

Lecture 3: Deep Neural Networks

Pranabendu Misra
Chennai Mathematical Institute

Advanced Machine Learning 2022

(based on slides by Madhavan Mukund)

- S is not
of sweet world
PAC low

Train set S
 $M \leftarrow$ learnt
from S
Gen Loss of M
is low

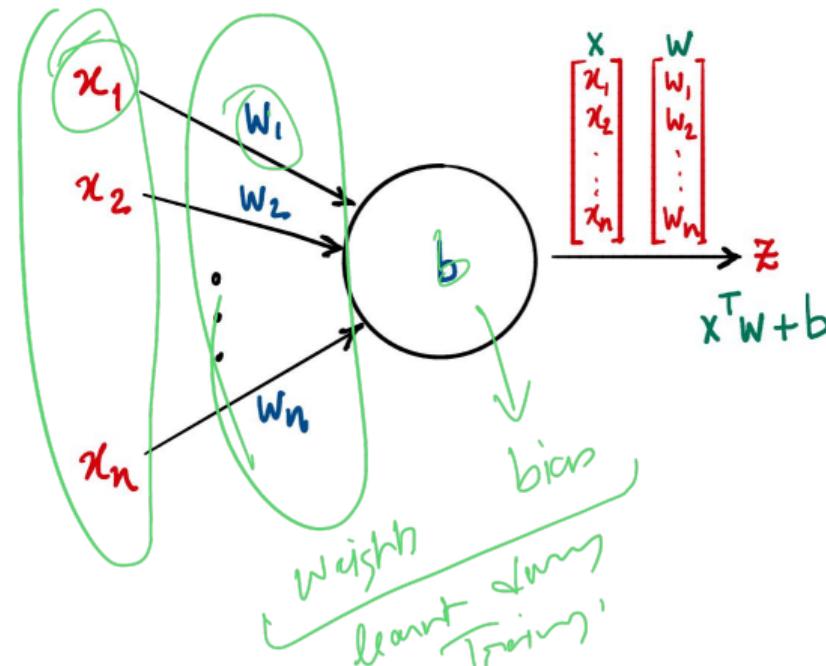
Linear separators and perceptrons

- ## ■ Perceptrons define linear separators

$$x^T w + b$$

- $x^T w + b > 0$, classify Yes (+1)
 - $x^T w + b < 0$, classify No (-1)

$$\begin{aligned} w &\in \mathbb{R}^n \\ w &= b + w \cdot x \end{aligned}$$



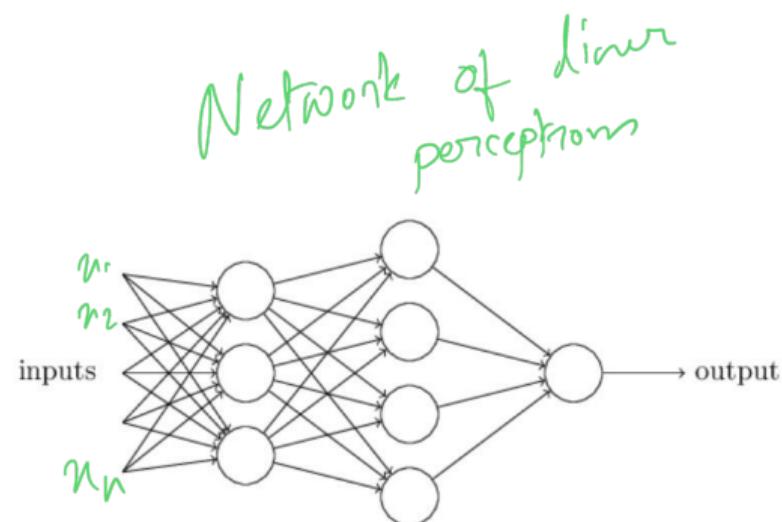
Linear separators and perceptrons

- Perceptrons define linear separators

$$x^T w + b$$

- $x^T w + b > 0$, classify Yes (+1)
- $x^T w + b < 0$, classify No (-1)

- **Unfortunately!** Network of perceptrons still defines only a linear separator



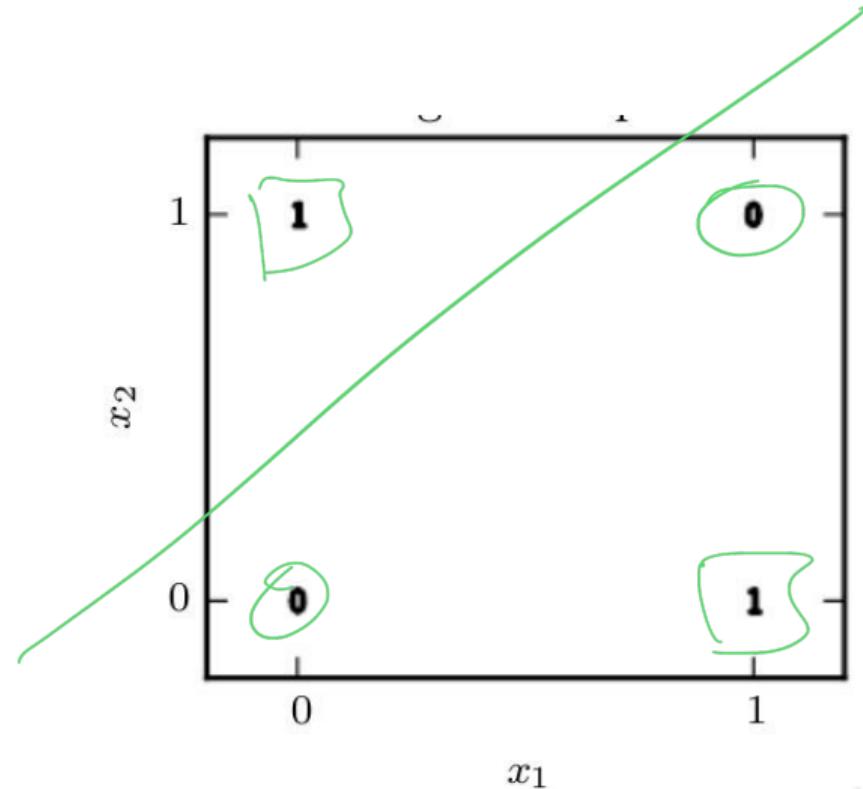
Linear separators and perceptrons

- Perceptrons define linear separators

$$x^T w + b$$

- $x^T w + b > 0$, classify Yes (+1)
- $x^T w + b < 0$, classify No (-1)

- **Unfortunately!** Network of perceptrons still defines only a linear separator
- Linear separators cannot describe XOR



