

COSI 165B

Deep Learning

Chuxu Zhang
Computer Science Department
Brandeis University

2/22/2021

□ Backpropagation

-
- 1: Compute the values of $z \in \mathbb{R}^m$, $a \in \mathbb{R}^m$, and $o \in \mathbb{R}$
 - 2: Compute

$$\delta^{[2]} \triangleq \frac{\partial J}{\partial o} = (o - y) \in \mathbb{R}$$

$$\delta^{[1]} \triangleq \frac{\partial J}{\partial z} = (W^{[2]\top} (o - y)) \odot 1\{z \geq 0\} \in \mathbb{R}^{m \times 1}$$

(by eqn. (3.12) and (3.13))

- 3: Compute

$$\frac{\partial J}{\partial W^{[2]}} = \delta^{[2]} a^\top \in \mathbb{R}^{1 \times m} \quad (\text{by eqn. (3.5)})$$

$$\frac{\partial J}{\partial b^{[2]}} = \delta^{[2]} \in \mathbb{R} \quad (\text{by eqn. (3.6)})$$

$$\frac{\partial J}{\partial W^{[1]}} = \delta^{[1]} x^\top \in \mathbb{R}^{m \times d} \quad (\text{by eqn. (3.7)})$$

$$\frac{\partial J}{\partial b^{[1]}} = \delta^{[1]} \in \mathbb{R}^m \quad (\text{as an exercise})$$

-
- 1: Compute and store the values of $a^{[k]}$'s and $z^{[k]}$'s for $k = 1, \dots, r$, and J .
 ▷ This is often called the “forward pass”
 - 2: .
 - 3: **for** $k = r$ to 1 **do** ▷ This is often called the “backward pass”
 - 4: **if** $k = r$ **then**
 - 5: compute $\delta^{[r]} \triangleq \frac{\partial J}{\partial z^{[r]}}$
 - 6: **else**
 - 7: compute

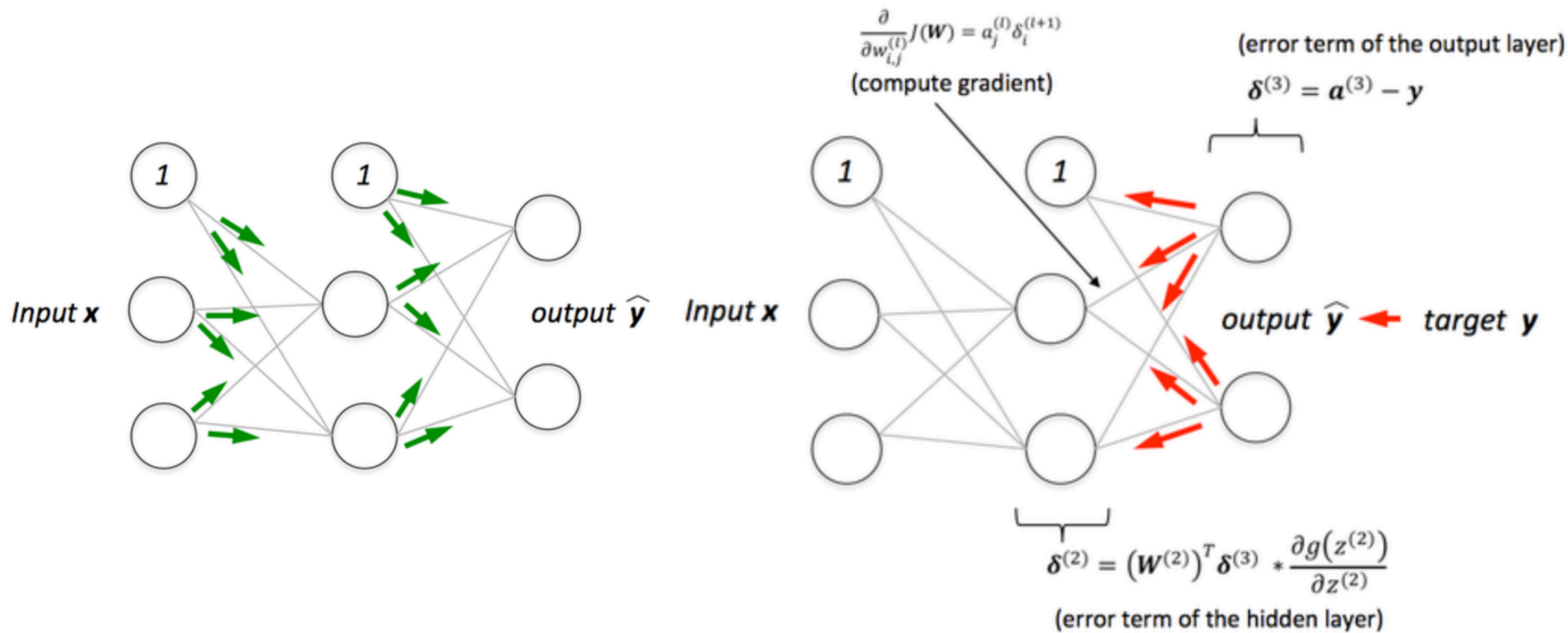
$$\delta^{[k]} \triangleq \frac{\partial J}{\partial z^{[k]}} = \left(W^{[k+1]\top} \delta^{[k+1]} \right) \odot \text{ReLU}'(z^{[k]})$$

