# Regularization for Deep Learning

ADVANCED ARTIFICIAL INTELLIGENCE

JUCHEOL MOON

# Regularization

- Regularization is any modification we make to a learning algorithm that is intended to reduce
  - its __generalization (test)__ error
  - but not its __training__ error
- Regularization strategies
  - put extra constraints on a machine learning model
  - add extra terms in the objective function

# Parameter Norm Penalties

- Limiting the capacity of models by adding a parameter norm penalty $\Omega(\theta)$ to the objective function $J$.
- $\tilde{J}(\vec{\theta}; \vec{X}, \vec{y}) = J(\vec{\theta}; \vec{X}, \vec{y}) + \alpha \Omega(\vec{\theta})$

- where $\alpha \in [0, \infty)$ is a <u>hyper-parameter</u>.
- $\alpha = 0$ results in no regularization,
- Large $\alpha$ correspond to more regularization.
- It is sometimes desirable to use a separate penalty with a different $\alpha$ coefficient for each layer of the network.

3

3

# $L^2$ Parameter Regularization

- Also known as <u>weight decay</u> and <u>ridge</u> regularization
- $\tilde{J}(\vec{w}; \vec{X}, \vec{y}) = J(\vec{w}; \vec{X}, \vec{y}) + \frac{\alpha}{2} \vec{w}^T \vec{w}$

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

- $\vec{w} \leftarrow \vec{w} - \varepsilon \left( \alpha \vec{w} + \nabla_{\vec{w}} J(\vec{w}; \vec{X}, \vec{y}) \right)$

- $\vec{w} \leftarrow \vec{w} - \varepsilon \alpha \vec{w} - \varepsilon \nabla_{\vec{w}} J(\vec{w}; \vec{X}, \vec{y}) = (1 - \varepsilon\alpha)\vec{w} - \varepsilon \nabla_{\vec{w}} J(\vec{w}; \vec{X}, \vec{y})$

- shrink the weight vector by a constant factor on each step
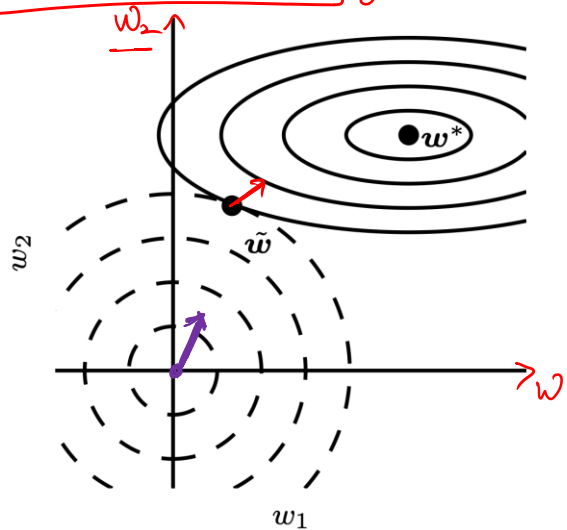
4

4

2

# Weight Decay as Constrained Optimization

■$\vec{w} \leftarrow \boxed{(1 - \epsilon\alpha)\vec{w}} \boxed{- \epsilon\nabla_{\overrightarrow{w}}J(\vec{w}; \vec{X}, \vec{y})}$ $\vec{g}$

$0 < \Sigma \ll 1$

$0 < \alpha \ll 1$
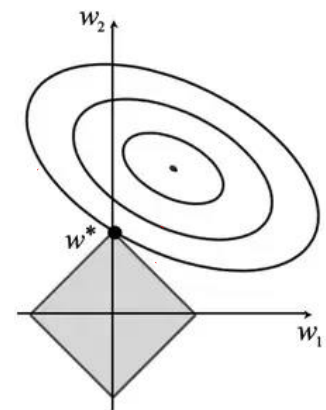
$0 < \Sigma\alpha \ll 1$