Linear separators and perceptrons

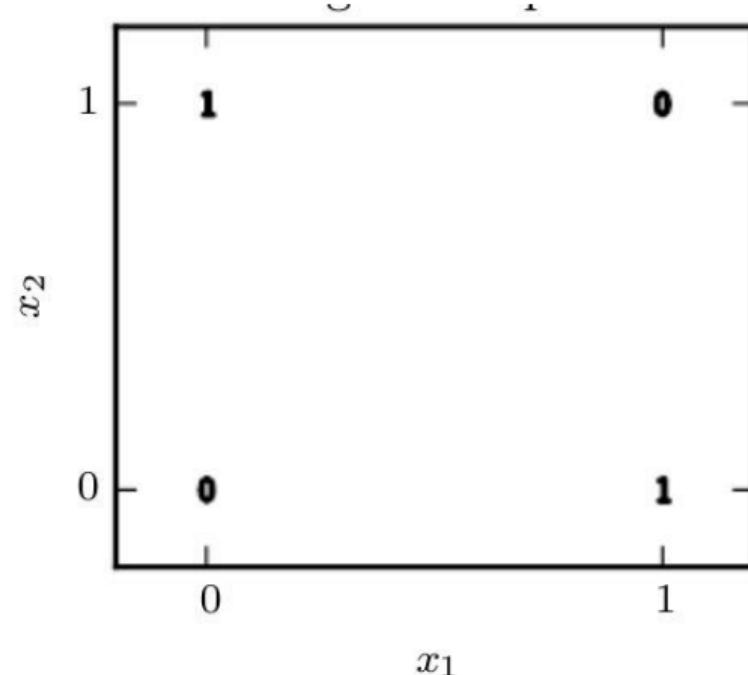
- Perceptrons define linear separators

$$x^T w + b$$

- $x^T w + b > 0$, classify Yes (+1)
- $x^T w + b < 0$, classify No (-1)

- **Unfortunately!** Network of perceptrons still defines only a linear separator
- Linear separators cannot describe XOR

We need non-linearity!



Linear separators and perceptrons

- Perceptrons define linear separators

$$x^T w + b$$

- $x^T w + b > 0$, classify Yes (+1)
- $x^T w + b < 0$, classify No (-1)

- **Unfortunately!** Network of perceptrons still defines only a linear separator
- Linear separators cannot describe XOR

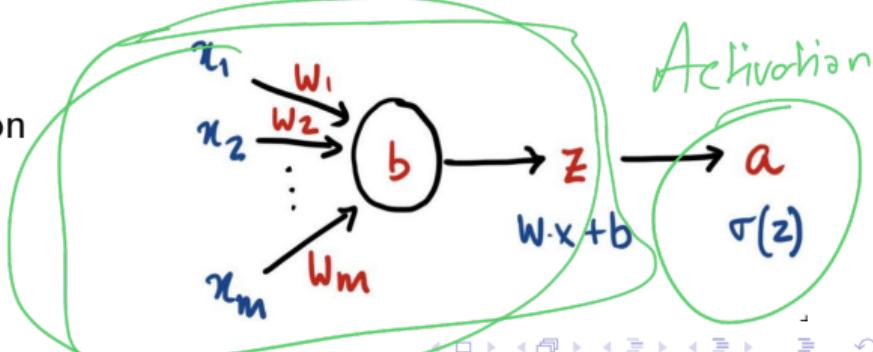
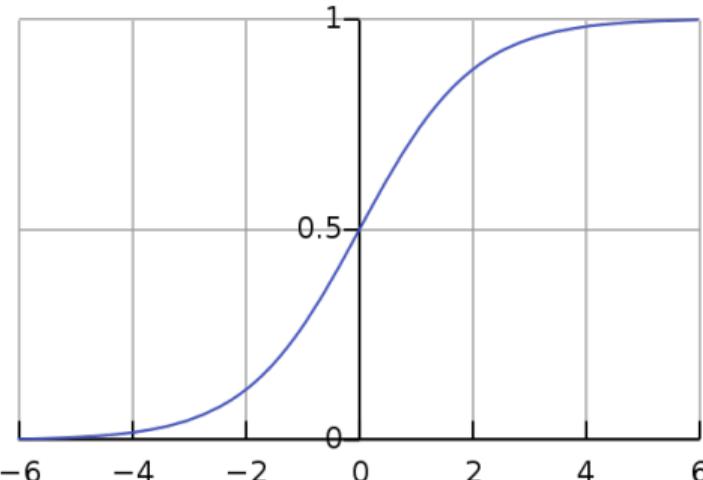
We need non-linearity!

- Introduce a non-linear **activation** function

- Traditionally sigmoid,

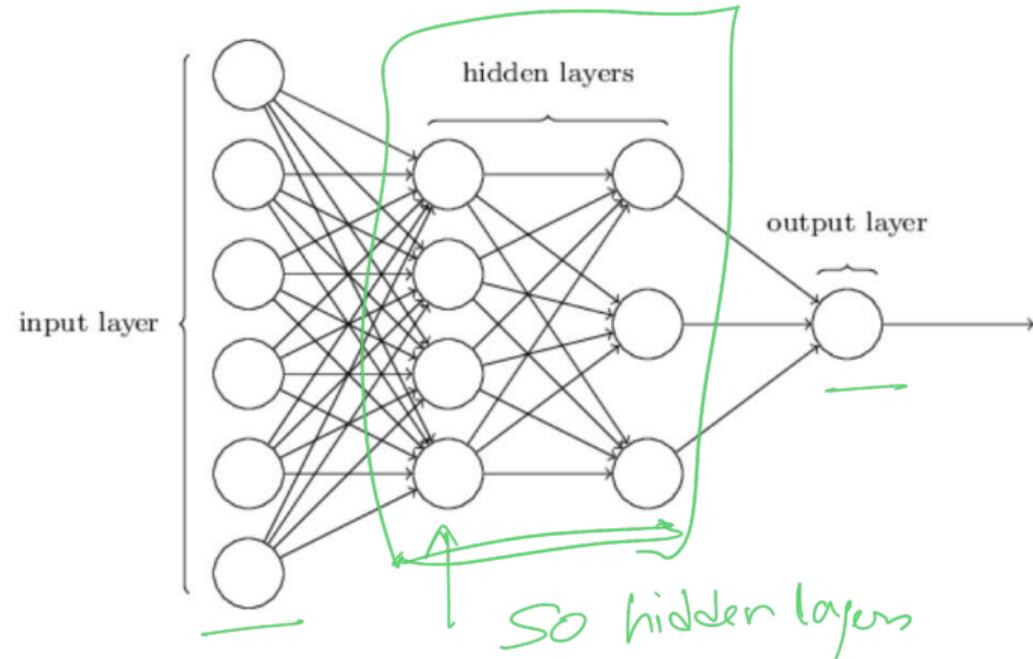
$$\sigma(z) = 1/(1 + e^{-z})$$

This is a neuron!