- 8: Compute

$$\frac{\partial J}{\partial W^{[k]}} = \delta^{[k]} a^{[k-1]\top}$$

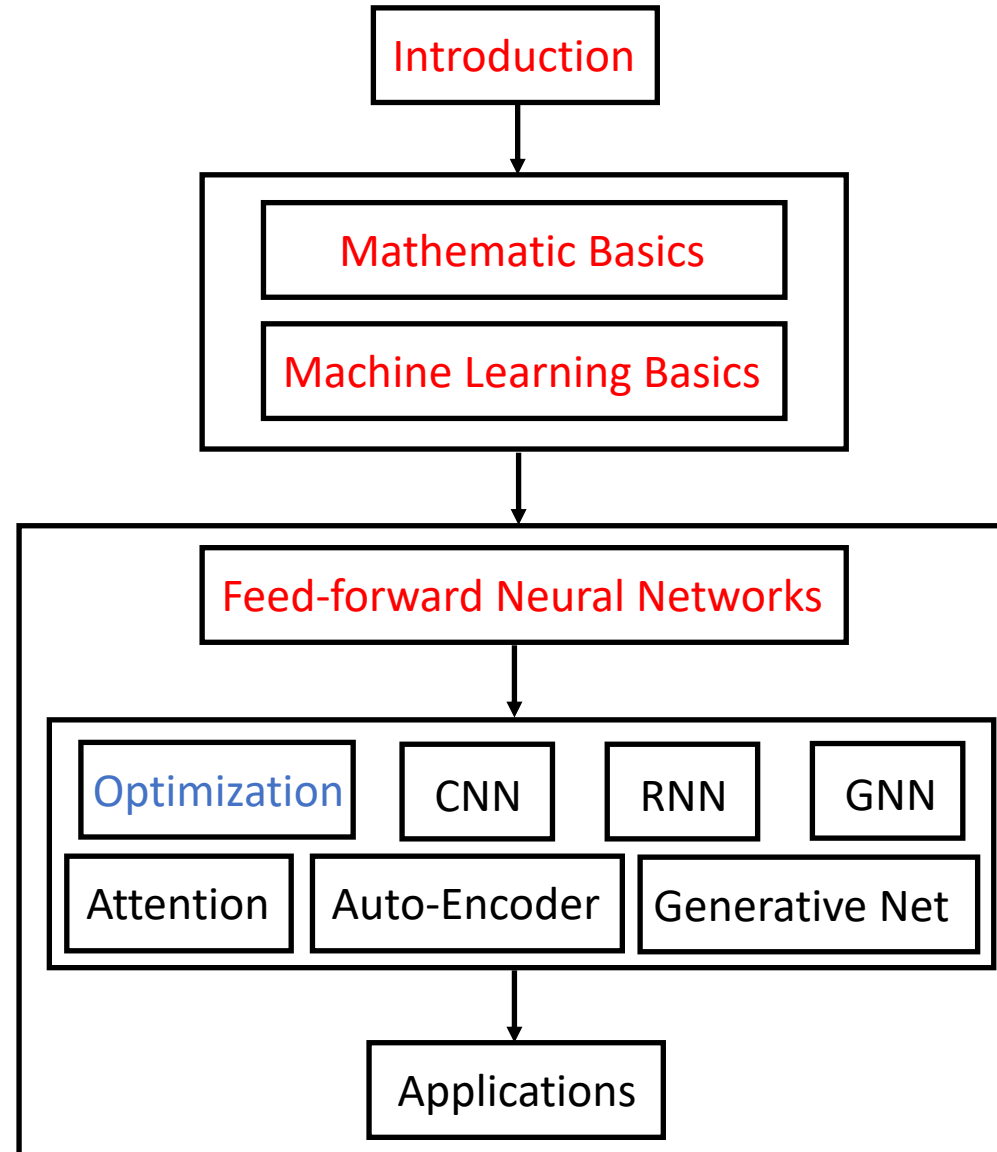
$$\frac{\partial J}{\partial b^{[k]}} = \delta^{[k]}$$

□ Backpropagation



In backpropagation, we “simply” backpropagate the error to update model parameters.

