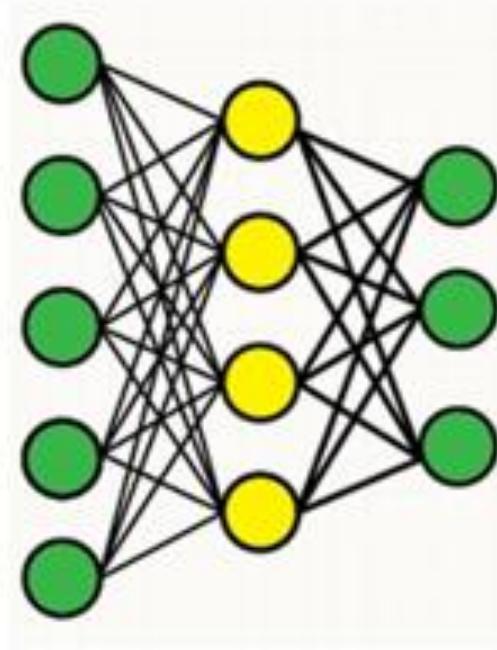
Co-adaptation



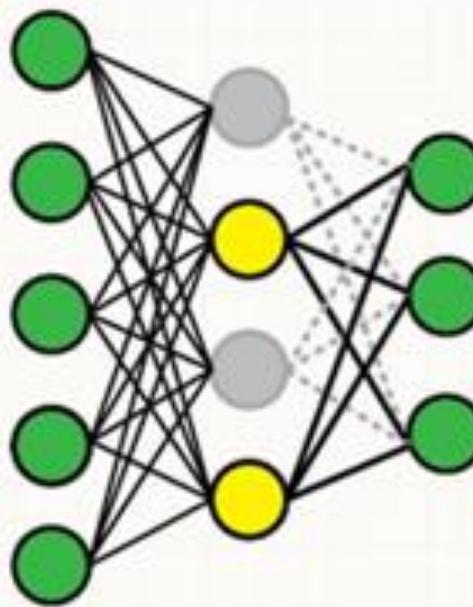
Bagging ~ Dropout



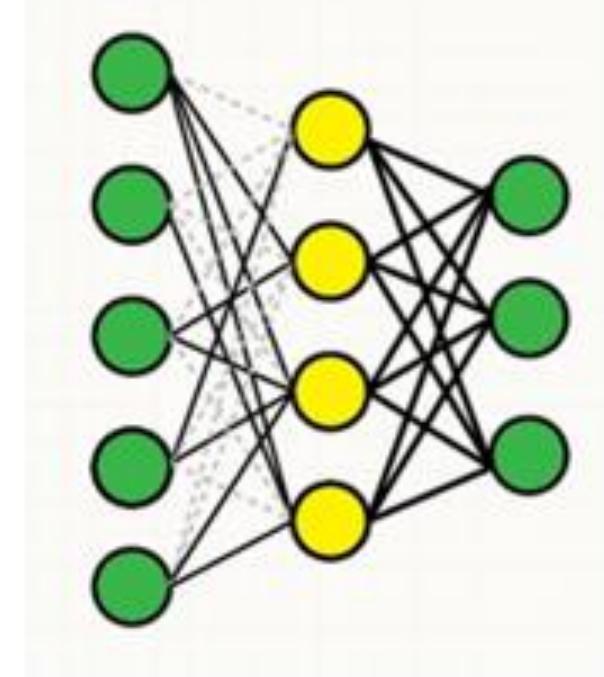
Generalization : Drop-connect



Original

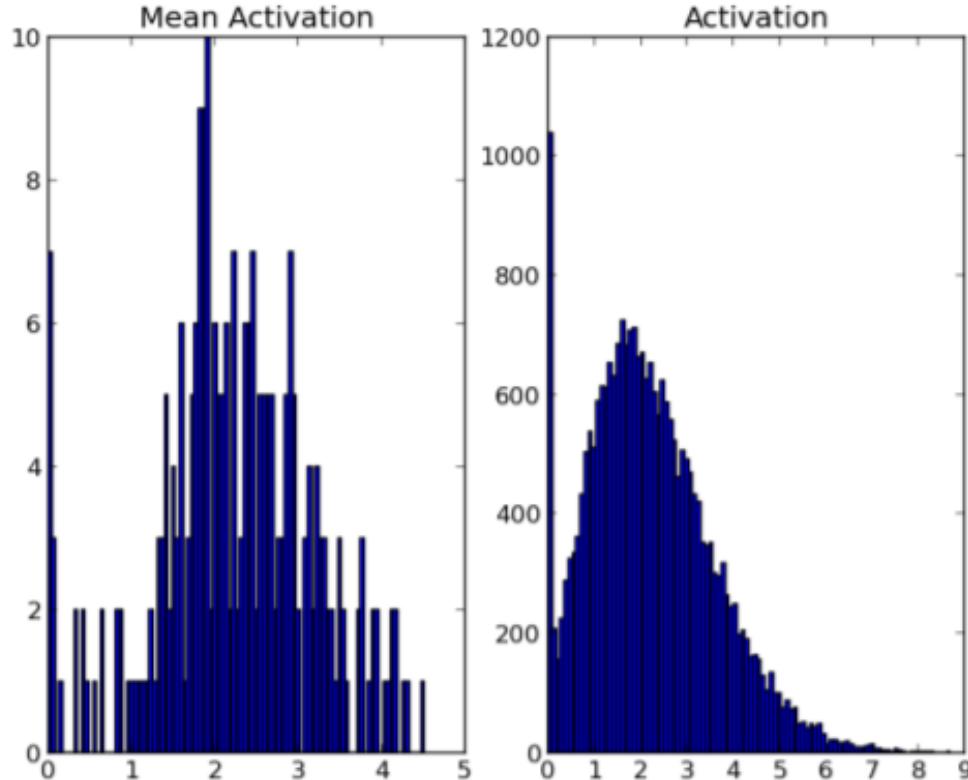


Dropout

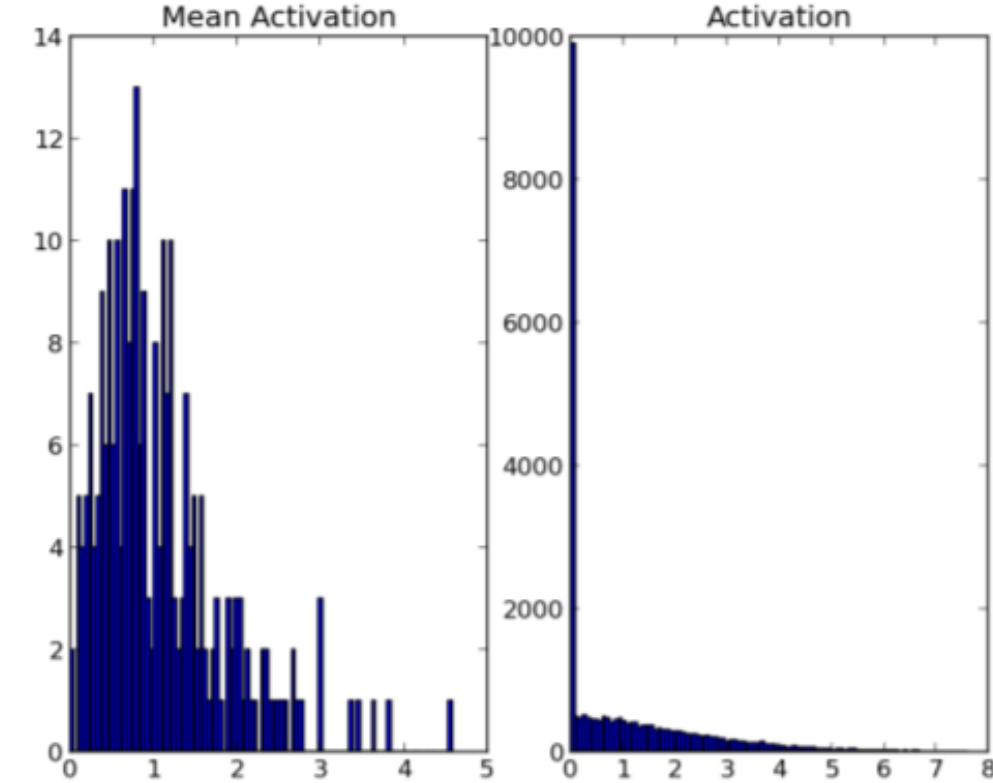


DropConnect