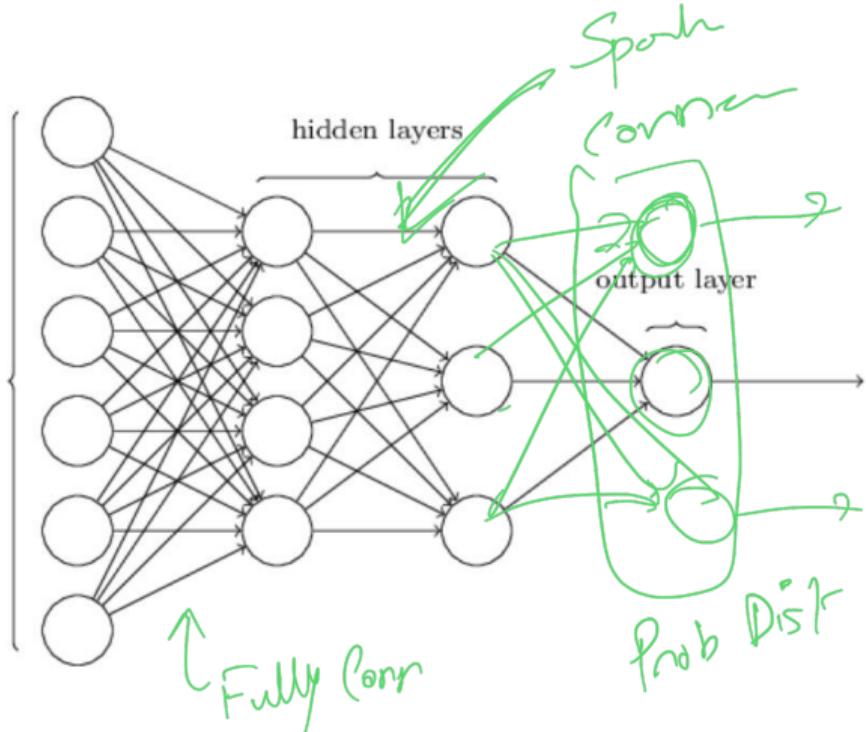
(Feed forward) Neural networks

- Acyclic network of perceptrons with non-linear activation functions
 - **Universal Approximation Theorem:** With just 1 hidden layer, a neural network can approximate any function for any degree of precision.
 - For a function f , it is possible to construct such a neural network. For example the **XOR** function.



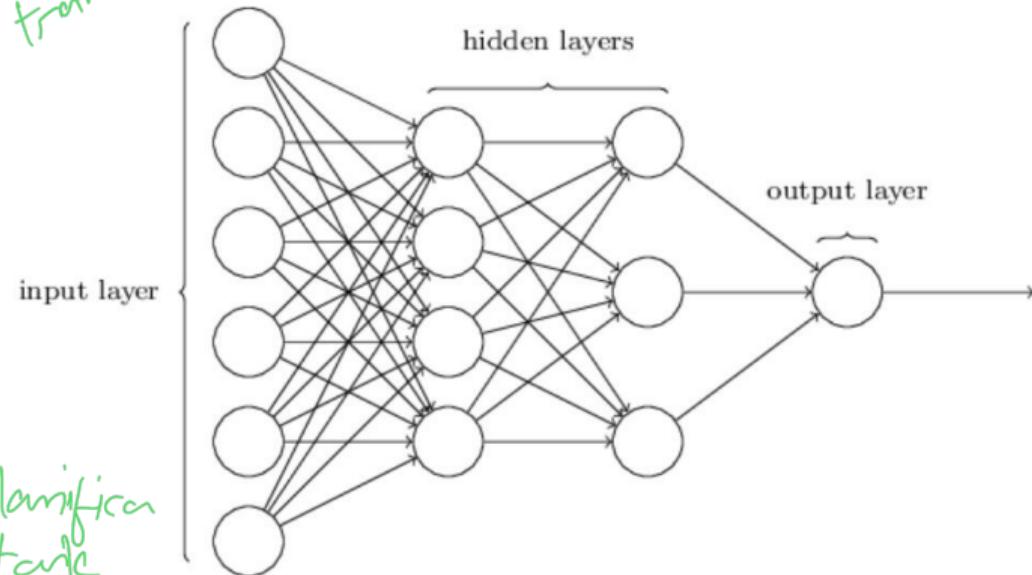
(Feed forward) Neural networks

- Acyclic network of perceptrons with non-linear activation functions \rightarrow neurons
 - The **Structure** of the network and the value of its **Parameters** θ (weights and biases of each neuron).
 - **Objective:** Given a training set S , compute a neural network with $\text{low generalization loss.}$ \leftarrow loss by training
 - We can estimate the generalization loss using a test set $T.$
- of the chosen structure*



(Feed forward) Neural networks

- How do we compute a neural network?
 - Choose the structure of the neural network, based on the ML task at hand.
 - Initially set the weights and biases of neurons to random numbers.
 - Choose a loss-function $\ell(\theta, S)$ on the output of the neural network, e.g. **Cross-Entropy Loss**
 - **Optimization problem:** Given S find the values for θ with least $\ell(\theta, S)$.
- Important* ↗ *also train* ↗



(Feed forward) Neural networks

- How do we compute a neural network?
 - Optimization problem:
Given S find the values for θ with least $\ell(\theta, S)$.

Highly Non-Trivial Problem!

(Feed forward) Neural networks

- How do we compute a neural network?

- Optimization problem:

Given S find the values for θ with least $\ell(\theta, S)$.

Parameter

HyperParameter

Highly Non-Trivial Problem!

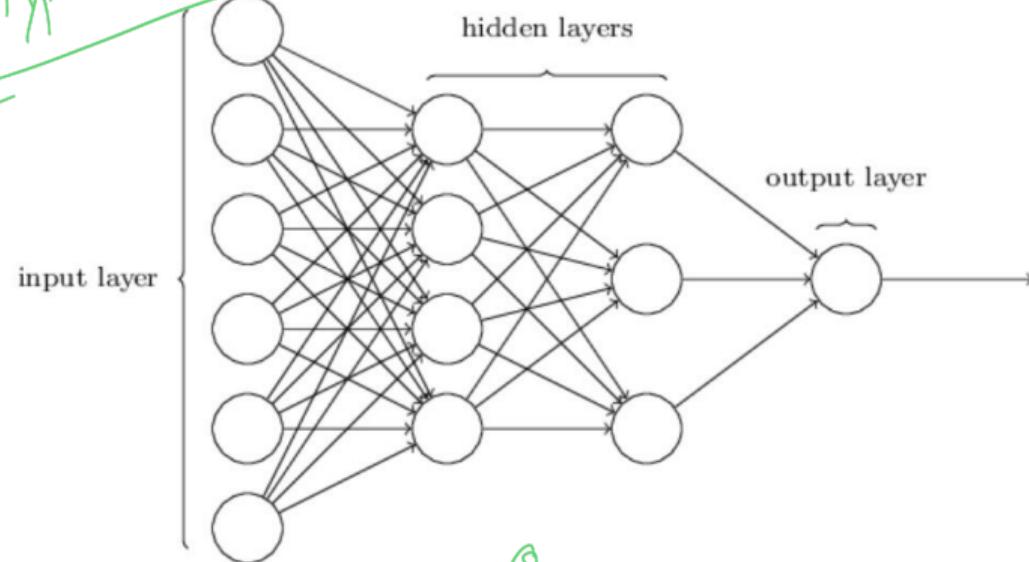
- We use the most basic optimization method:

Gradient Descent

- Update θ so that $\ell(\theta, S)$ decreases

$$\theta \leftarrow \theta - \alpha \cdot \frac{d\ell(\theta, S)}{d\theta}$$

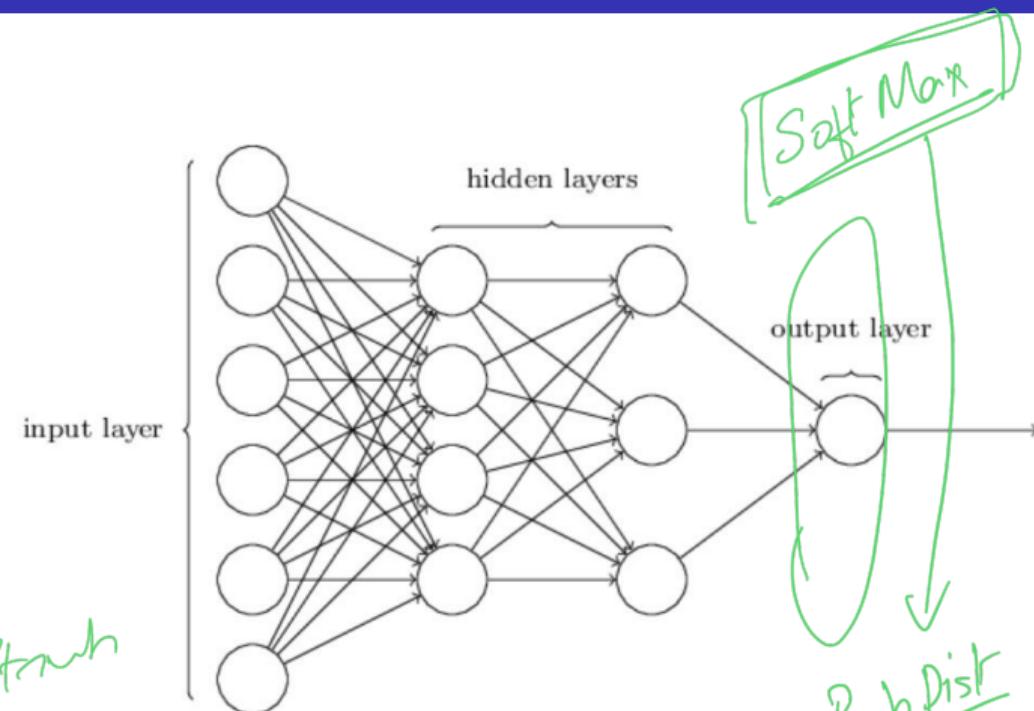
- α is the learning rate



(Feed forward) Neural networks

- Acyclic network of perceptrons with non-linear activation functions
- Ingredients

- Output layer activation function
- Loss function for gradient descent *on S*
- Hidden layer activation functions
- Network architecture: Interconnection of layers
- Initial values of weights and biases



Training a neural network

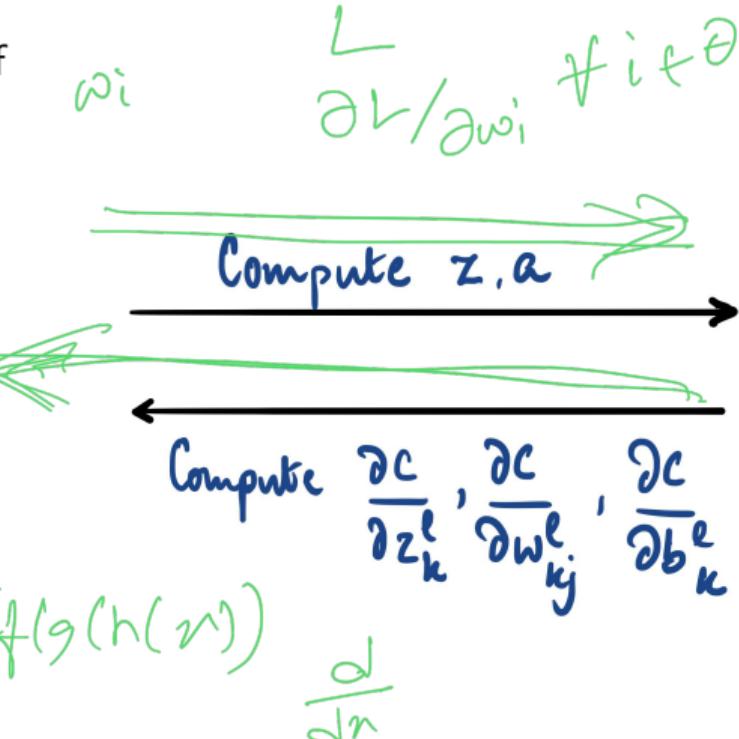
- **Backpropagation** — efficient implementation of gradient descent for neural networks

$$\theta \leftarrow \theta - \alpha \frac{\partial l(\theta, s)}{\partial \theta}$$

Update for w_i

$$w_i \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$$

How do we compute this?



$$\frac{\partial}{\partial n} (f(g(h(n)))$$

$$\frac{\partial}{\partial n}$$

Training a neural network

- Backpropagation — efficient implementation of gradient descent for neural networks

- Forward pass, compute outputs, activation values

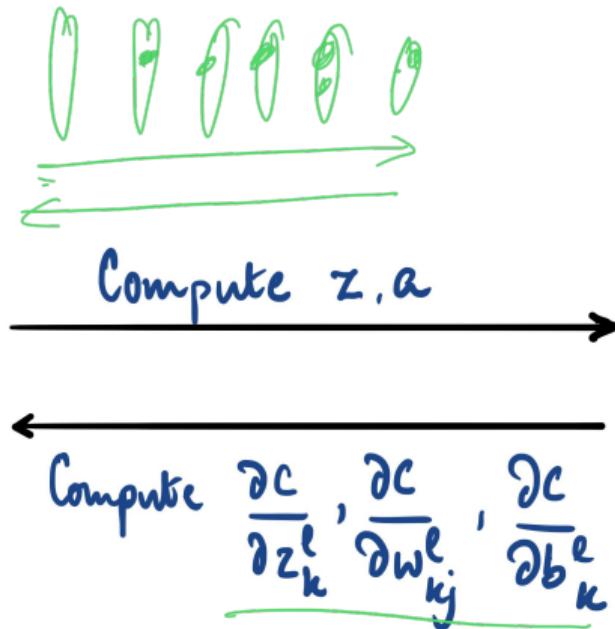
- Backward pass, use chain rule to compute all gradients in one scan

$$f(n) = \alpha n + \beta$$

$$g(n) = \gamma n + \delta$$

$$l(n) = g(f(n))$$

$$\frac{\partial l}{\partial \alpha} = \underbrace{\frac{\partial g}{\partial f}}_{\gamma} \cdot \underbrace{\frac{\partial f}{\partial \alpha}}_{\alpha} = \gamma \cdot n$$



Training a neural network

- Backpropagation — efficient implementation of gradient descent for neural networks
- Forward pass, compute outputs, activation values
- Backward pass, use chain rule to compute all gradients in one scan