Structure of This Course



❑ Parameter Norm Penalties

$$\tilde{J}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = J(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \alpha \Omega(\boldsymbol{\theta}) \quad \text{ridge regression}$$

❑ L2 Regularization

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

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

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

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

Single step: multiplicatively shrink the weight vector by a constant factor on each step, just before performing the usual gradient update

□ L2 Regularization

Whole process: analysis by quadratic approximation

$$\hat{J}(\boldsymbol{\theta}) = J(\boldsymbol{w}^*) + \frac{1}{2}(\boldsymbol{w} - \boldsymbol{w}^*)^\top \boldsymbol{H}(\boldsymbol{w} - \boldsymbol{w}^*)$$

Hessian matrix

set gradient to zero

$$\nabla_{\boldsymbol{w}} \hat{J}(\boldsymbol{w}) = \boldsymbol{H}(\boldsymbol{w} - \boldsymbol{w}^*)$$

add weight
decay gradient

$$\begin{aligned}\alpha \tilde{\boldsymbol{w}} + \boldsymbol{H}(\tilde{\boldsymbol{w}} - \boldsymbol{w}^*) &= 0 \\ (\boldsymbol{H} + \alpha \boldsymbol{I}) \tilde{\boldsymbol{w}} &= \boldsymbol{H} \boldsymbol{w}^* \\ \tilde{\boldsymbol{w}} &= (\boldsymbol{H} + \alpha \boldsymbol{I})^{-1} \boldsymbol{H} \boldsymbol{w}^*.\end{aligned}$$

decomposition

$$\boldsymbol{H} = \underset{\text{lambda}}{\boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^\top}.$$

$$\begin{aligned}\tilde{\boldsymbol{w}} &= (\boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^\top + \alpha \boldsymbol{I})^{-1} \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^\top \boldsymbol{w}^* \\ &= \left[\boldsymbol{Q} (\boldsymbol{\Lambda} + \alpha \boldsymbol{I}) \boldsymbol{Q}^\top \right]^{-1} \boldsymbol{Q} \boldsymbol{\Lambda} \boldsymbol{Q}^\top \boldsymbol{w}^* \\ &= \boldsymbol{Q} (\boldsymbol{\Lambda} + \alpha \boldsymbol{I})^{-1} \boldsymbol{\Lambda} \boldsymbol{Q}^\top \boldsymbol{w}^*.\end{aligned}$$

The effect of weight decay is to rescale \boldsymbol{w}^* along the axes defined by the eigenvectors

Along the directions where the eigenvalues of \boldsymbol{H} are relatively large, the effect of regularization is relatively small

□ L1 Regularization

$$\Omega(\boldsymbol{\theta}) = \|\boldsymbol{w}\|_1 = \sum_i |w_i|, \quad \text{sum of absolute values}$$

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

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

$$\hat{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = J(\boldsymbol{w}^*; \boldsymbol{X}, \boldsymbol{y}) + \sum_i \left[\frac{1}{2} H_{i,i} (\boldsymbol{w}_i - \boldsymbol{w}_i^*)^2 + \alpha |\boldsymbol{w}_i| \right] \quad \text{quadratic approximation}$$

analytical
solution

$$w_i = \text{sign}(w_i^*) \max \left\{ |w_i^*| - \frac{\alpha}{H_{i,i}}, 0 \right\}$$

two cases

$$w_i^* \leq \frac{\alpha}{H_{i,i}} \quad \quad w_i^* > \frac{\alpha}{H_{i,i}}$$

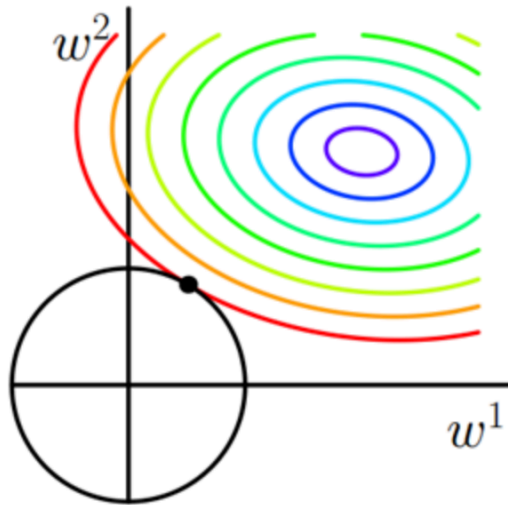
zero shift \boldsymbol{w} by a distance

□ L2 and L1 Comparison (Regularization as Constraint)

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

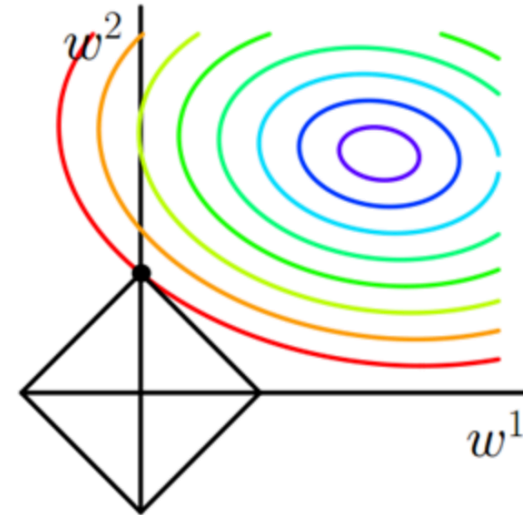
$$\tilde{J}(\mathbf{w}; \mathbf{X}, \mathbf{y}) = \alpha \|\mathbf{w}\|_1 + J(\mathbf{w}; \mathbf{X}, \mathbf{y})$$