Effect of Drop-connect

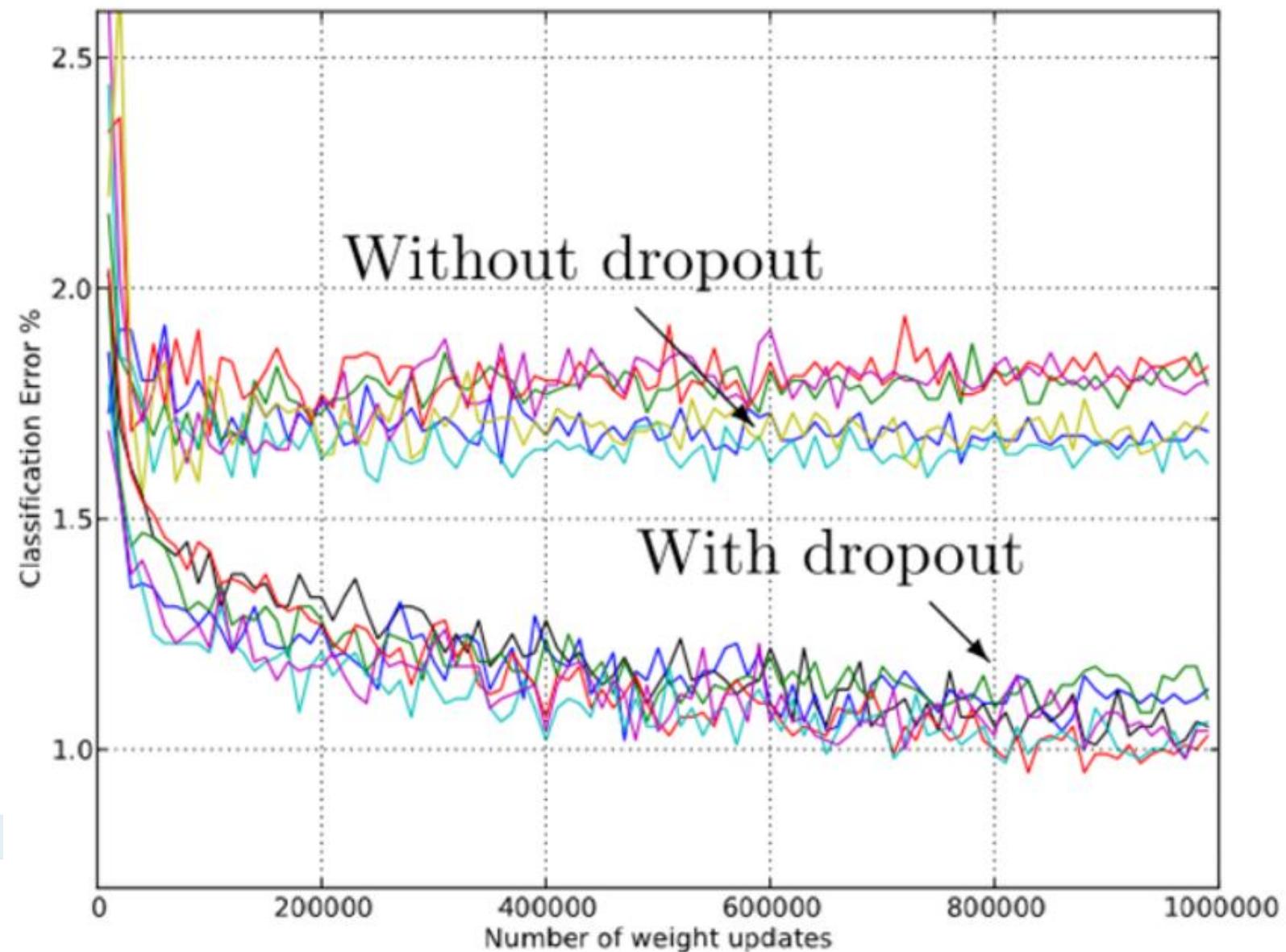


(a) Without dropout

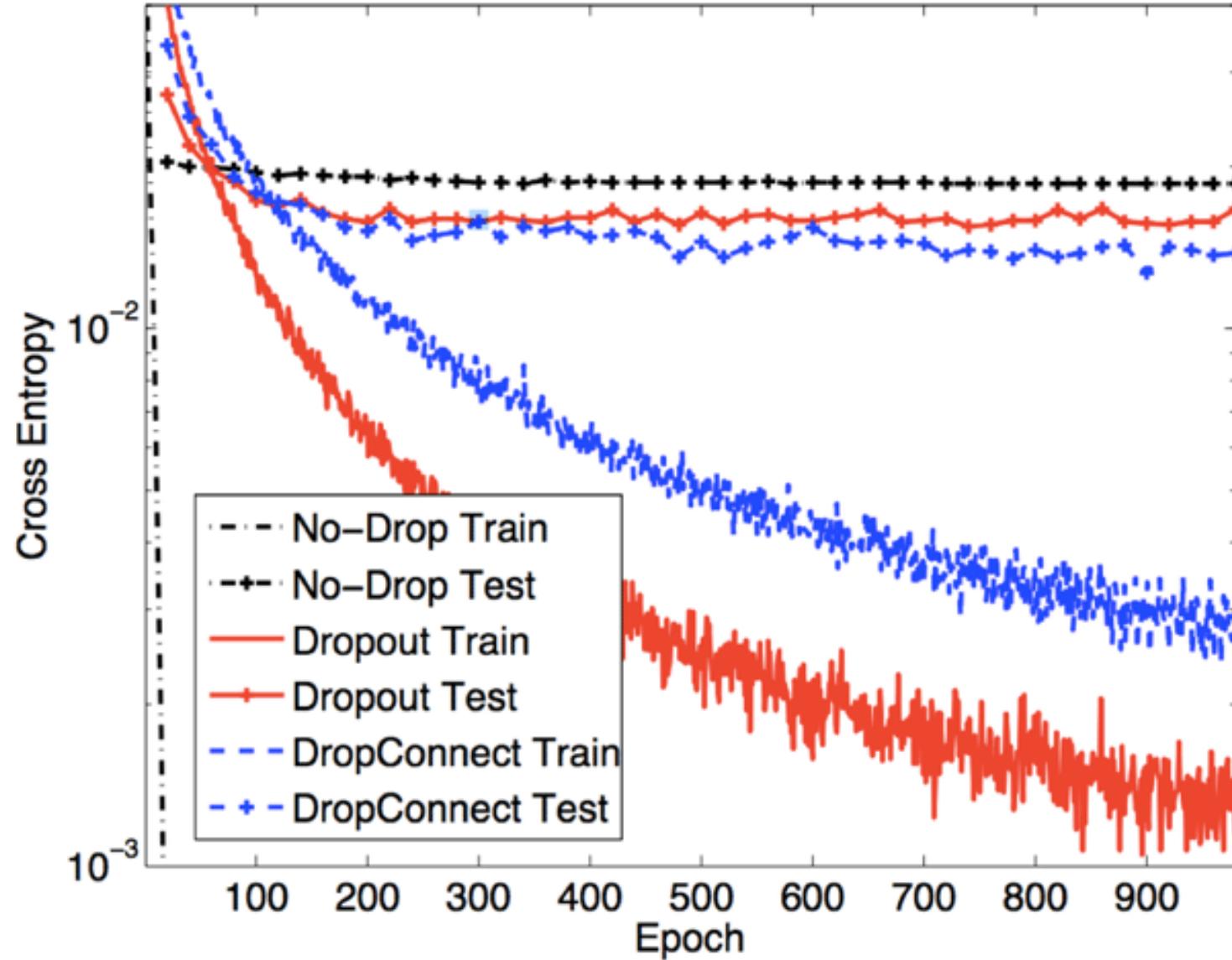


(b) Dropout with $p = 0.5$.

Effect of Drop-connect



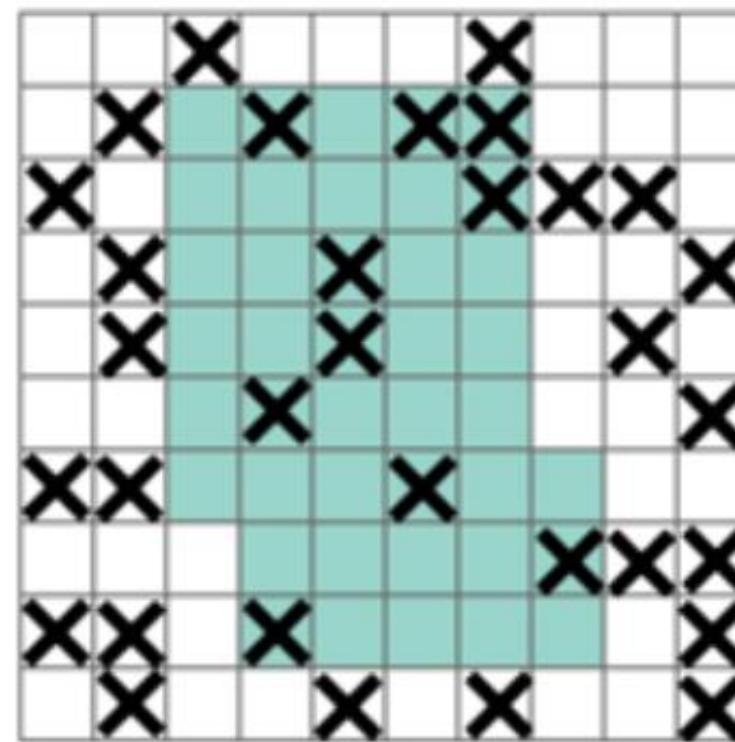
Effect of Drop-connect



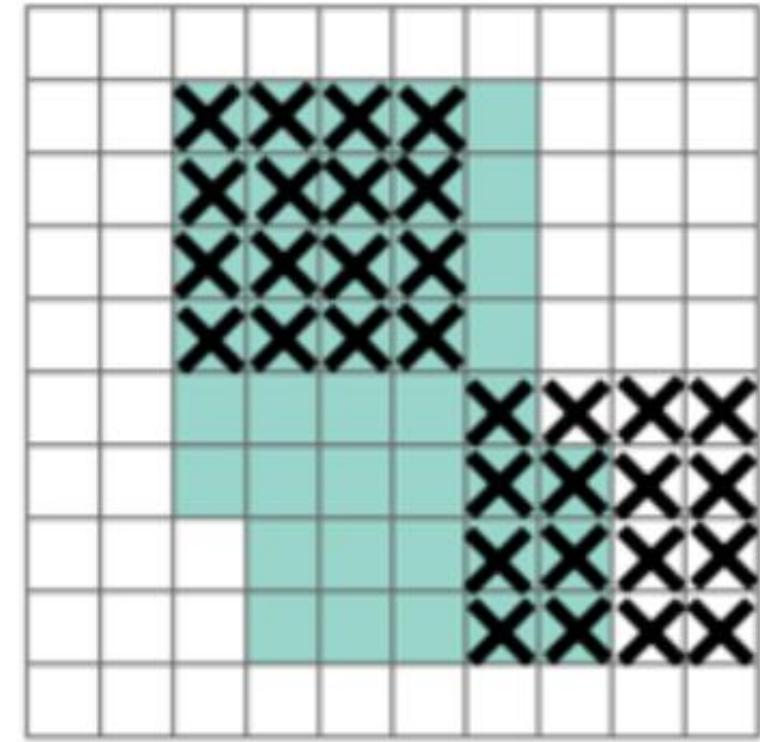
Dropout in CNN



(a)



(b)