- Stochastic gradient descent (SGD)
 - Update parameters in **minibatches**

- **Epoch:** set of minibatches that covers entire training data

$$\frac{\partial L(\theta, s')}{\partial \theta}$$

$$\frac{\partial L(\theta, s)}{\partial \theta}$$

small ^{batch} ^{dom}
sample taken
from S

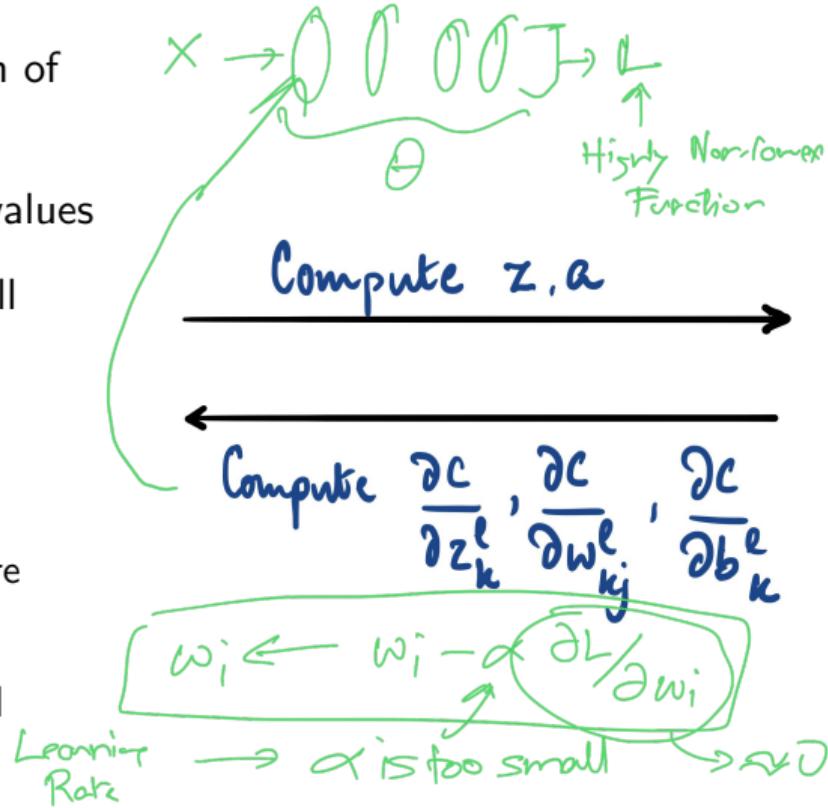
Compute z, a

Compute $\frac{\partial C}{\partial z_k^e}, \frac{\partial C}{\partial w_{kj}^e}, \frac{\partial C}{\partial b_k^e}$

10000 data points

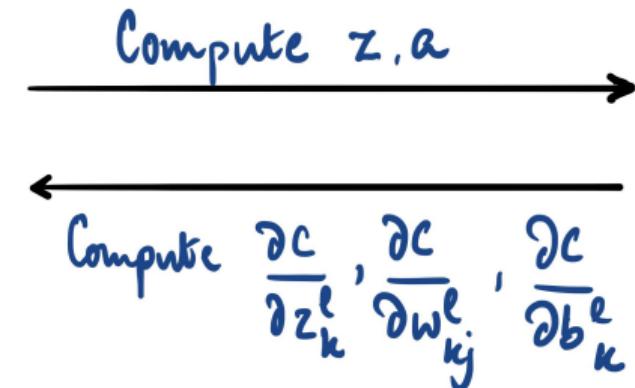
Training a neural network

- Backpropagation — efficient implementation of gradient descent for neural networks
- Forward pass, compute outputs, activation values
- Backward pass, use chain rule to compute all gradients in one scan
- Stochastic gradient descent (SGD)
 - Update parameters in minibatches
 - Epoch: set of minibatches that covers entire training data
- Difficulties: slow convergence, vanishing and exploding gradients



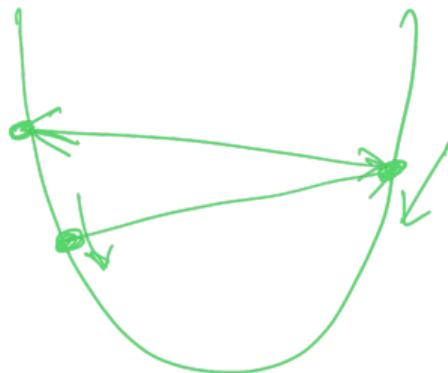
Unstable gradients

- **Vanishing gradients** — gradients become smaller towards lower layers (closer to input)
 - Gradient descent updates leave these layers' parameters virtually unchanged



Unstable gradients

- **Vanishing gradients** — gradients become smaller towards lower layers (closer to input)
 - Gradient descent updates leave these layers' parameters virtually unchanged
- Also exploding gradients, recurrent neural networks with feedback edges



$\frac{\partial L}{\partial w_i}$ is too long

$w_i \leftarrow w_i - \underbrace{\frac{\partial L}{\partial w_i}}$

Compute z, a

←

Compute $\frac{\partial C}{\partial z_k^e}, \frac{\partial C}{\partial w_{kj}^e}, \frac{\partial C}{\partial b_k^e}$

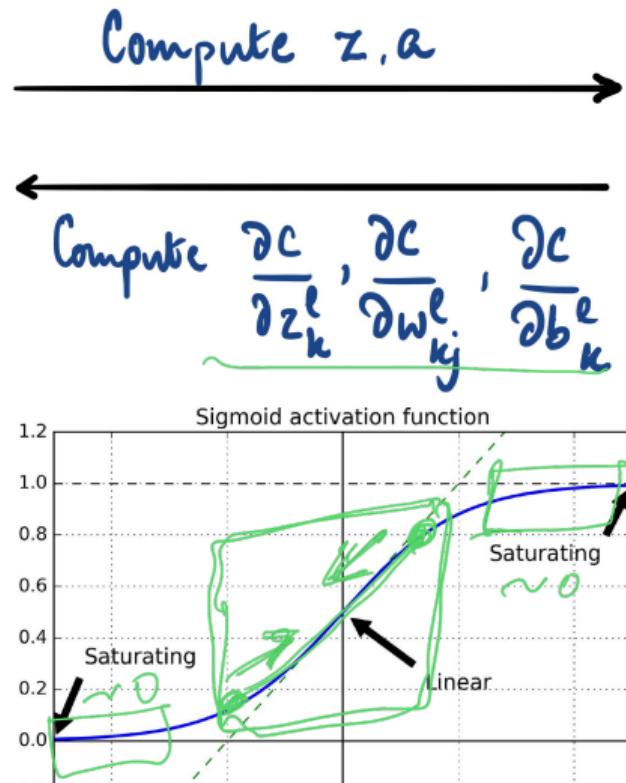
Unstable gradients

- **Vanishing gradients** — gradients become smaller towards lower layers (closer to input)
 - Gradient descent updates leave these layers' parameters virtually unchanged
- Also exploding gradients, recurrent neural networks with feedback edges
- In general, unstable gradients, different layers learn at different speeds

Compute z, a 
Compute $\frac{\partial C}{\partial z_k^e}, \frac{\partial C}{\partial w_{kj}^e}, \frac{\partial C}{\partial b_k^e}$ 

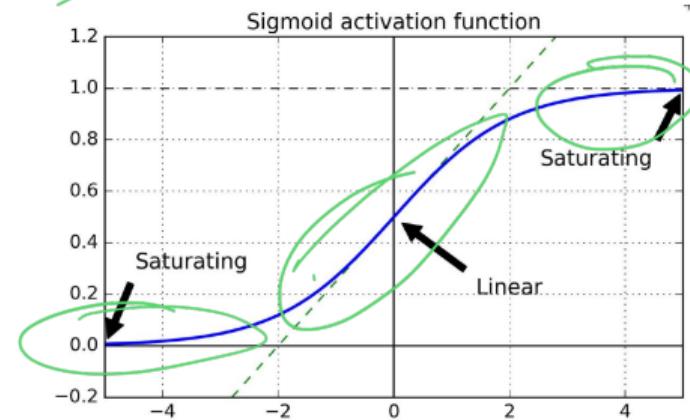
Unstable gradients

- **Vanishing gradients** — gradients become smaller towards lower layers (closer to input)
 - Gradient descent updates leave these layers' parameters virtually unchanged
- Also exploding gradients, recurrent neural networks with feedback edges
- In general, unstable gradients, different layers learn at different speeds
- **[Xavier Glorot and Joshua Bengio, 2010]**
 - Random initialization, traditionally Gaussian distribution $\mathcal{N}(0, 1)$
 - Variance keeps increasing going forward
 - Saturating sigmoid function