constraint
optimization



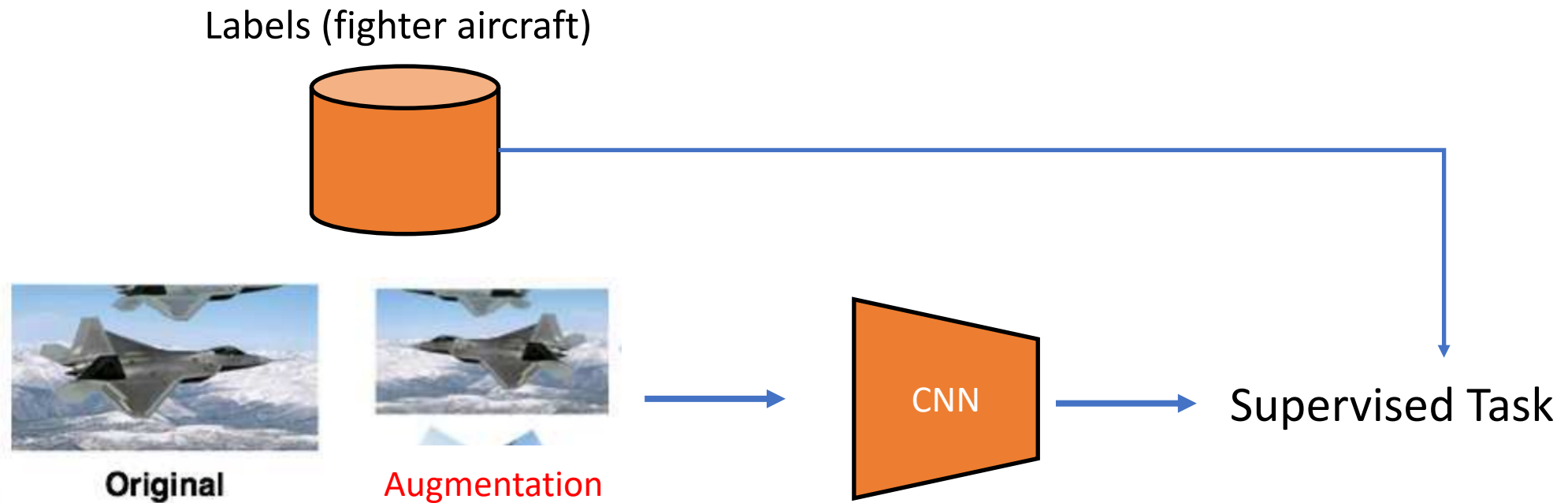
non-sparse weights
(most are non-zero)

tangent points

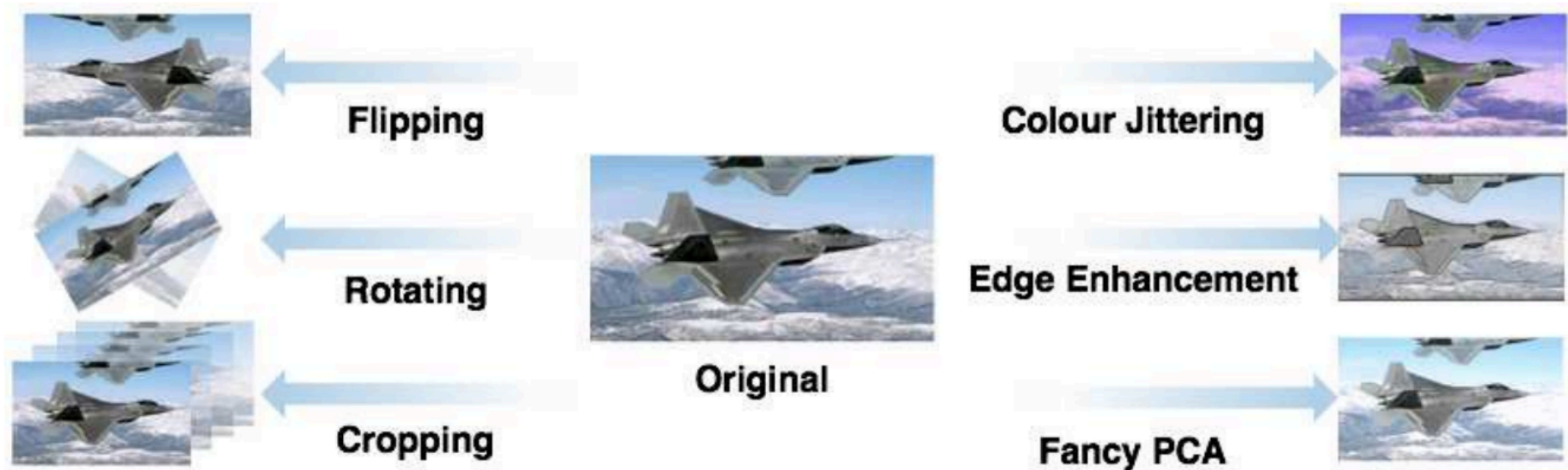


sparse weights
(most are zero)