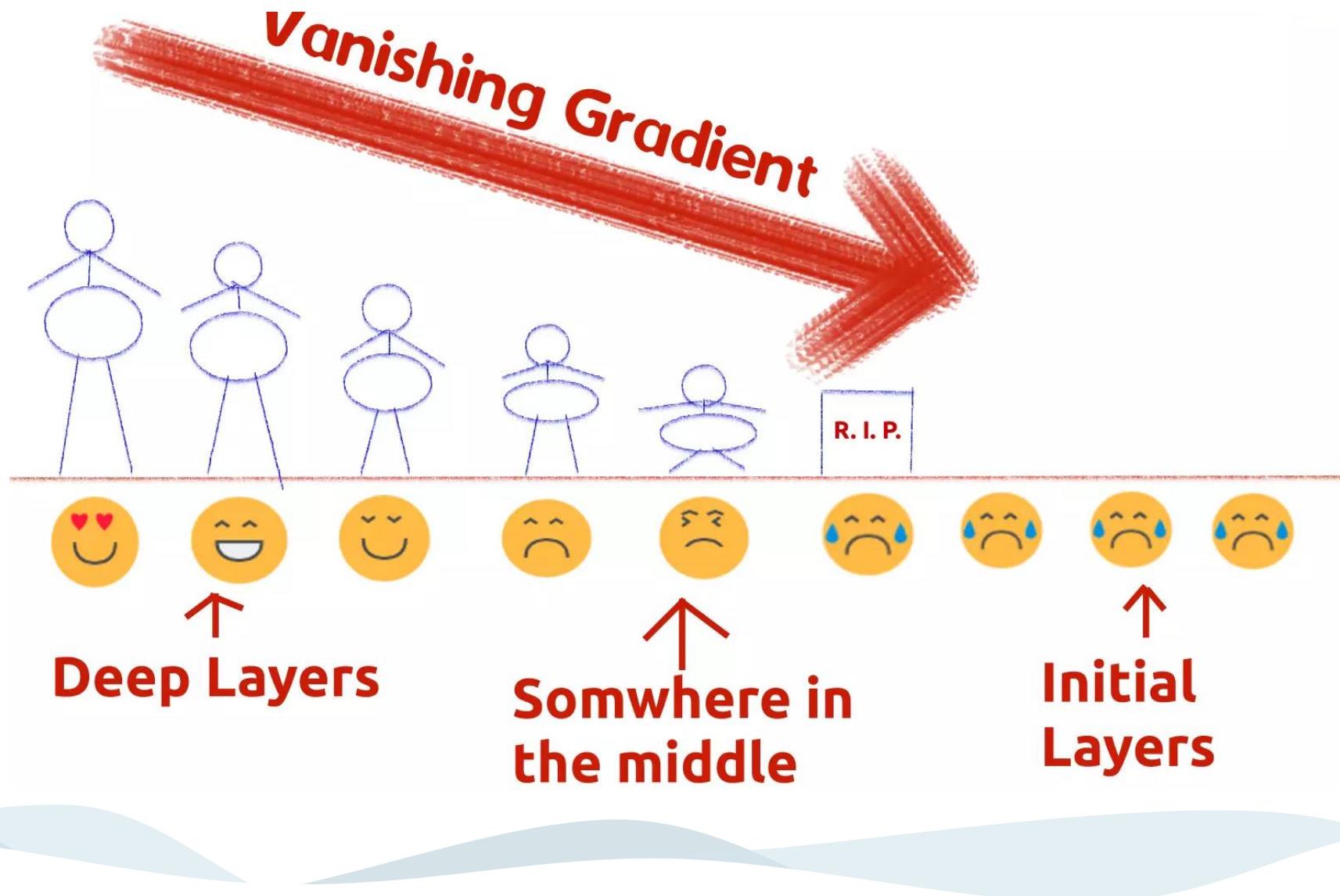


(c)

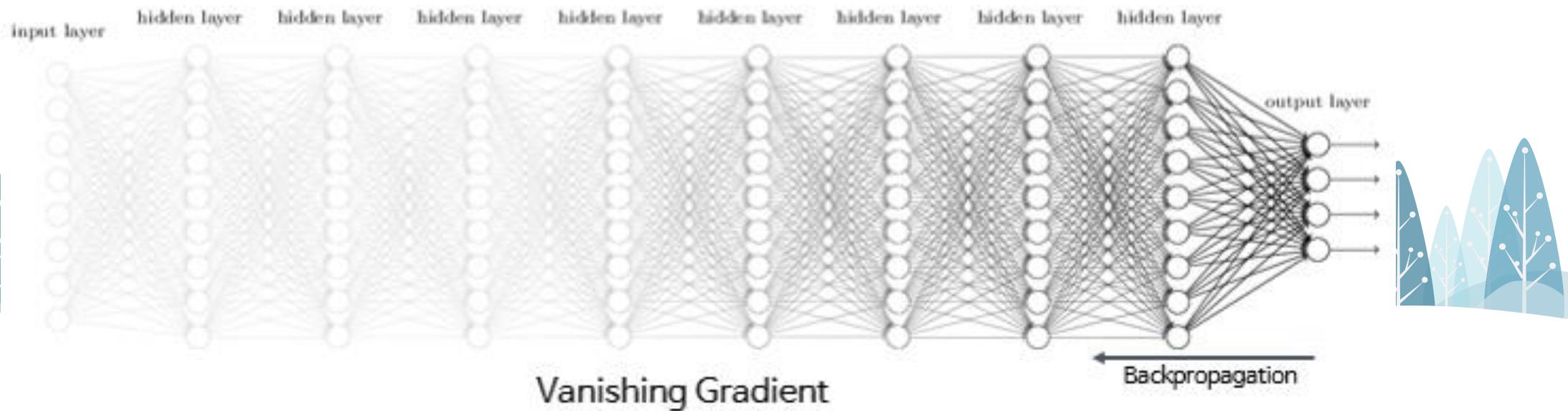
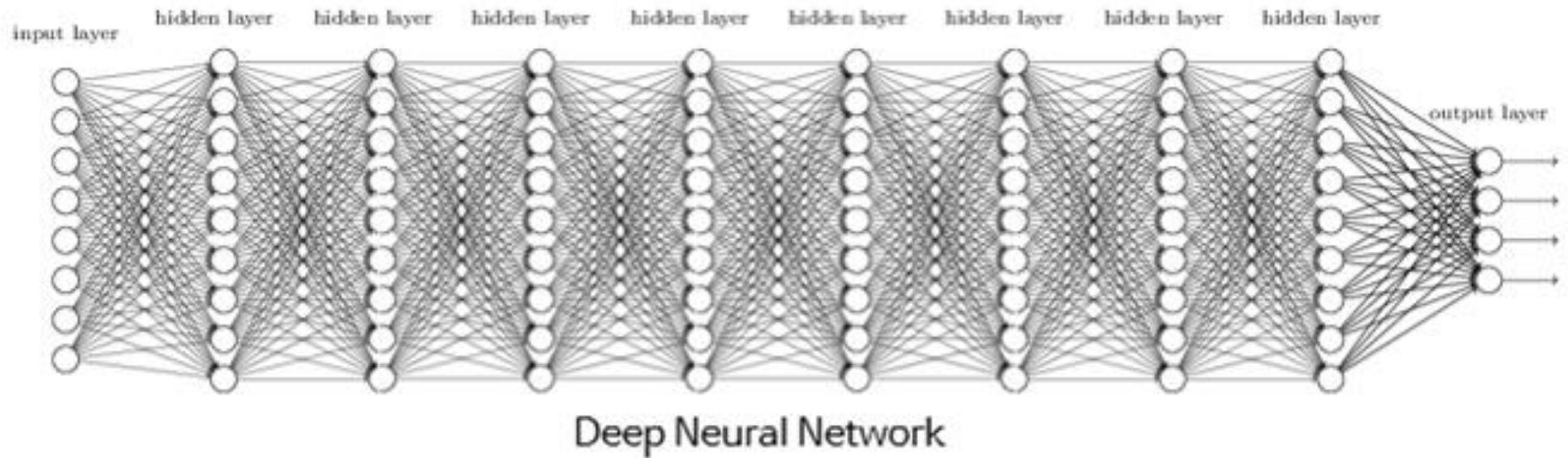
Under-fitting Model에 사용 X

1. Neuron을 확률적으로 누락시키는 기법
2. 만약, under-fitting model에 사용한다면, 성능 악화
3. CNN Model에서 사용의 제한점
4. 수렴 속도 저하의 문제 & Training error는 증가할 수 있음

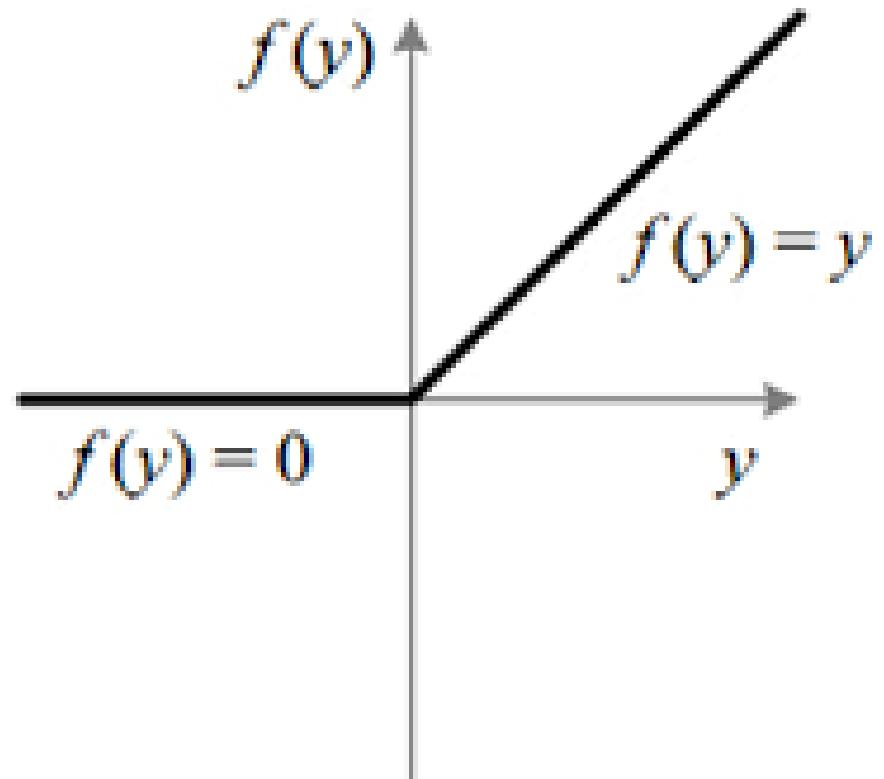
3. Vanishing / Exploding Gradient



3. Vanishing / Exploding Gradient



3. Vanishing / Exploding Gradient



Weight initialization method

- The right weight initialization method can speed up time-to-convergence considerably. The choice of your initialization method depends on your activation function. Some things to try:
 - When using ReLU or leaky RELU, use [He initialization](#)
 - When using SELU or ELU, use [LeCun initialization](#)
 - When using softmax, logistic, or tanh, use [Glorot initialization](#)
 - Most initialization methods come in uniform and normal distribution flavors.