$1 - \Sigma\alpha < 1$

# $L^1$ Parameter Regularization

■Also known as LASSO regularization

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

■$\nabla_{\overrightarrow{w}}\tilde{J}(\vec{w}; \vec{X}, \vec{y}) = D_{\overrightarrow{w}}J(\vec{w}; \vec{X}, \vec{y}) + \alpha\, sign(\vec{w})$

$\vec{g}$

■$\vec{w} \leftarrow \vec{w} \underline{- \epsilon\alpha\, sign(\vec{w})} \boxed{- \epsilon\nabla_{\overrightarrow{w}}J(w; \vec{X}, \vec{y})}$

# $L^2$ Regularization in Neural Network

- $\tilde{J}(\overrightarrow{W}, \vec{b}) = \frac{1}{m}\sum_i^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\alpha}{2}\sum_l^L \|\overrightarrow{W}^{(l)}\|_F^2$

  $\underbrace{\phantom{\frac{1}{m}\sum_i^m L(\hat{y}^{(i)}, y^{(i)})}}_{J} \qquad \text{layer} \quad \underbrace{\phantom{\frac{\alpha}{2}\sum}}_{\Omega}$

- Frobenius norm

  - $\|\overrightarrow{W}^{(l)}\|_F^2 = \sum_i \sum_j \left(w_{ij}^{(l)}\right)^2$

$\ell=1 \qquad W_{ij}^{(2)} \quad \ell=2 \qquad\qquad \ell=3$

$W_{jk}^{(3)}$

# $L^2$ Regularization in Neural Network

- $J(\overrightarrow{W}, \vec{b}) = \frac{1}{m}\sum_i^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\alpha}{2}\sum_l^L \|\overrightarrow{W}^{(l)}\|_F^2$

- In deep neural network
  - $\alpha \uparrow$: underfitting, $\alpha \downarrow$: overfitting

$w_2$

$\vec{w}^{\circ}$

$(1-\varepsilon\alpha)\vec{w}$

$\vec{w} \leftarrow (1-\varepsilon\alpha)\vec{w} - \varepsilon D_{\vec{w}}J$

$w_1$
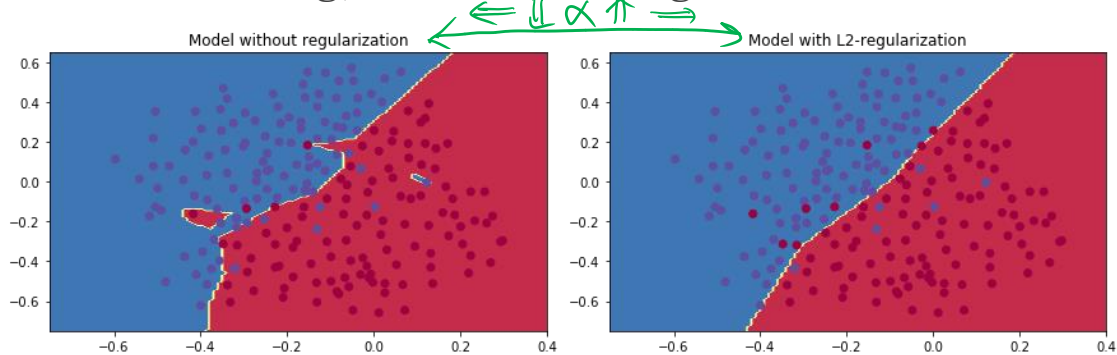
underfitted boundary

overfitted boundary

$\sigma$

$\omega$

# $L^2$ Regularization in Neural Network

- $J(\vec{W}, \vec{b}) = \frac{1}{m} \sum_i^m L(\hat{y}^{(i)}, y^{(i)}) + \frac{\alpha}{2} \sum_l^L \left\| \vec{W}^{(l)} \right\|_F^2$