More Data, Better Generalization

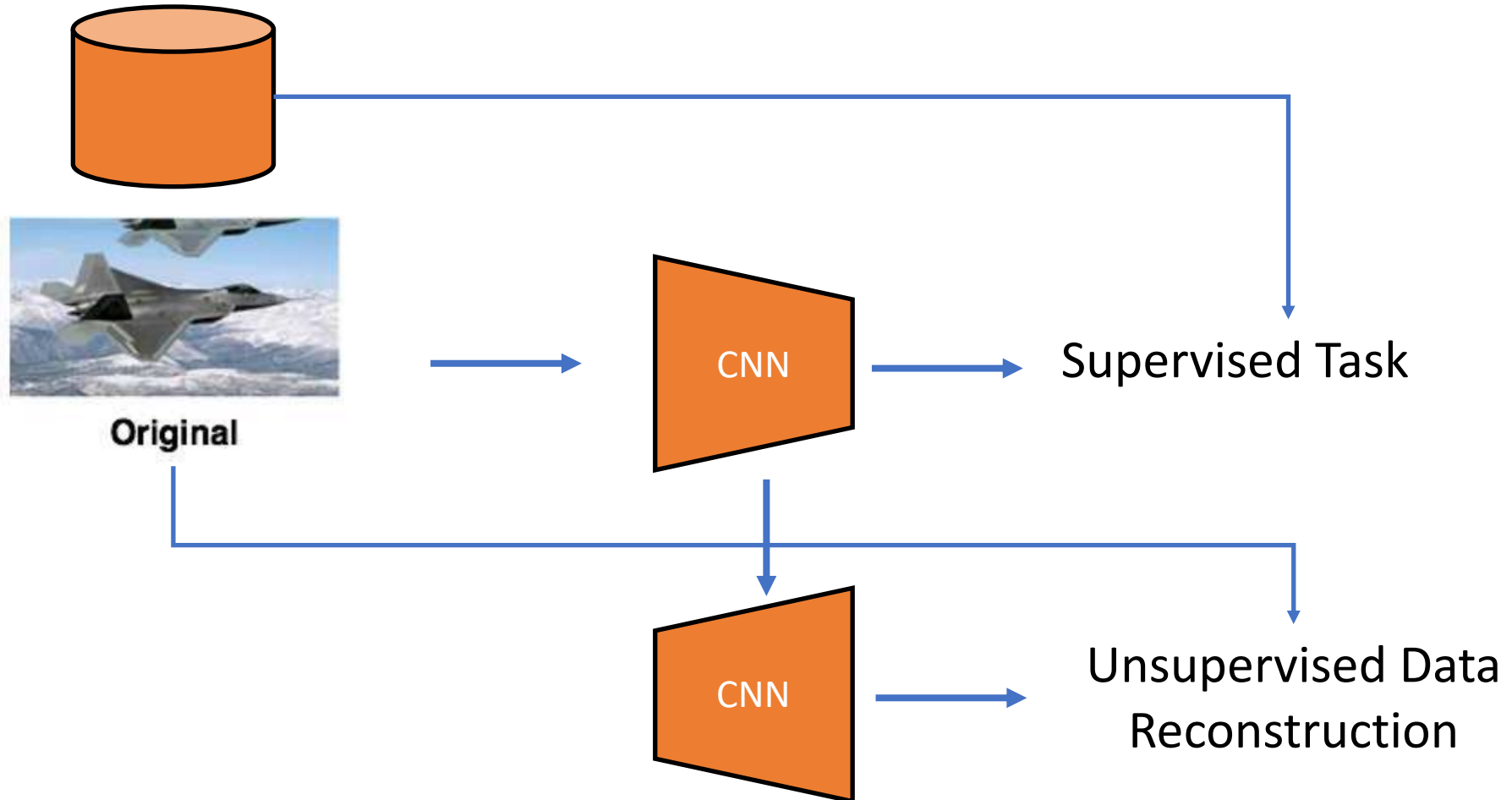


□ Example (Image)



□ Supervised Task + Unsupervised Data Reconstruction

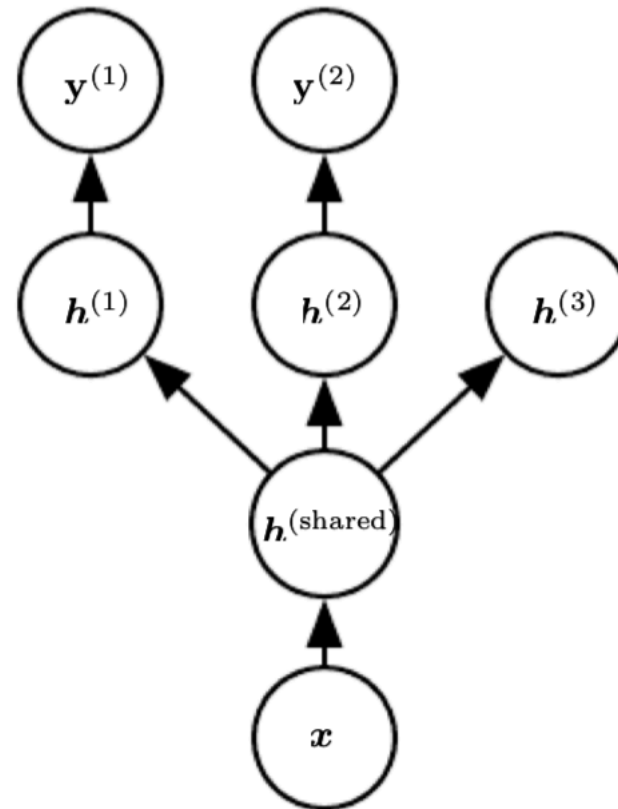
Labels (fighter aircraft)