Activation Function

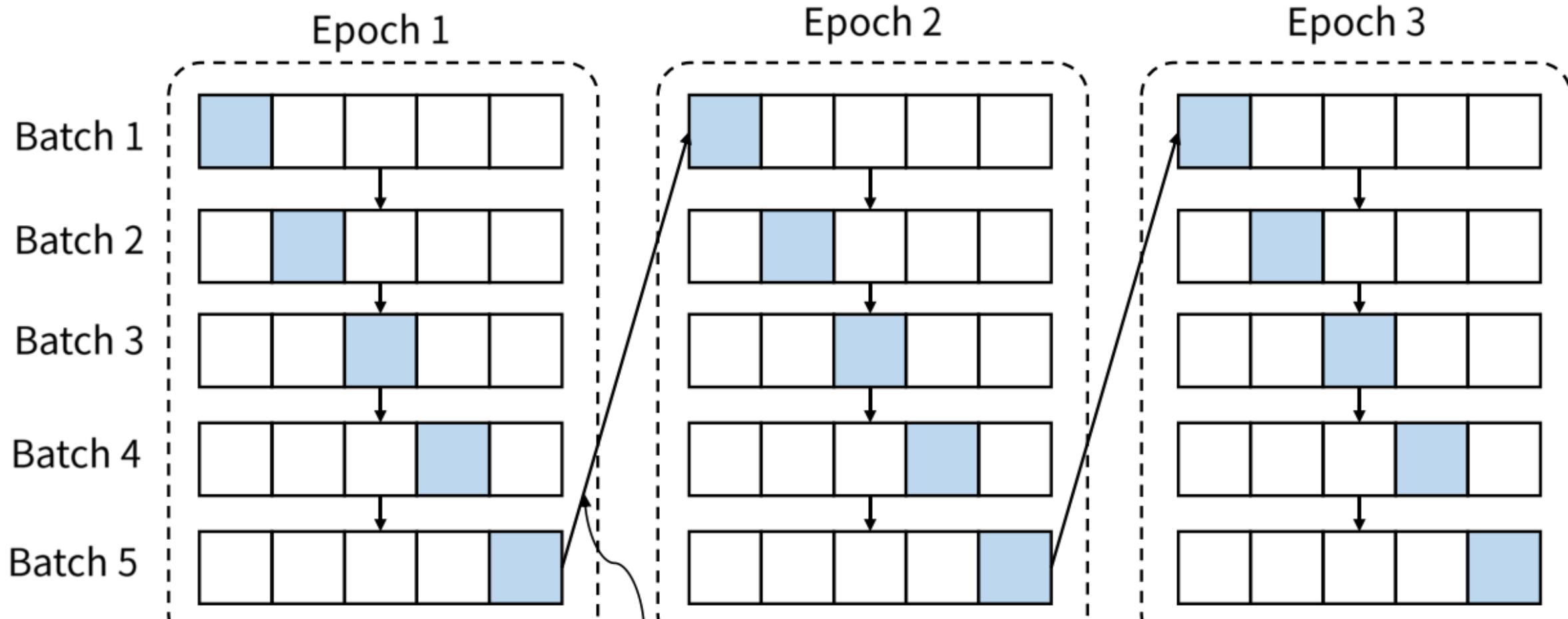
1. ReLu의 문제점에 대해

2. ReLu를 개선 or 응용한 activation function

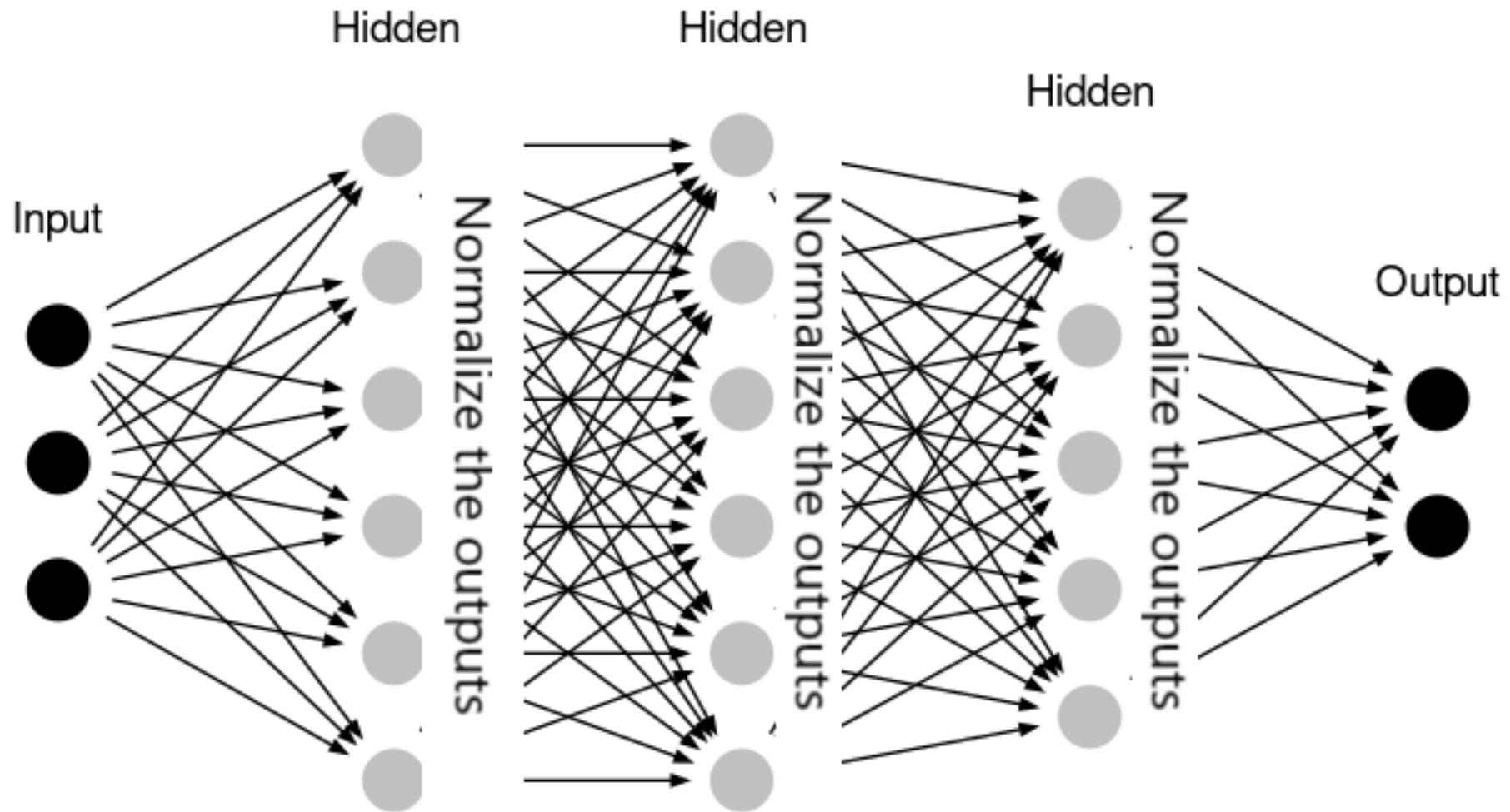
LIST

Leaky ReLu / SELU / ELU /
Swish / GELU

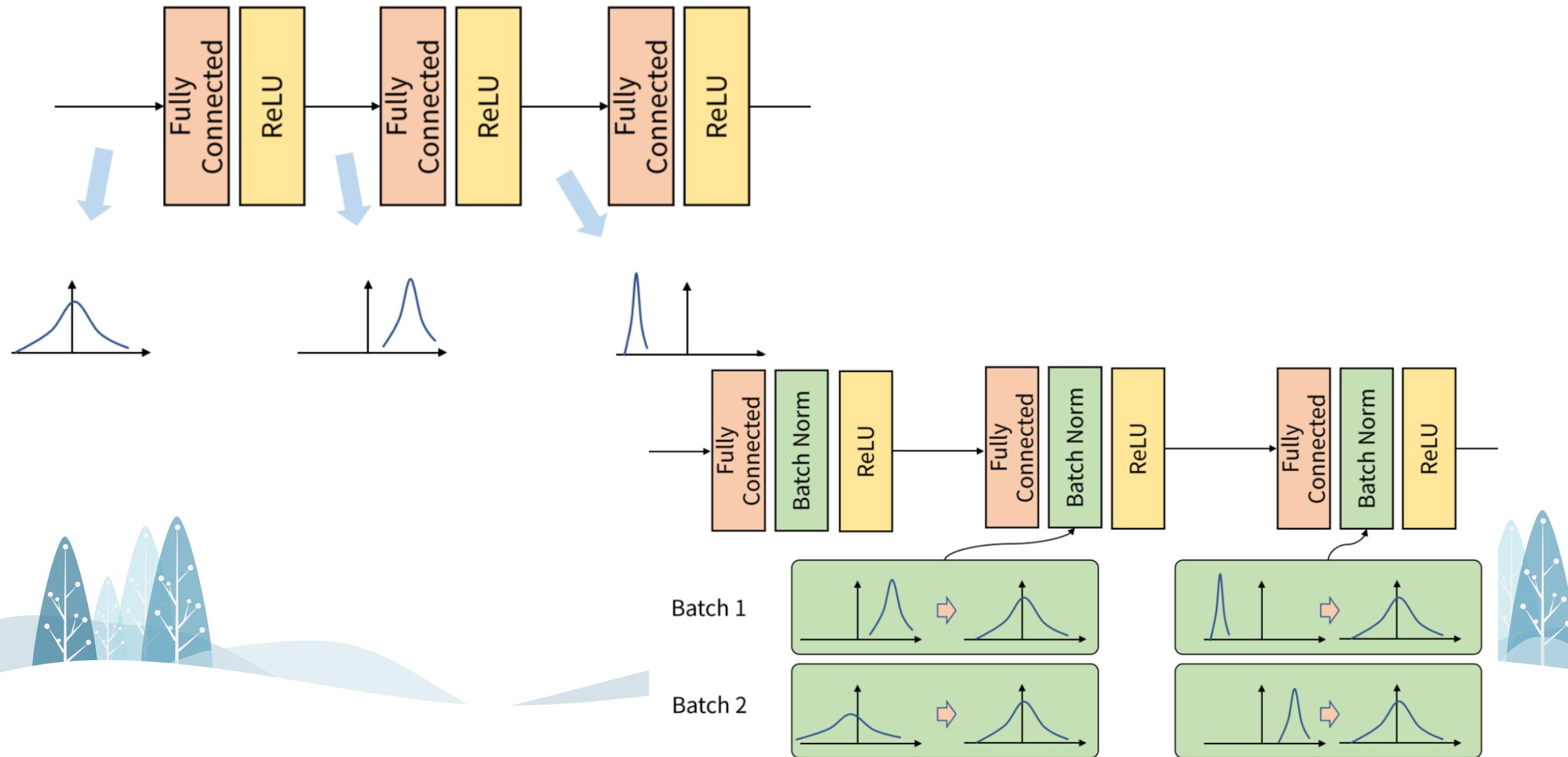
Review of Batch & Epoch



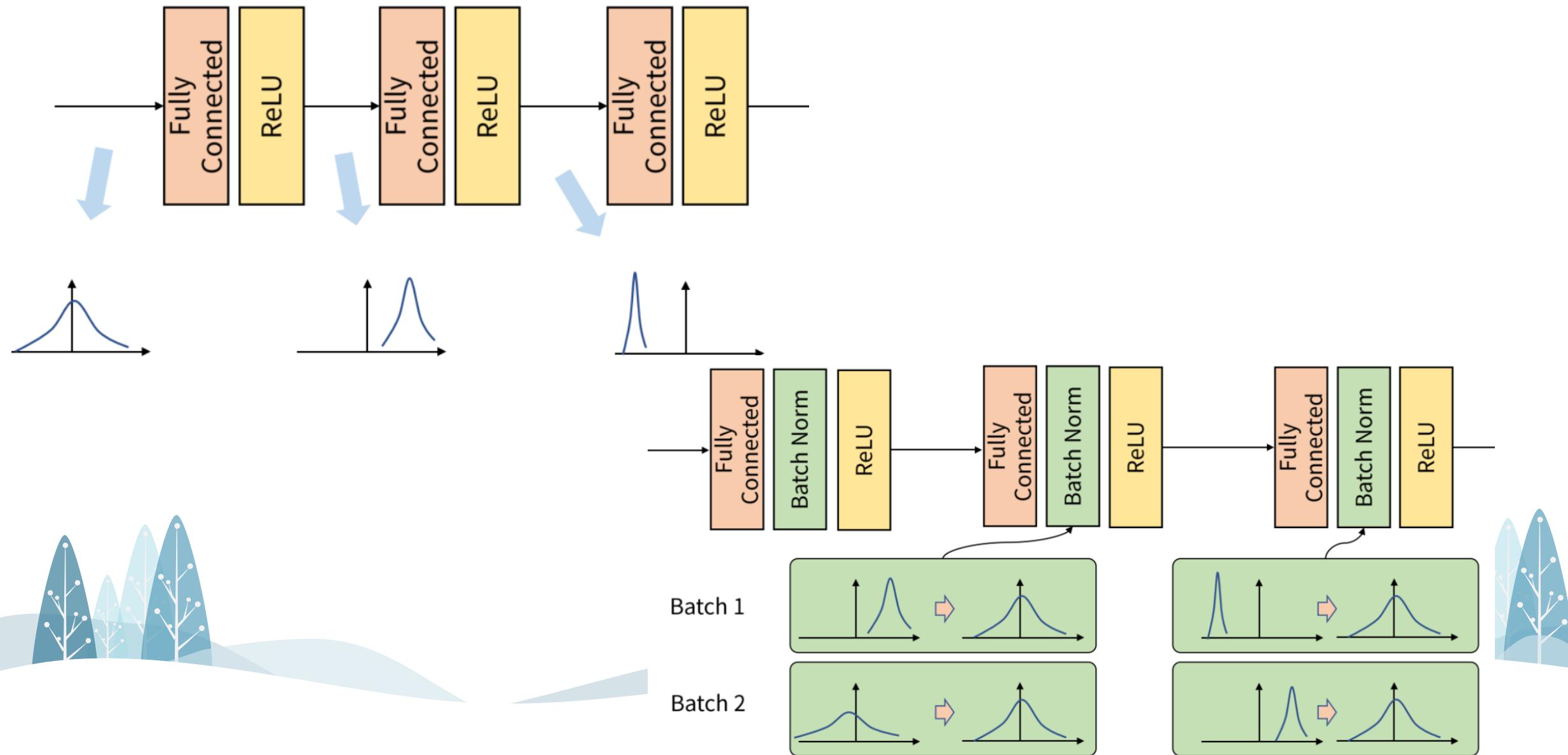