- In deep neural network
  - $\alpha \uparrow$: underfitting, $\alpha \downarrow$: overfitting

$$\Longleftarrow \Downarrow \alpha \Uparrow \Longrightarrow$$
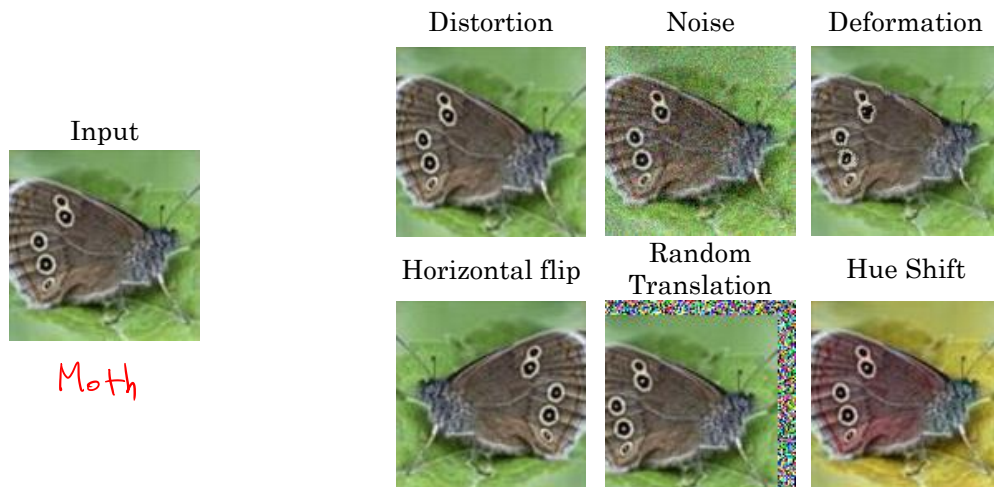


Model without regularization

Model with L2-regularization

9

# Dataset Augmentation

- The best way to make a machine learning model generalize better is to train it on more data.