□ Multiple Tasks

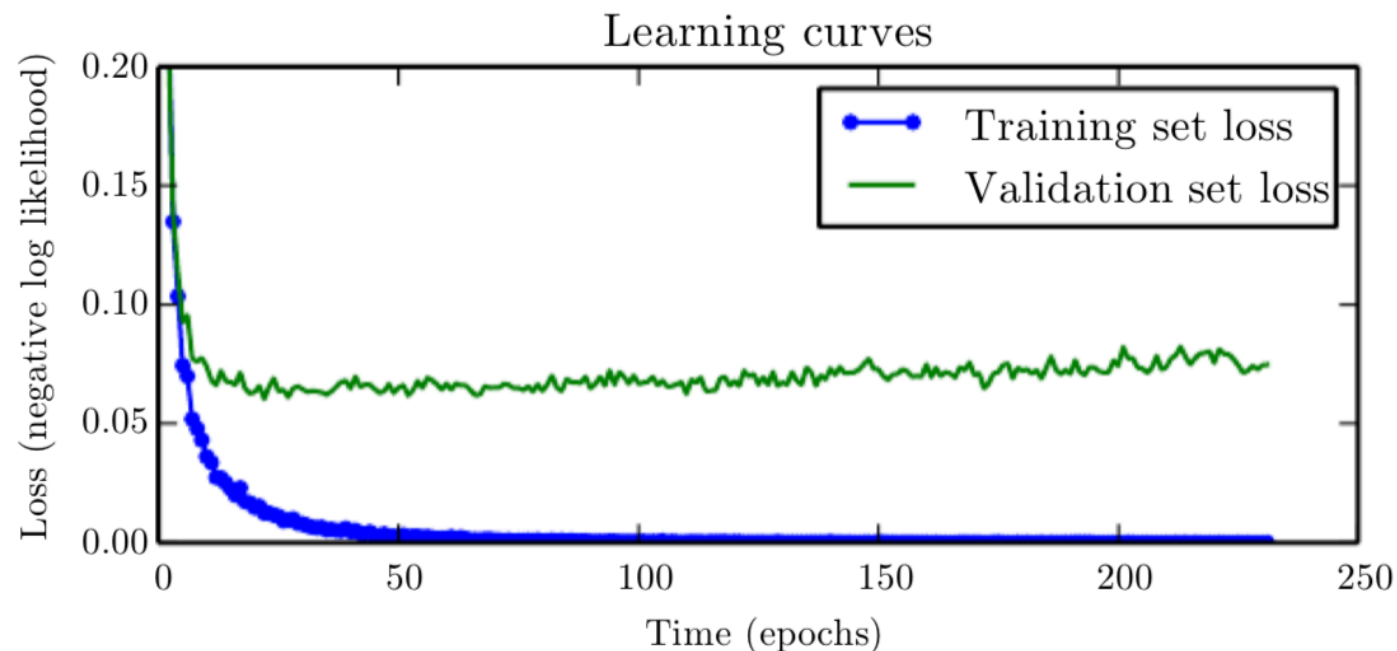
Different Tasks

Task-specific
parameters



Additional training examples put more pressure on the parameters of the model towards values that generalize well, when part of a model is shared across tasks, that part of the model is more constrained towards good values, often yielding better generalization.

□ Learning Curve



Training error decreases steadily over time, but validation set error begins to rise again

Instead of running model until it reaches minimum validation error, run it until the error on validation set has not improved for some amount of time.

□ Algorithm

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.

$\theta \leftarrow \theta_o$

$i \leftarrow 0$

$j \leftarrow 0$

$v \leftarrow \infty$

$\theta^* \leftarrow \theta$

$i^* \leftarrow i$

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^*

□ Algorithm

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

Algorithm 7.3 Meta-algorithm using early stopping to determine at what objective value we start to overfit, then continue training until that value is reached.