3. Batch Normalization



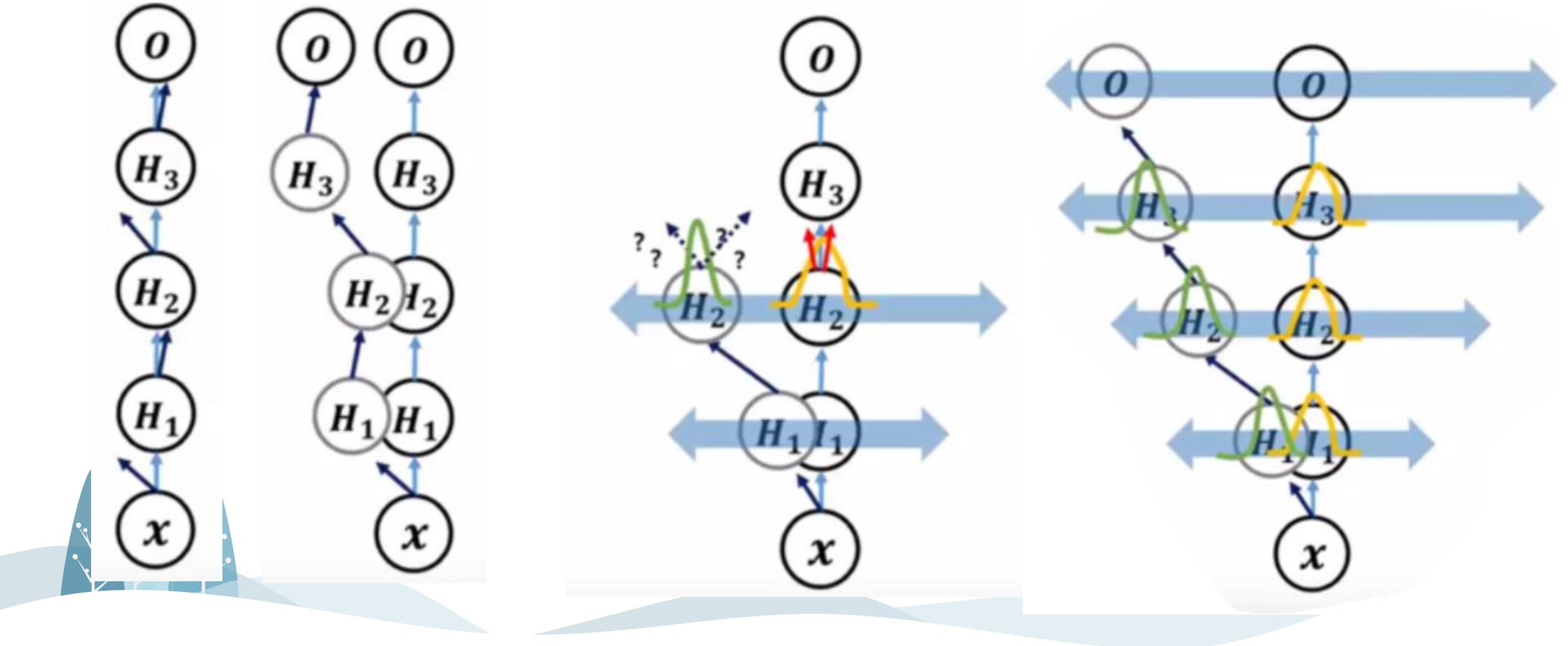
Internal Covariance Shift



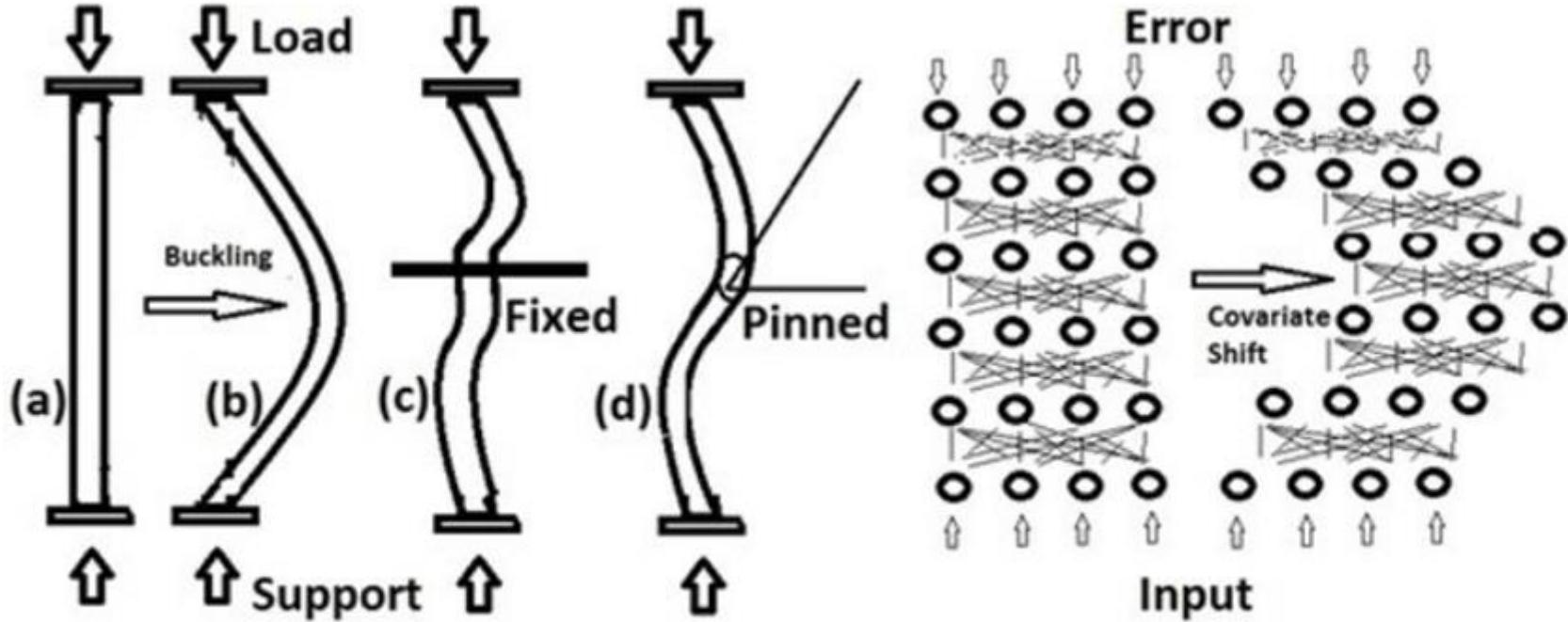
Internal Covariance Shift



Internal Covariance Shift



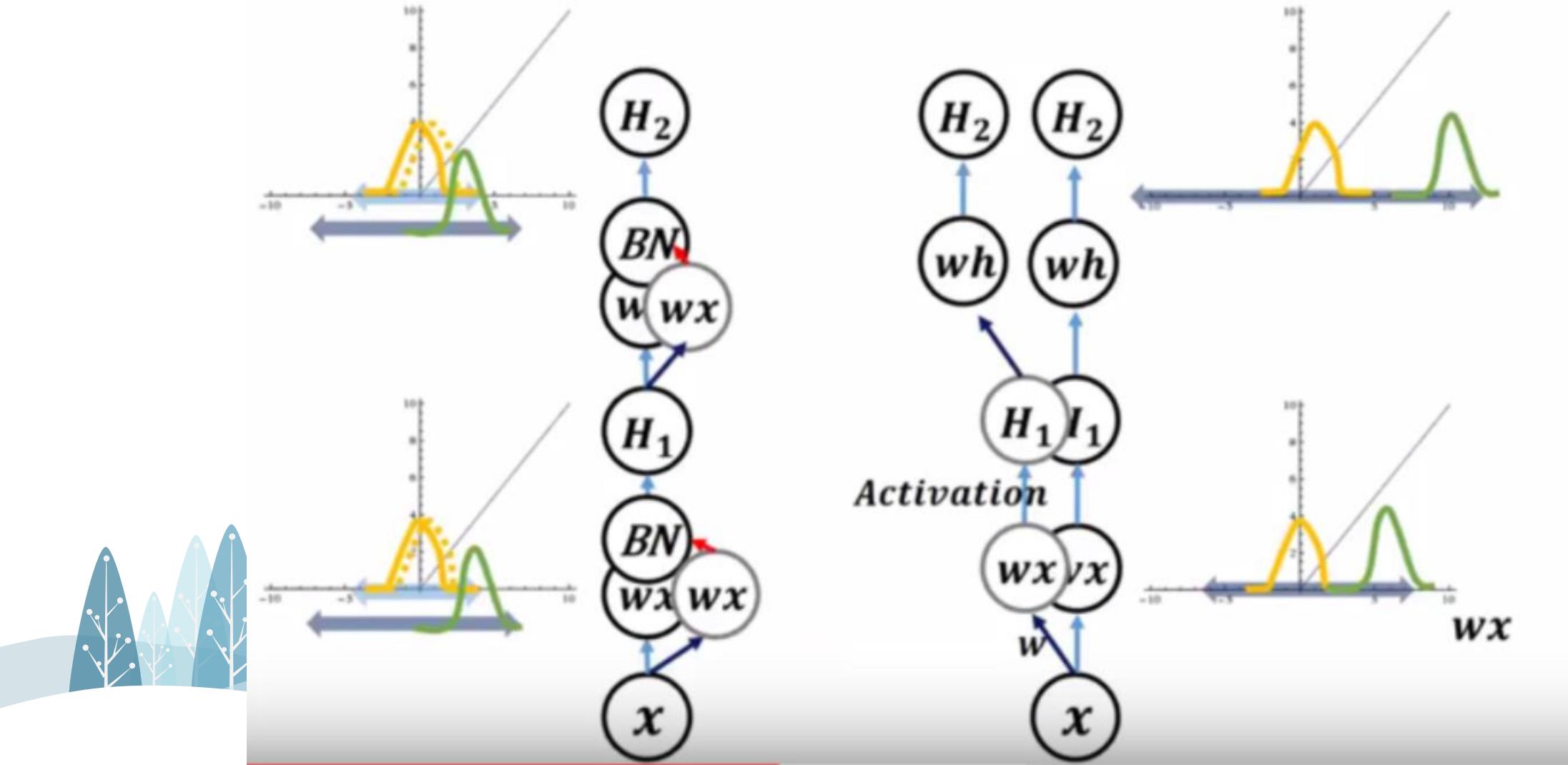
Internal Covariance Shift



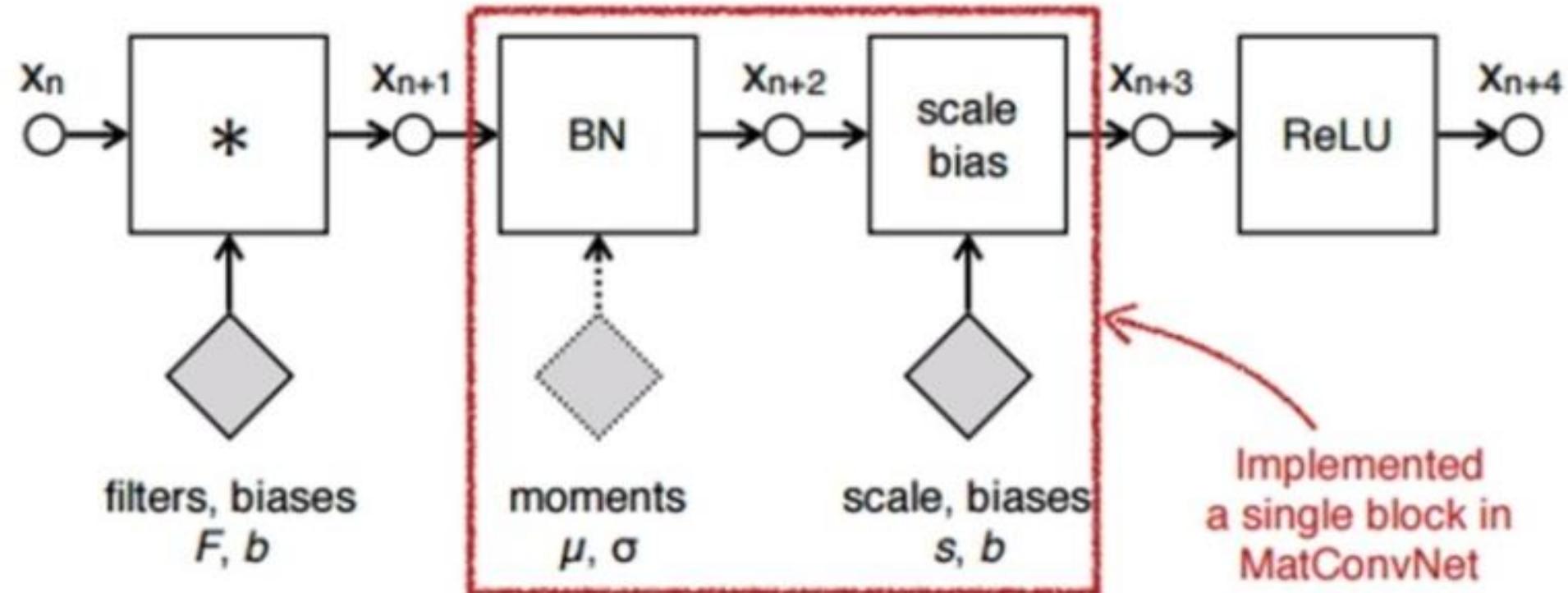