Distortion    Noise    Deformation

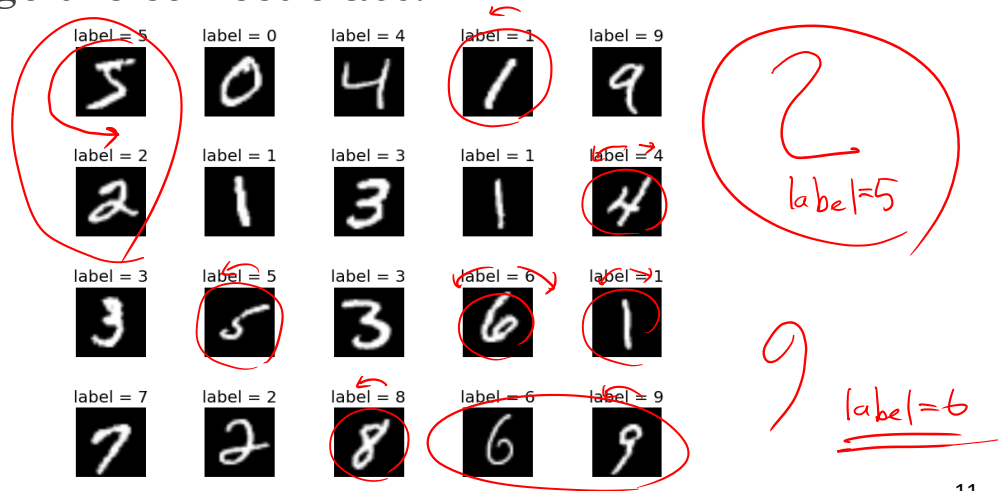Input

Moth

Horizontal flip    Random Translation    Hue Shift

10

# Dataset Augmentation

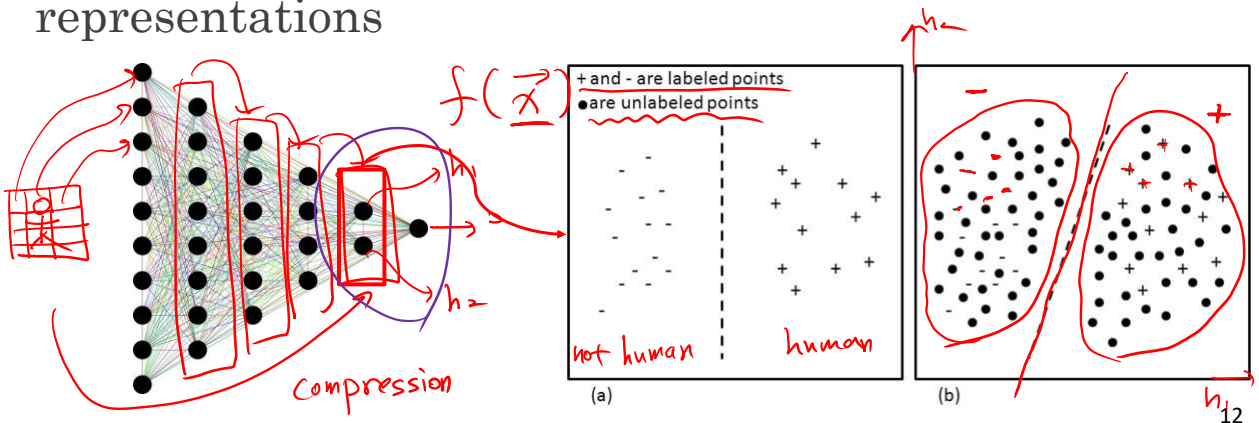- One must be careful not to apply transformations that would change the correct class.

# Semi-Supervised Learning

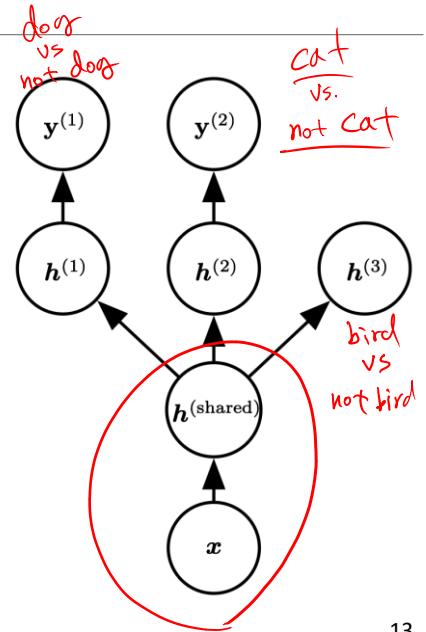- Semi-supervised learning usually refers to learning a representation $h = f(x)$
  - Examples from the same class have similar representations

# Multi-Task Learning

▪The model can generally be divided into two kinds of parts and associated parameters:

  ▪Task-specific parameters (which only benefit from the examples of their task to achieve good generalization)

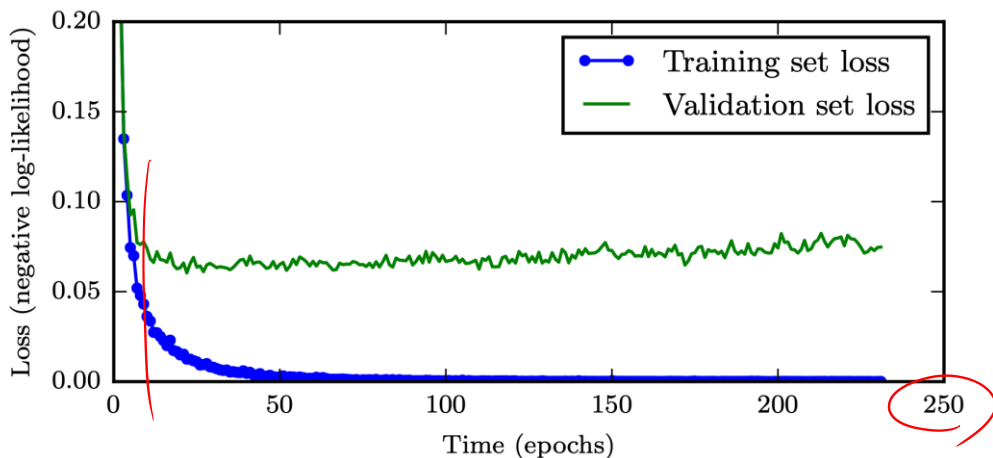  ▪Generic parameters, shared across all the tasks (which benefit from the pooled data of all the tasks).