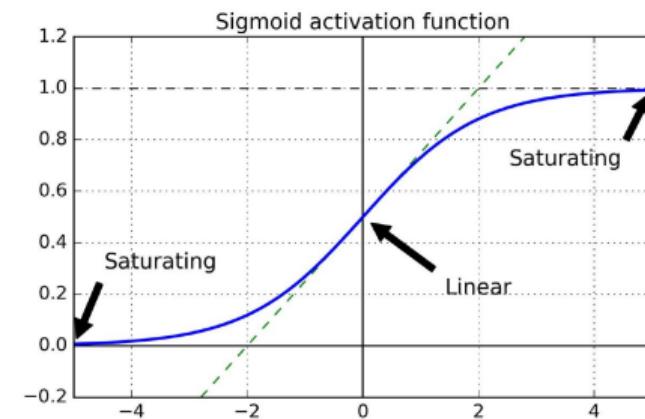
Initializing neural networks

- Want “signal” to flow well in both directions during backpropagation
 - Signal should not die out, explode, saturate



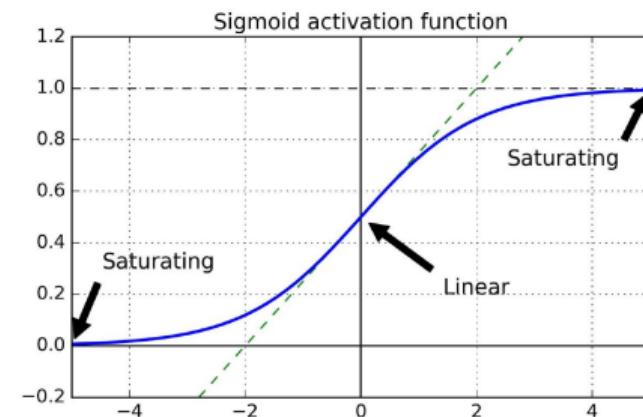
Initializing neural networks

- Want “signal” to flow well in both directions during backpropagation
 - Signal should not die out, explode, saturate
- [Glorot,Bengio] Gradients should have equal variance before and after flowing through a layer in both directions
 - Equal variance requires $fan_{in} = fan_{out}$



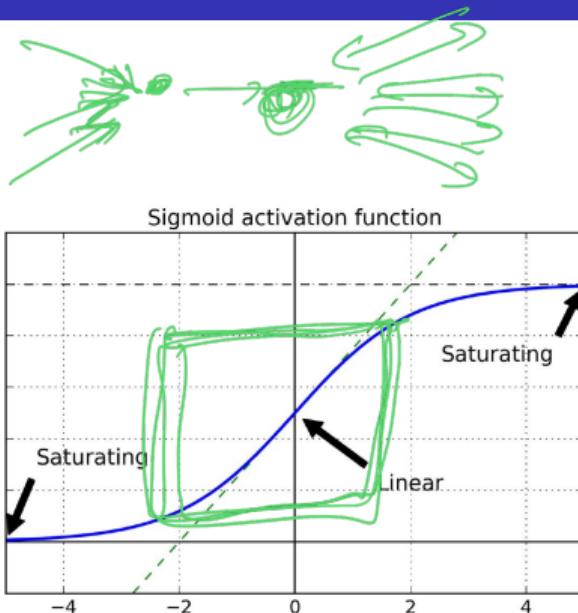
Initializing neural networks

- Want “signal” to flow well in both directions during backpropagation
 - Signal should not die out, explode, saturate
- [Glorot,Bengio] Gradients should have equal variance before and after flowing through a layer in both directions
 - Equal variance requires $fan_{in} = fan_{out}$
- Let $fan_{avg} = (fan_{in} + fan_{out})/2$



Initializing neural networks

- Want “signal” to flow well in both directions during backpropagation
 - Signal should not die out, explode, saturate
- [Glorot, Bengio] Gradients should have equal variance before and after flowing through a layer in both directions
 - Equal variance requires $fan_{in} = fan_{out}$
- Let $fan_{avg} = (fan_{in} + fan_{out})/2$
- Initialize with
 - Gaussian, $\mathcal{N}(0, 1/fan_{avg})$
 - Uniform, $\mathcal{U}(-r, r)$, $r = \sqrt{\frac{3}{fan_{avg}}}$



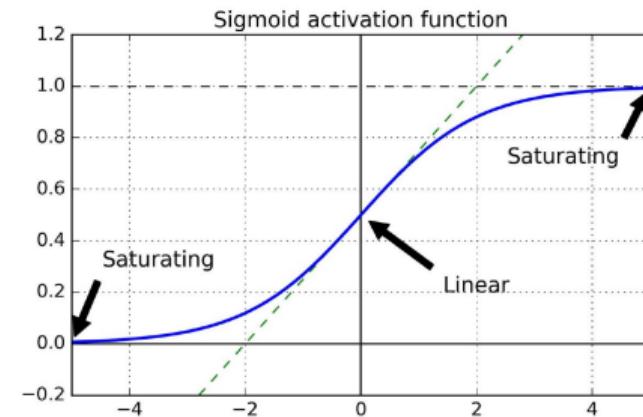
Initializing neural networks

- Let $fan_{avg} = (fan_{in} + fan_{out})/2$

- Initialize with

- Gaussian, $\mathcal{N}(0, 1/fan_{avg})$

- Uniform, $\mathcal{U}(-r, r)$, $r = \sqrt{\frac{3}{fan_{avg}}}$



Initializing neural networks

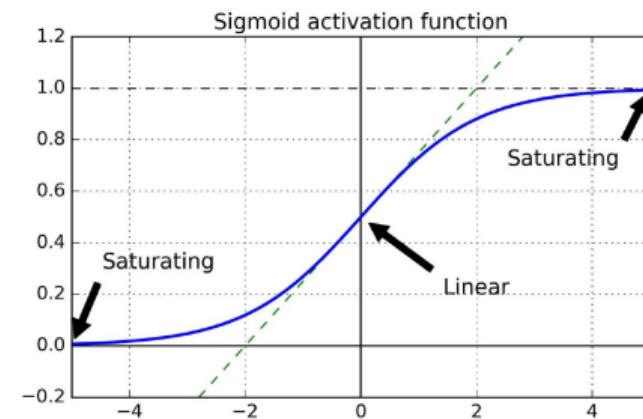
- Let $fan_{avg} = (fan_{in} + fan_{out})/2$