For both, Buckling or Co-Variate Shift a small perturbation leads to a large change in the later.

Debiprasad Ghosh, PhD, Uses AI in Mechanics

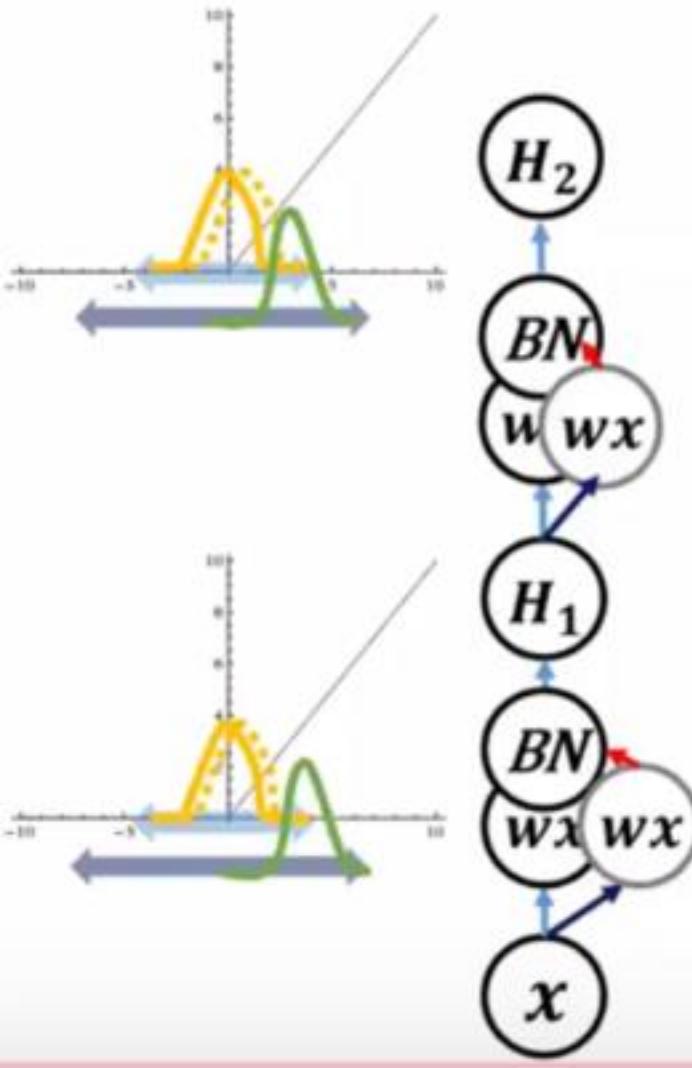
Internal Covariance Shift



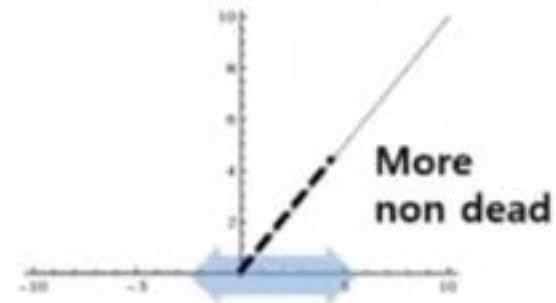
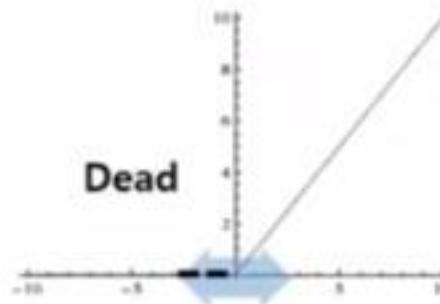
Internal Covariance Shift