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 θ 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 updates θ .
 $\epsilon \leftarrow J(\theta, \mathbf{X}^{(\text{subtrain})}, \mathbf{y}^{(\text{subtrain})})$
while $J(\theta, \mathbf{X}^{(\text{valid})}, \mathbf{y}^{(\text{valid})}) > \epsilon$ **do**
 Train on $\mathbf{X}^{(\text{train})}$ and $\mathbf{y}^{(\text{train})}$ for n steps.
end while

□ How It Acts As Regularizer

$$\hat{J}(\theta) = J(w^*) + \frac{1}{2}(w - w^*)^\top H(w - w^*)$$

quadratic approximation

$$\nabla_w \hat{J}(w) = H(w - w^*)$$

$$w^{(\tau)} = w^{(\tau-1)} - \epsilon \nabla_w J(w^{(\tau-1)})$$

$$= w^{(\tau-1)} - \epsilon H(w^{(\tau-1)} - w^*)$$

$$w^{(\tau)} - w^* = (I - \epsilon H)(w^{(\tau-1)} - w^*)$$

$$H = Q\Lambda Q^\top. \quad w^{(\tau)} - w^* = (I - \epsilon Q\Lambda Q^\top)(w^{(\tau-1)} - w^*)$$

$$Q^\top(w^{(\tau)} - w^*) = (I - \epsilon\Lambda)Q^\top(w^{(\tau-1)} - w^*)$$

$$Q^\top w^{(\tau)} = [I - (I - \epsilon\Lambda)^\tau]Q^\top w^*$$

L2 regularization

$$Q^\top \tilde{w} = (\Lambda + \alpha I)^{-1} \Lambda Q^\top w^*$$

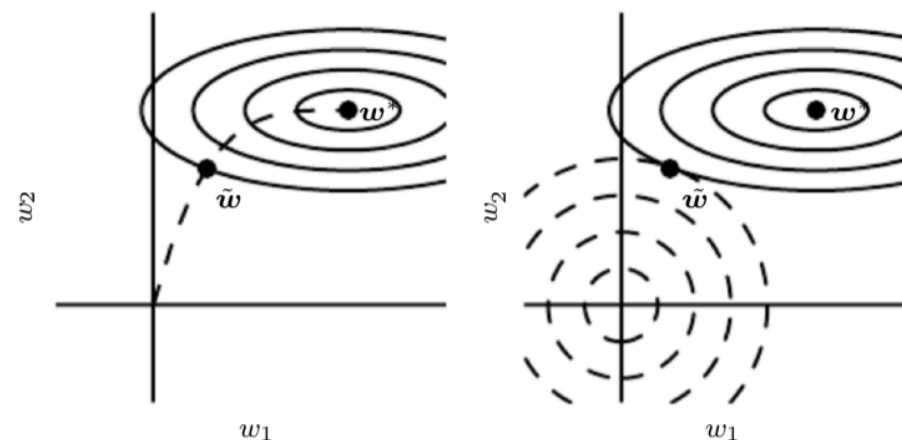
$$Q^\top \tilde{w} = [I - (\Lambda + \alpha I)^{-1} \alpha] Q^\top w^*$$

$$(I - \epsilon\Lambda)^\tau = (\Lambda + \alpha I)^{-1} \alpha$$

$$\tau \approx \frac{1}{\epsilon\alpha},$$

$$\alpha \approx \frac{1}{\tau\epsilon}.$$

L2 regularization and early stopping can be seen to be equivalent.



□ Ensemble Method

How: train several different models separately, then have all of the models vote on the output for test examples

Expected square error $\mathbb{E} \left[\left(\frac{1}{k} \sum_i \epsilon_i \right)^2 \right] = \frac{1}{k^2} \mathbb{E} \left[\sum_i \left(\epsilon_i^2 + \sum_{j \neq i} \epsilon_i \epsilon_j \right) \right]$ errors drawn from normal distribution

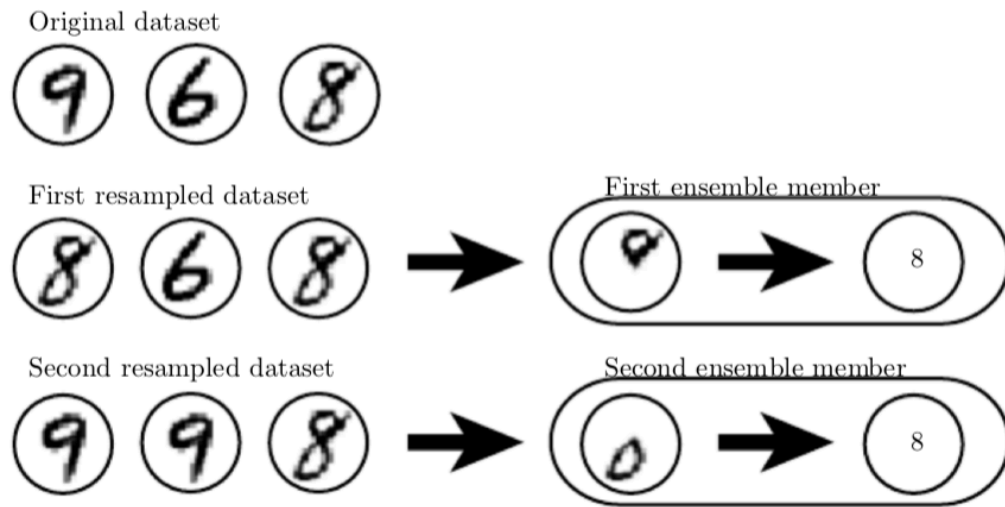
$$= \frac{1}{k} v + \frac{k-1}{k} c.$$

When errors are perfectly correlated and $c = v$, the mean squared error reduces to v , so the model averaging does not help at all.