- Initialize with

- Gaussian, $\mathcal{N}(0, 1/fan_{avg})$

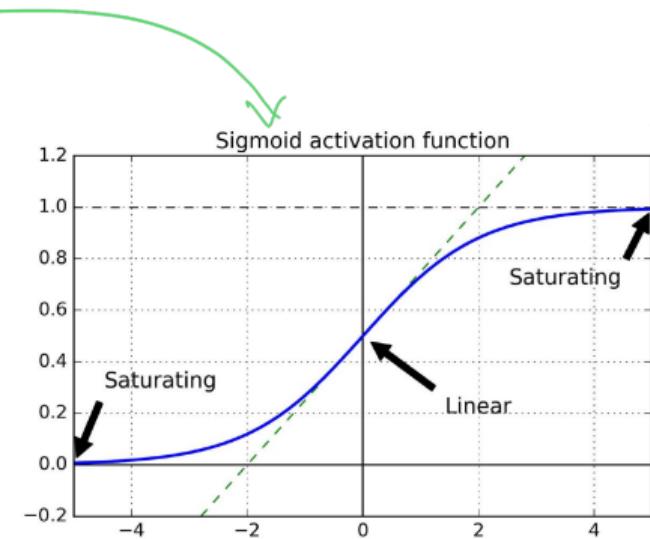
- Uniform, $\mathcal{U}(-r, r)$, $r = \sqrt{\frac{3}{fan_{avg}}}$

- [Yann LeCun, 1990s] earlier proposed the same with fan_{avg} replaced by fan_{in}
 - Equivalent if $fan_{in} = fan_{out}$



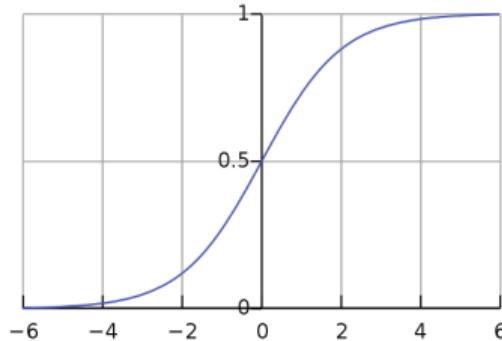
Initializing neural networks

- Let $fan_{avg} = (fan_{in} + fan_{out})/2$
- Initialize with
 - Gaussian, $\mathcal{N}(0, 1/fan_{avg})$
 - Uniform, $\mathcal{U}(-r, r)$, $r = \sqrt{\frac{3}{fan_{avg}}}$
- [Yann LeCun, 1990s] earlier proposed the same with fan_{avg} replaced by fan_{in}
 - Equivalent if $fan_{in} = fan_{out}$
- Other choices for specific activation function
 - ReLU, [He et al, 2015], $\mathcal{N}(0, 2/fan_{in})$



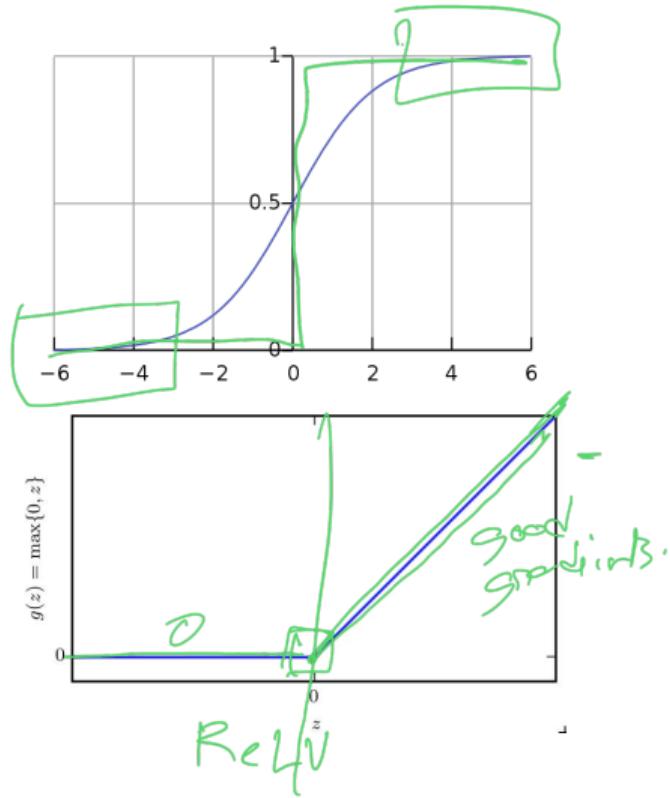
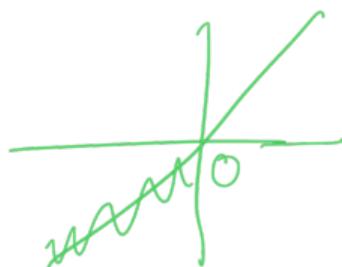
Non-saturating activation functions

- Sigmoid was initially chosen as a “smooth” step



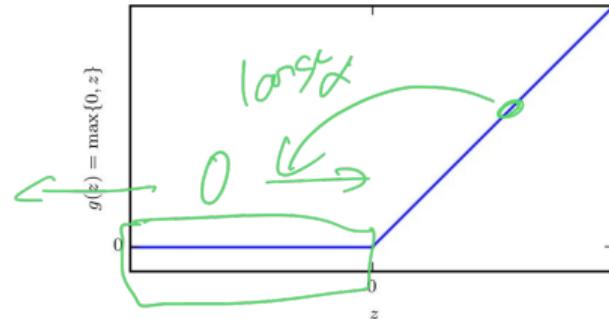
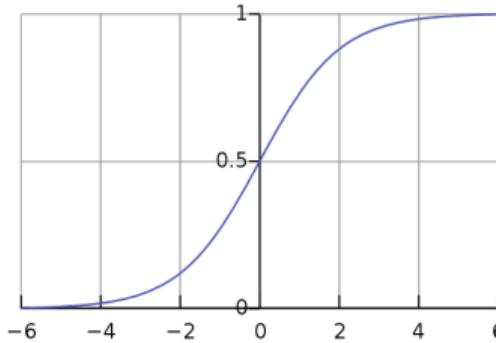
Non-saturating activation functions

- Sigmoid was initially chosen as a “smooth” step
- Rectified linear unit (ReLU):
$$g(z) = \max(0, z)$$
 - Fast to compute
 - Non-differentiable point not a bottleneck



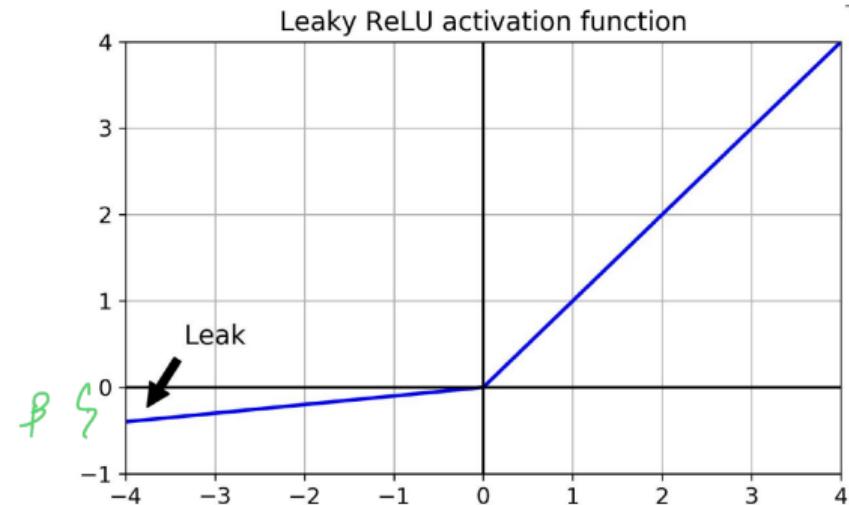
Non-saturating activation functions

- Sigmoid was initially chosen as a “smooth” step
- Rectified linear unit (ReLU):
$$g(z) = \max(0, z)$$
 - Fast to compute
 - Non-differentiable point not a bottleneck
- “Dying ReLU”
 - Neuron dies — weighted sum of outputs is negative for all training samples
 - With a large learning rate, half the network may die!