13

# Early Stopping

▪We often observe that training error decreases steadily over time, but validation set error begins to rise again.



14

# Early Stopping version 1

- Every time the error on the validation set improves, we store a copy of the model parameters.

- When the training algorithm terminates, we return these parameters, rather than the latest parameters.

- An additional cost to early stopping is the need to maintain a copy of the best parameters.
  - This cost is generally negligible.

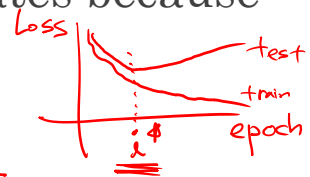- Early stopping requires a validation set, which means some training data is not fed to the model.

# Early Stopping version 2

- Perform extra training after the initial training with early stopping has completed.

- In this second training pass, we train for the same number of steps as the early stopping procedure determined was optimal in the first pass.

- On the second round of training, each pass through the dataset will require more parameter updates because the training set is bigger.

# Early Stopping version 2

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

Let $\boldsymbol{X}^{(\text{train})}$ and $\boldsymbol{y}^{(\text{train})}$ be the training set.
Split $\boldsymbol{X}^{(\text{train})}$ and $\boldsymbol{y}^{(\text{train})}$ into $(\boldsymbol{X}^{(\text{subtrain})}, \boldsymbol{X}^{(\text{valid})})$ and $(\boldsymbol{y}^{(\text{subtrain})}, \boldsymbol{y}^{(\text{valid})})$ respectively.
Run early stopping (algorithm 7.1) starting from random $\boldsymbol{\theta}$ using $\boldsymbol{X}^{(\text{subtrain})}$ and $\boldsymbol{y}^{(\text{subtrain})}$ for training data and $\boldsymbol{X}^{(\text{valid})}$ and $\boldsymbol{y}^{(\text{valid})}$ for validation data. This returns $i^*$, the optimal number of steps.
Set $\boldsymbol{\theta}$ to random values again.
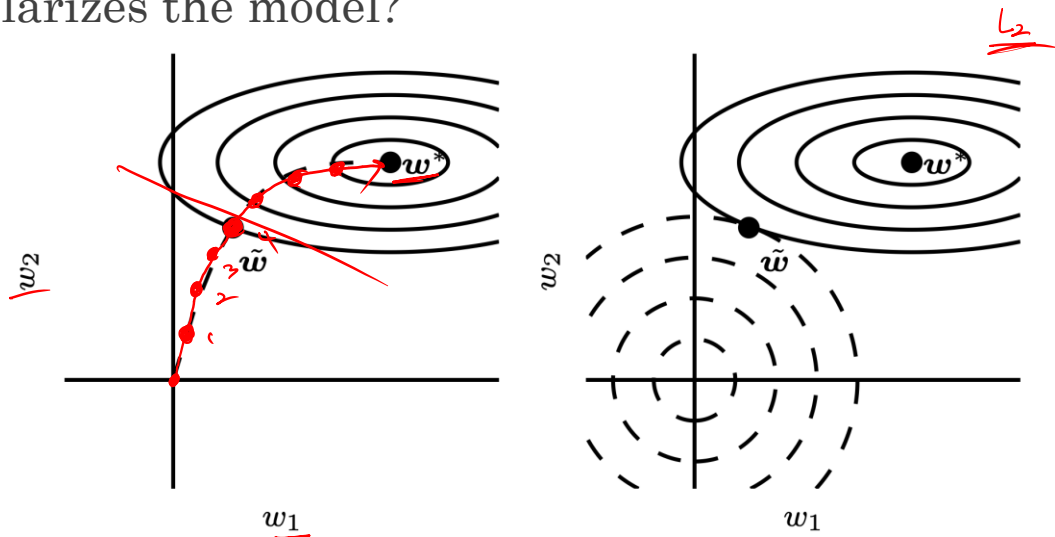Train on $\boldsymbol{X}^{(\text{train})}$ and $\boldsymbol{y}^{(\text{train})}$ for $i^*$ steps.