Scale & Shifting Parameters



0 mean and 1 variance is preferred
But it is not best

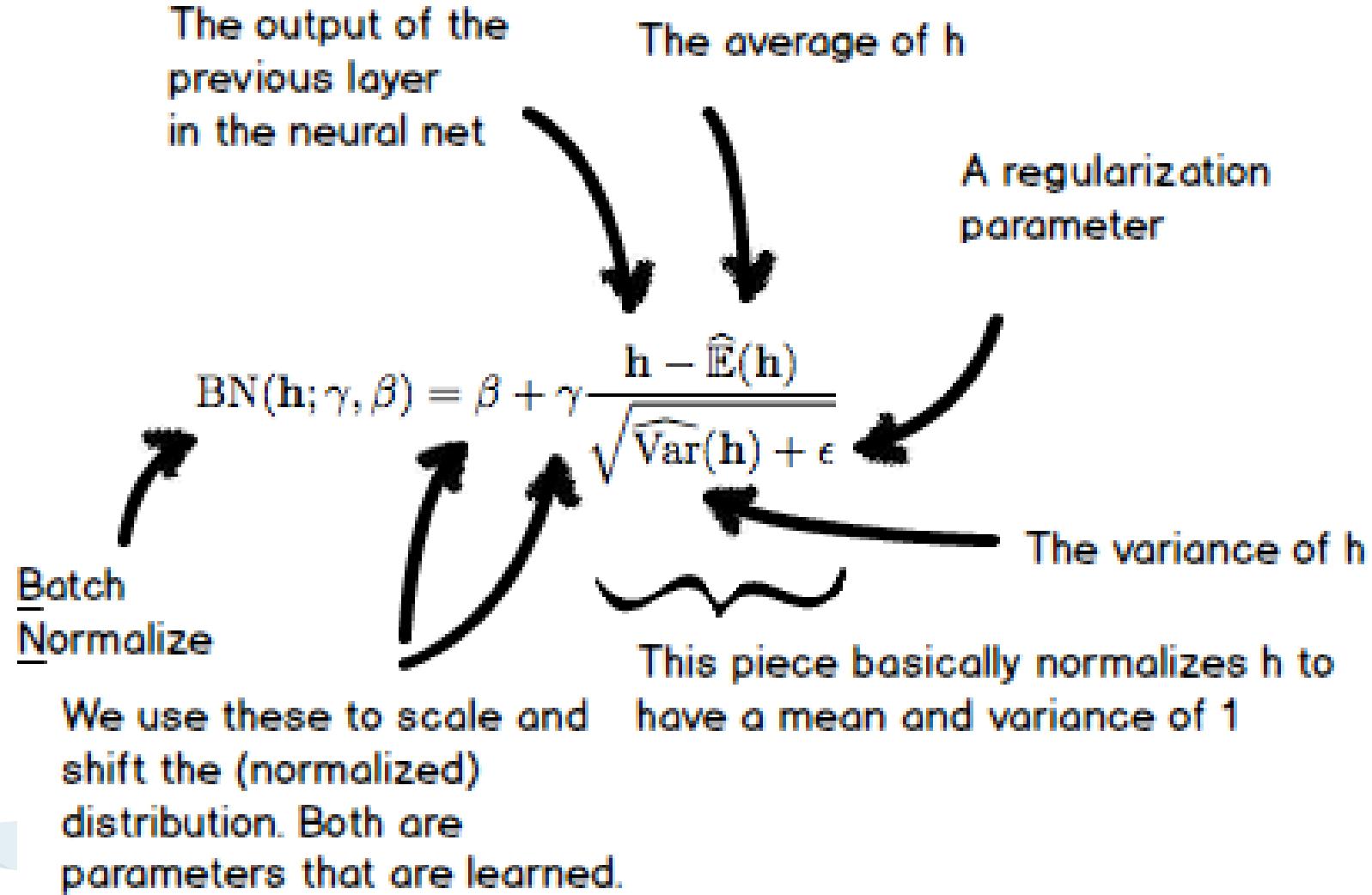


$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i)$$

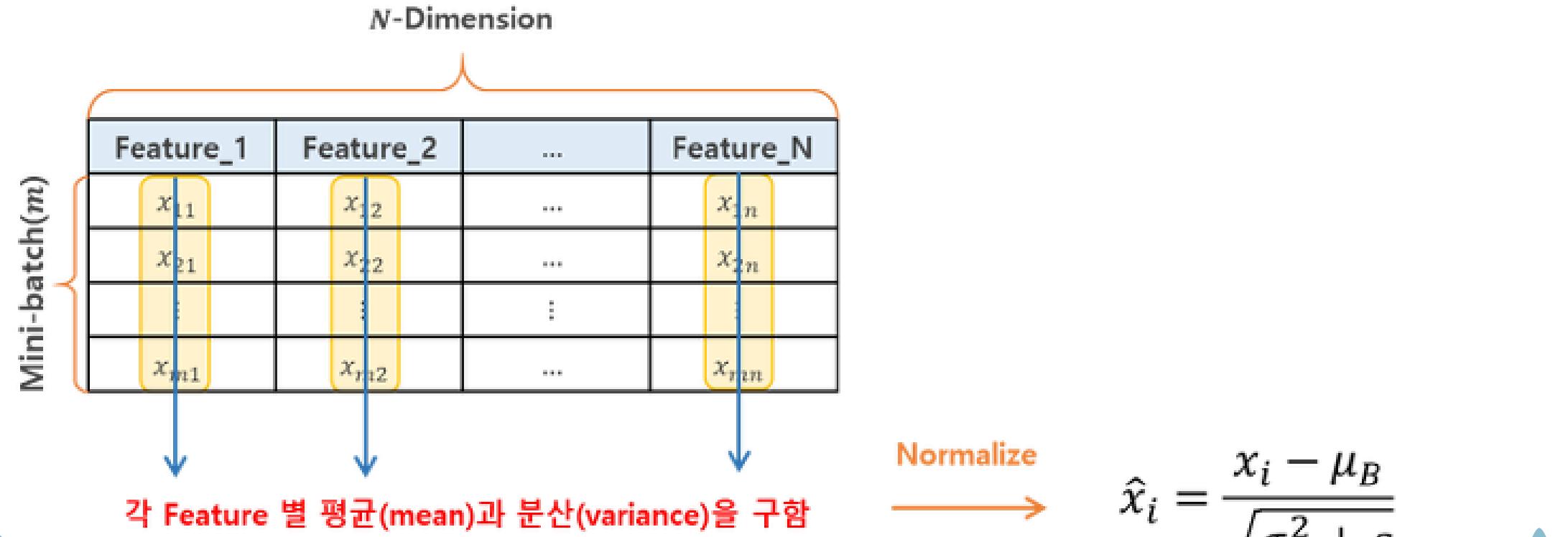
// scale and shift



Batch Normalization



Batch Normalization



Batch Normalization

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$;

Parameters to be learned: γ, β

Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{scale and shift}$$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Effect of Batch Normalization