When errors are perfectly uncorrelated and $c = 0$, the expected squared error of the ensemble is only v/k . This means that the expected squared error of the ensemble k decreases linearly with the ensemble size.

Conclusion: if the members make independent errors, the ensemble will perform significantly better than its members.

□ Example

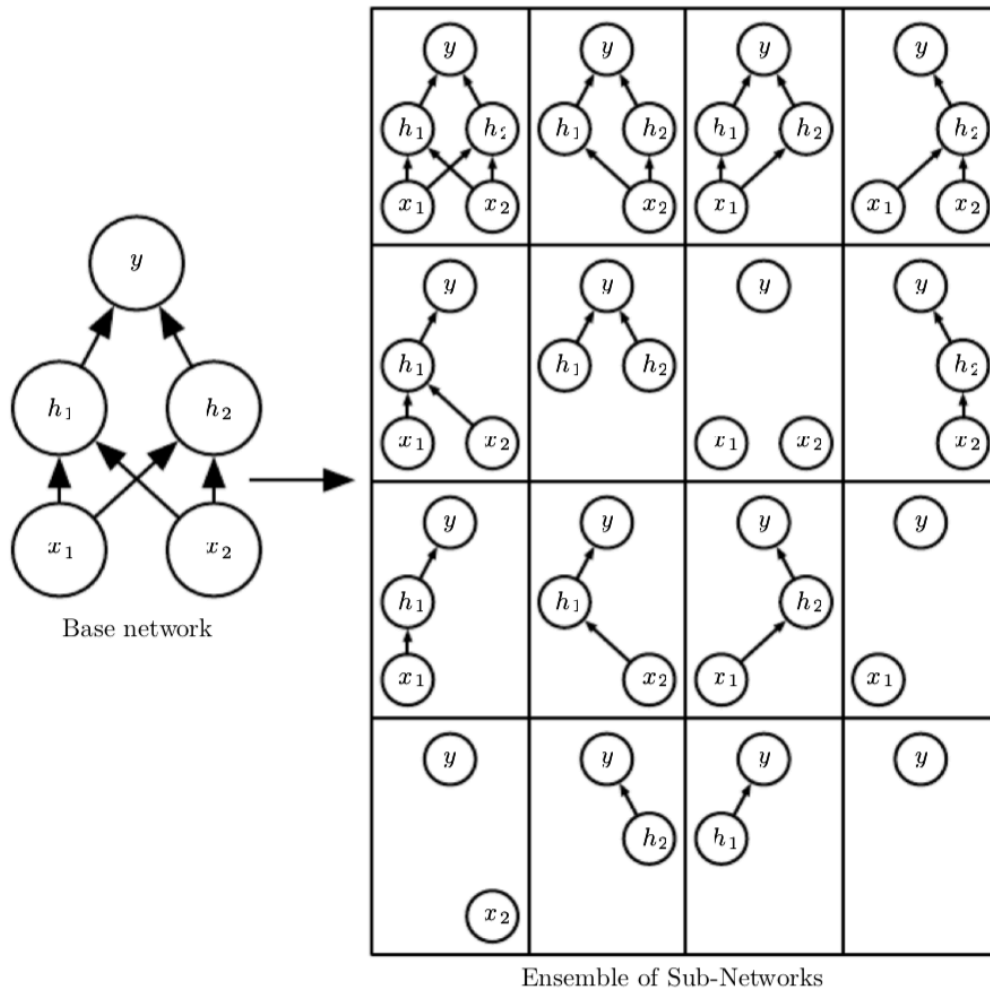


K different datasets. Each one has the same number of examples as the original dataset, but each dataset is constructed by sampling with replacement from the original dataset.

Model i is then trained on dataset i .

Missing some of the examples from the original dataset and also contains several duplicate examples.

Method



How: Randomly some portions of weights to zero
(e.g., 50%)

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

□ Strengths

Very computationally cheap. requires only $O(n)$ computation per example per update, to generate n random binary numbers and multiply them by the state. It may also require $O(n)$ memory to store these binary numbers until the back-propagation stage.

Does not significantly limit the type of model or training procedure that can be used. It works well with nearly any model that uses a distributed representation and can be trained with stochastic gradient descent.

Can work with other regularization methods.

Require sufficient training data.

□ Bagging vs Dropout

Bagging: k different models, construct k different datasets, train model i on dataset i.

Dropout: approximate this process, but with an exponentially large number of neural networks.

Bagging: the models are all independent, each model is trained to convergence on its respective training set

Dropout: the models share parameters, with each model inheriting a different subset of parameters from the parent neural network

Bagging: accumulate votes from all of its members for prediction

Dropout: each sub-model defined by mask vector u defines a probability distribution

$$\sum_{\mu} p(\mu) p(y \mid x, \mu)$$

exponential number of terms, it is intractable

approximate the inference with sampling, by averaging together the output from many masks.

Q & A