# Early Stopping

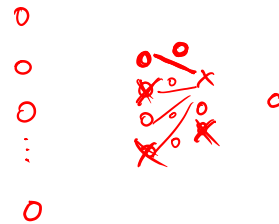- What is the actual mechanism by which early stopping regularizes the model?

# Sparse Representations

- Meaning of $L^1$ (or $L^2$) regularization?
  - Many of the parameters become zero (or close to zero)
- $\hat{y} = f(\underline{h}) = f(f(x))$

- Representational sparsity describes a representation where many of the elements of the representation are zero (or close to zero)
- $\tilde{J}(\theta; X, y) = J + \alpha \Omega(\underline{h})$
  - For example, $\Omega(h) = \|h\|_1$

# Ensemble Methods

- Consider for example a set of $k$ regression models.
  - Suppose that each model makes an error $\epsilon_i$ on each example
  - with variances $\mathbb{E}[\epsilon_i^2] = \underline{v}$, covariances $\mathbb{E}[\epsilon_i \epsilon_j] = c$
  - The error made by the average prediction: $\frac{1}{k}\sum_i \epsilon_i$

- The expected squared error
  - $\mathbb{E}\left[\left(\frac{1}{k}\sum_i \epsilon_i\right)^2\right] = \frac{1}{k^2}\mathbb{E}\left[\left(\sum_i \epsilon_i\right)^2\right] = \frac{1}{k^2}\mathbb{E}\left[\sum_i\left(\epsilon_i^2 + \sum_{i\neq j}\epsilon_i\epsilon_j\right)\right]$

  $= \frac{1}{k^2}\left(\sum_i^k \mathbb{E}[\epsilon_i^2] + \sum_i\sum_{i\neq j}\mathbb{E}[\epsilon_i\epsilon_j]\right) = \frac{1}{k}v + \frac{1}{k^2}k(k-1)c$

# Ensemble Methods

- The expected squared error
  - $\mathbb{E}\left[\left(\frac{1}{k}\sum_i \epsilon_i\right)^2\right] = \frac{1}{k}v + \frac{k-1}{k}c$
- When the errors are perfectly correlated: $c = v$
  - $\mathbb{E}\left[\left(\frac{1}{k}\sum_i \epsilon_i\right)^2\right] = \frac{1}{k}\left(2v + (k-1)v\right) = v$
- When the errors are perfectly uncorrelated: $c = 0$
  - $\mathbb{E}\left[\left(\frac{1}{k}\sum_i \epsilon_i\right)^2\right] = \frac{v}{k}$
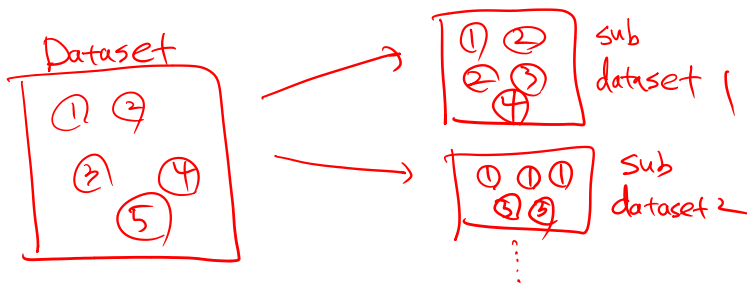- The expected squared error of the ensemble decreases linearly with the ensemble size.