Fast Convergence

High Learning rate

빠른 학습 속도 구현 가능,
Optimum에 더 빠른 수렴

Not sensitive to Weight Initialization

Initial value에 덜 민감

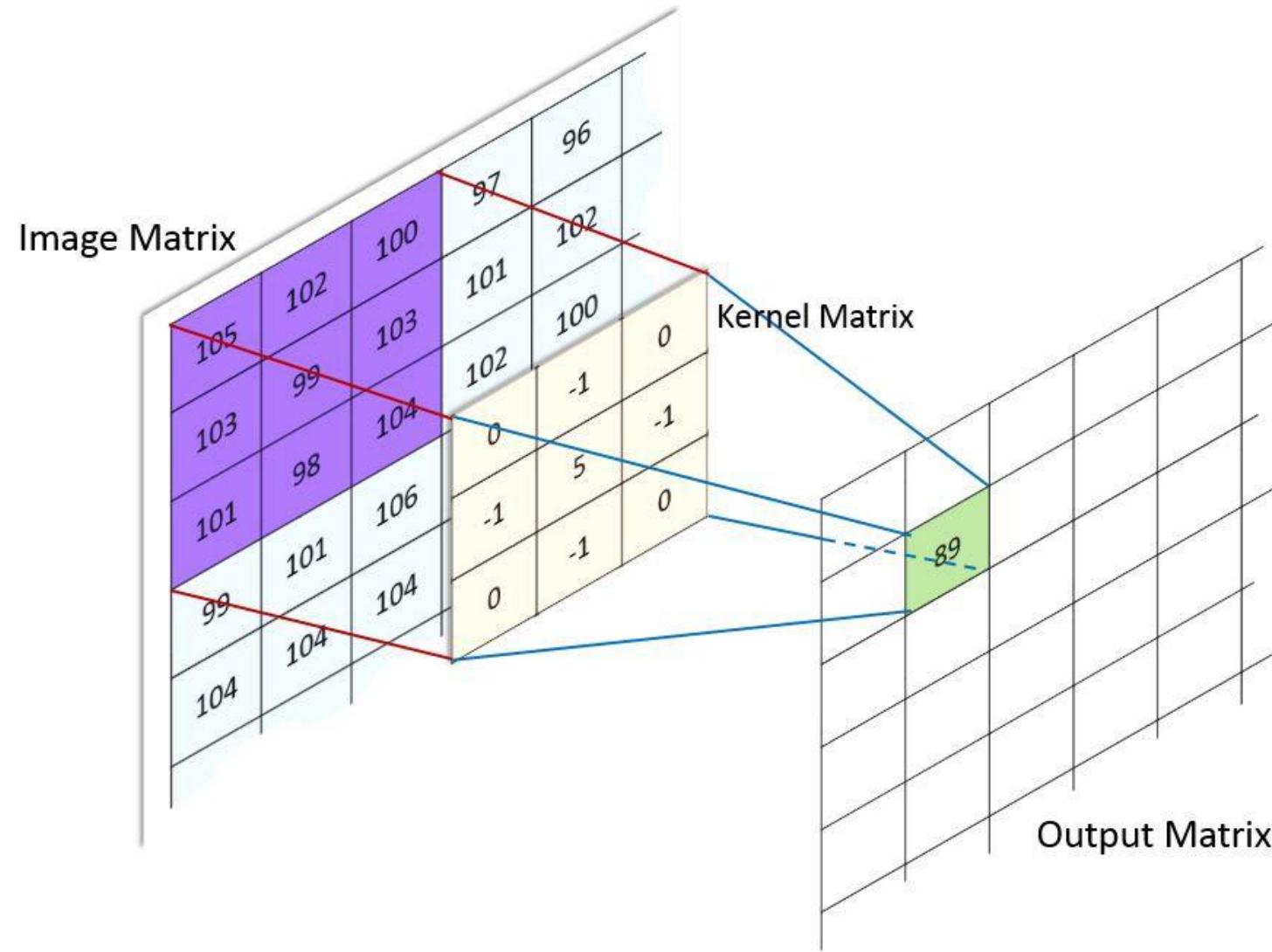
Weight initial Setting 문제를
덜 고민하는 것이 가능

Regularization

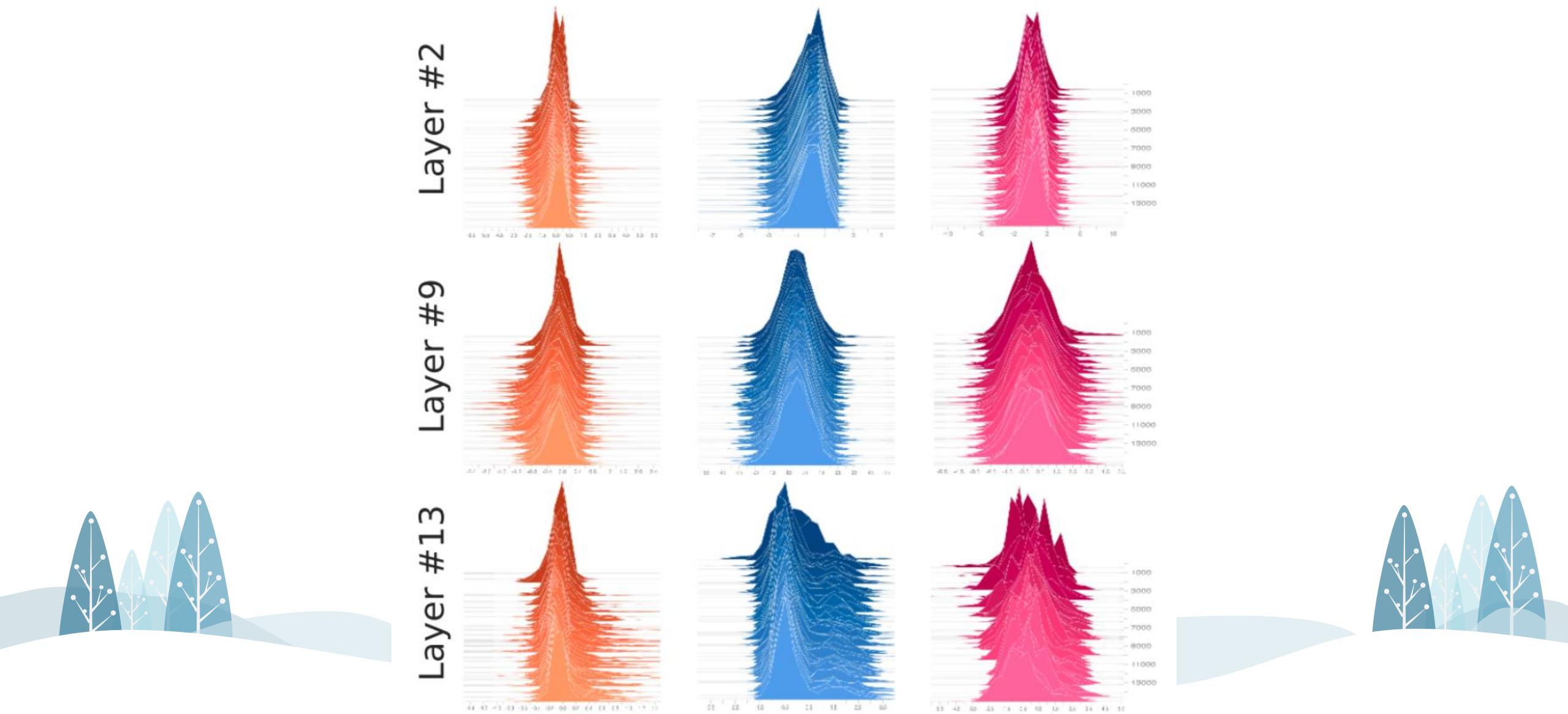
Dropout과 비슷한 효과

Dropout과 Batch Normalization을
동시에 사용하지는 않음!
둘 모두 비슷한 효과!!

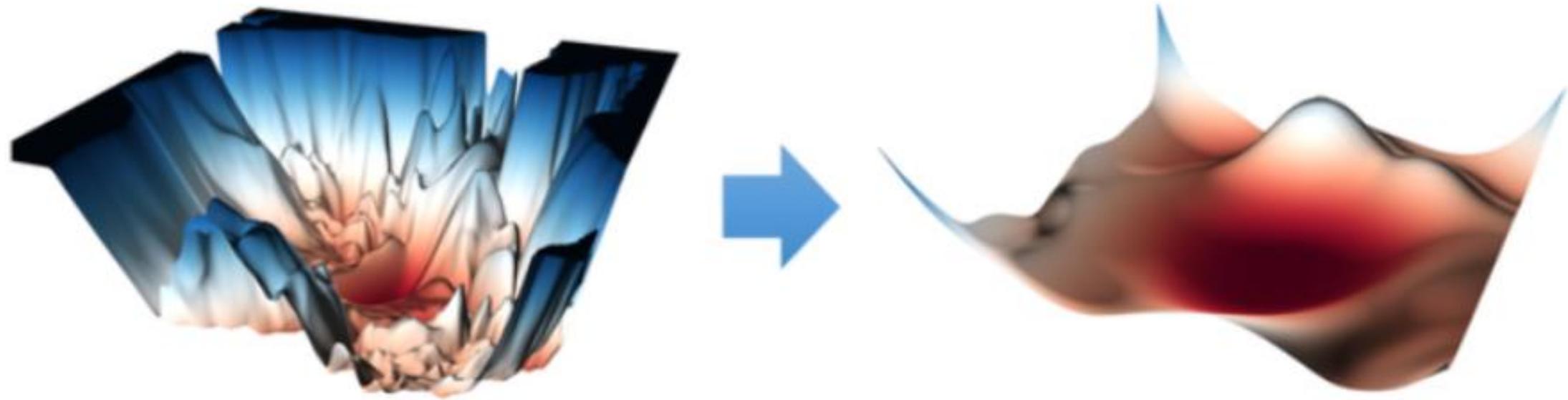
Batch-Normalization in CNN



Does Batch Normalization help to solve Internal Covariance Shifting?



Does Batch Normalization help to solve Internal Covariance Shifting?



Layer 배치 순서

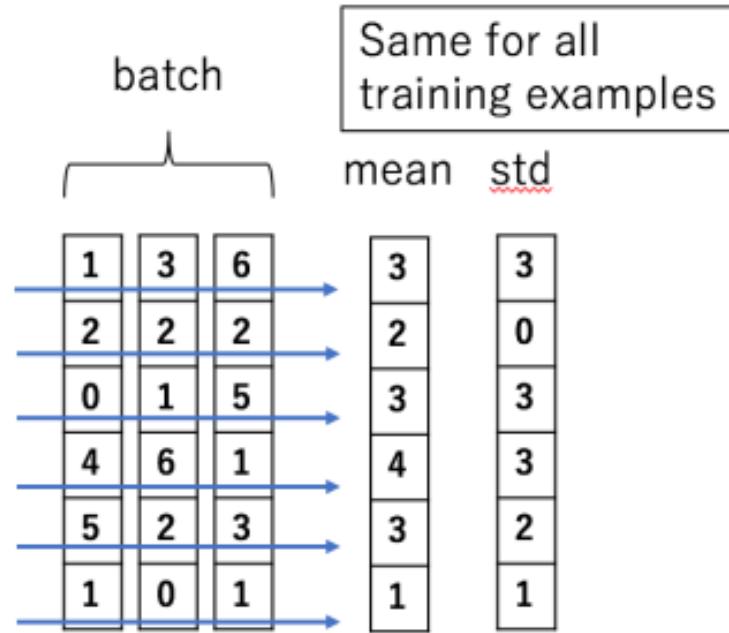
1. Convolution layer
2. Batch normalization
3. Activation
4. Dropout
5. Maxpooling

1. Dense layer
2. Batch normalization
3. Activation(+Layer normalization)
4. Dropout

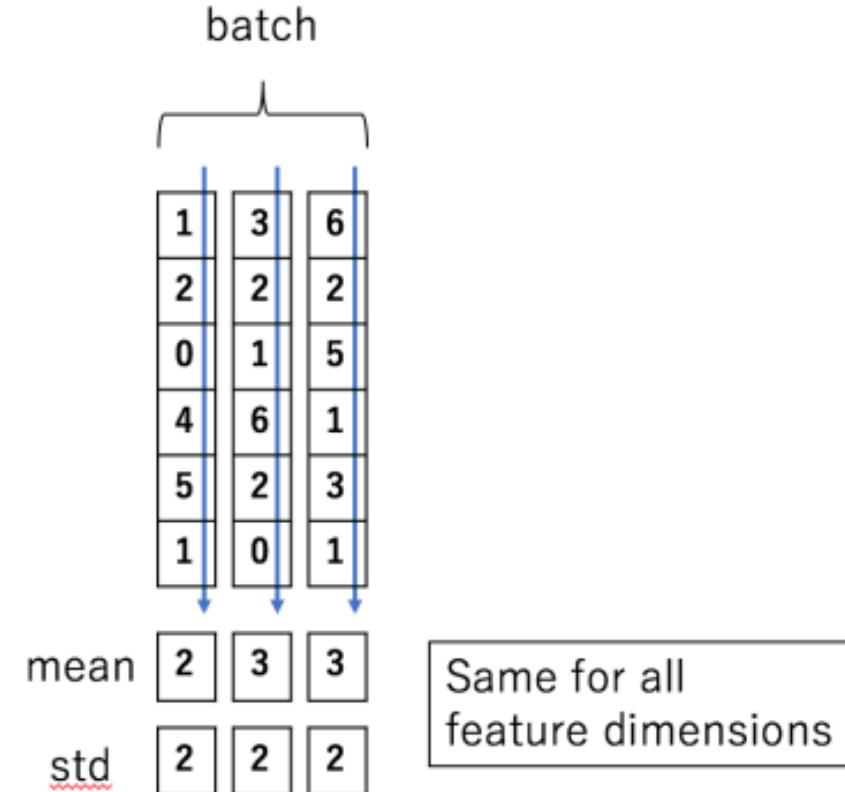
물론, 저자들마다 하는 얘기가 다름!!

Layer Normalization / Batch-Renormalization

Batch Normalization



Layer Normalization





Q & A

감사합니다