$\frac{v}{k} \longrightarrow v \longrightarrow$ error

# Bagging

- Bagging involves constructing $k$ different datasets.
  - Each dataset has the same number of examples as the original dataset
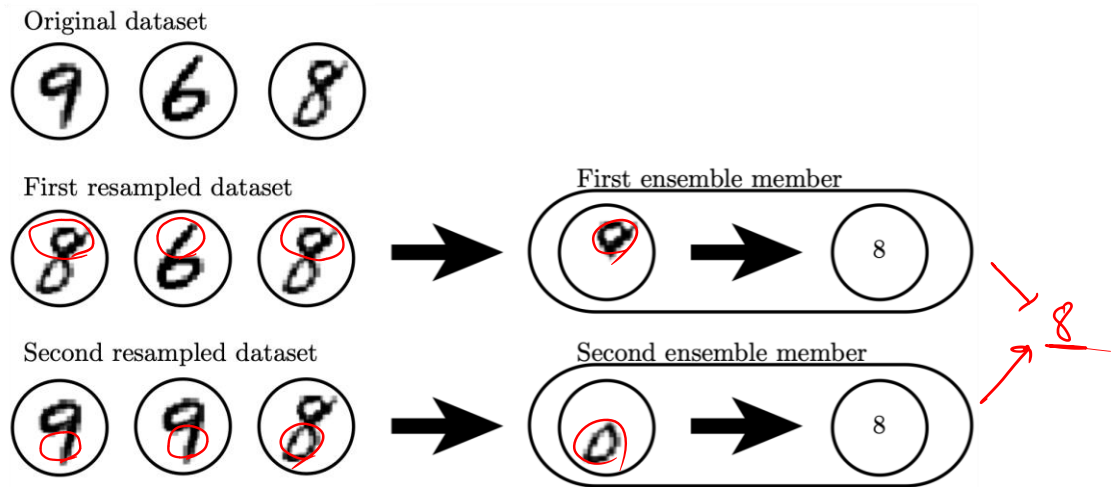  - but each dataset is constructed by sampling with replacement from the original dataset



Dataset → sub dataset 1, sub dataset 2 ...

# Bagging

- How bagging works

Original dataset

First resampled dataset | First ensemble member

Second resampled dataset | Second ensemble member

# Ensemble in Neural networks

- Neural networks reach a wide enough variety of solution points that they can often benefit from model averaging even if all of the models are trained on the **same dataset**.
  - differences in random initialization
  - random selection of minibatches,
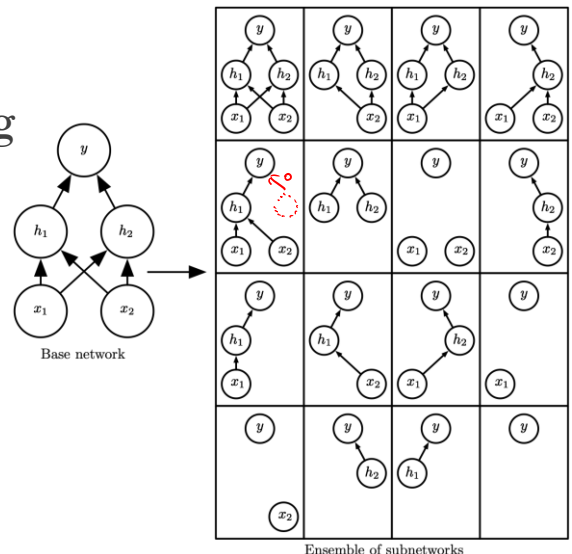  - differences in hyperparameters

# Dropout

- Bagging involves training multiple models, and evaluating multiple models on each test example.
  - impractical when each model is a large neural network, since training and evaluating such networks is costly in terms of runtime and memory

- Dropout provides an inexpensive approximation to training and evaluating a bagged ensemble of exponentially many neural networks

# Dropout

- Dropout trains the ensemble consisting of all sub-networks that can be formed by removing non-output units from an underlying base network.

- In most modern neural networks, we can effectively remove a unit from a network by multiplying its output value by zero



Base network

Ensemble of subnetworks

# Dropout

- Learn with bagging
  - we define $k$ different models
  - construct $k$ different datasets by sampling from the training set with replacement
  - then train model $i$ on dataset $i$.

- Dropout aims to approximate this process
  - Each time we load an example into a minibatch,
  - we randomly sample a different binary mask to apply to all of the input and hidden units in the network.
  - The mask for each unit is sampled independently.
  - The probability of sampling a mask value of one is a **hyperparameter** fixed before training begins

# Dropout

- The models share parameters, with each model inheriting a different subset of parameters from the parent neural network.

- The models can have different dropout probabilities for the layers.

- Predictions at test phase,
  - No dropout! ( base network )

- Significant advantage of dropout is that it does not significantly limit the type of model or training procedure that can be used.