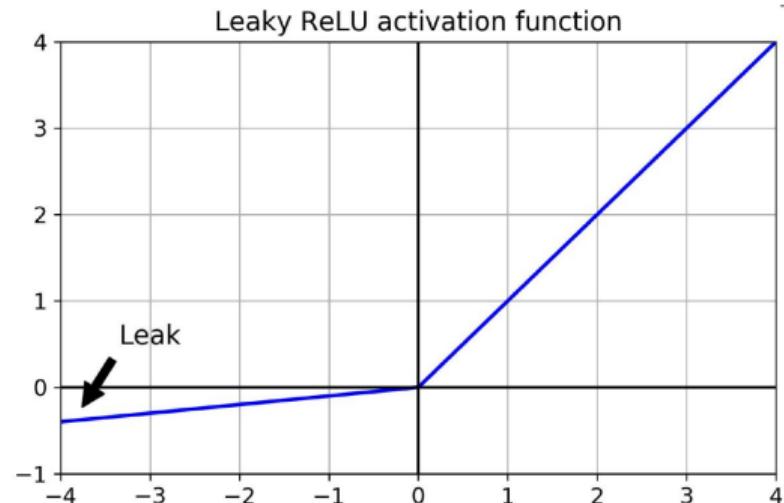
Non-saturating activation functions

- Leaky ReLU, $\max(\beta z, z)$
 - “Leak” β is a hyperparameter



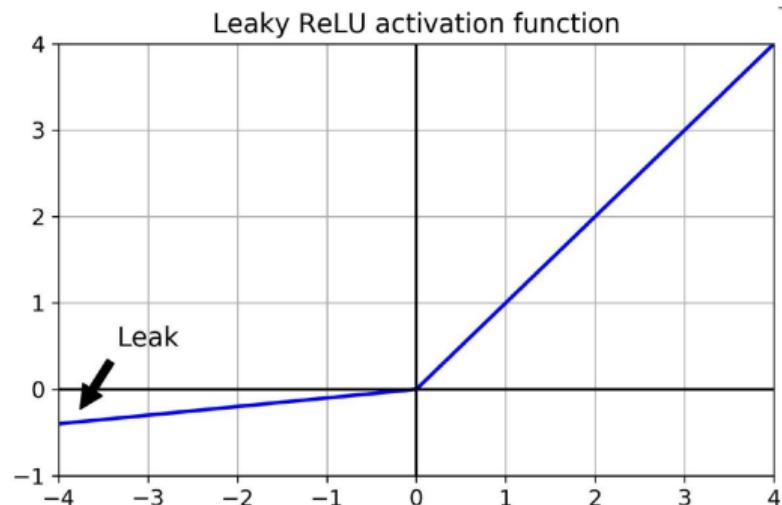
Non-saturating activation functions

- Leaky ReLU, $\max(\alpha z, z)$
 - “Leak” α is a hyperparameter
- RReLU — random leak
 - Pick α from a random range during training
 - Fix to an average value when testing
 - Seems to work well, act as a regularizer



Non-saturating activation functions

- Leaky ReLU, $\max(\alpha z, z)$
 - “Leak” α is a hyperparameter
- RReLU — random leak
 - Pick α from a random range during training
 - Fix to an average value when testing
 - Seems to work well, act as a regularizer
- PReLU — parametric ReLU [He et al, 2015]
 - α is learned during training
 - Often outperforms ReLU, but could lead to overfitting



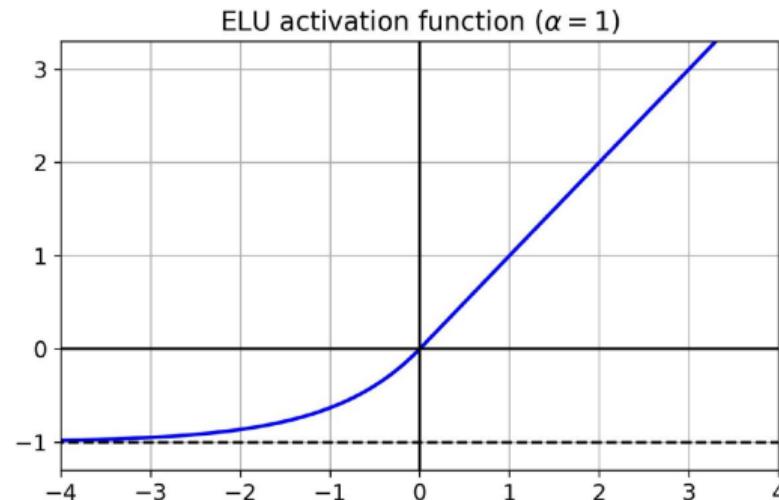
Non-saturating activation functions

■ ELU — Exponential Linear Unit

[Clevert et al, 2015]

$$ELU_{\alpha}(z) = \begin{cases} \alpha(e^z - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

- Training converges faster
- Computing exponential is slower
- In practice, slower than ReLU



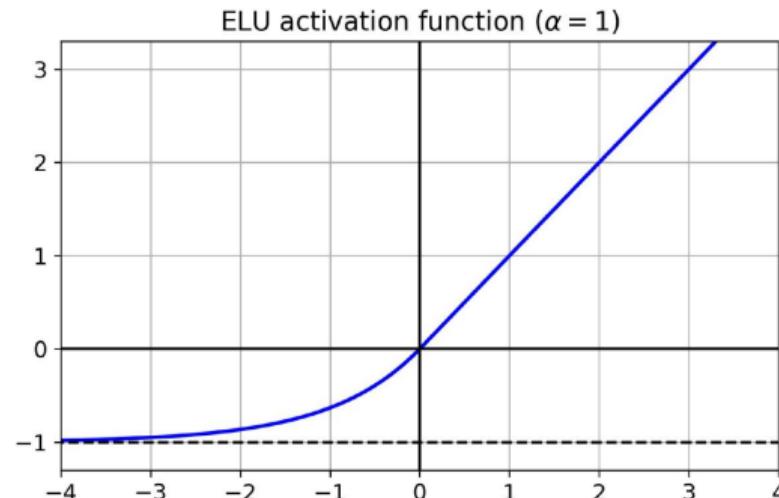
Non-saturating activation functions

- SELU — Scaled ELU

[Klambauer et al, 2017]

$$SELU_{\alpha}(z) = \lambda \begin{cases} \alpha(e^z - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$

- **Self-normalizing** — output of each layer preserves mean **0** and standard deviation **1** during training
- Use LeCun initialization, $\mathcal{N}(0, 1/fan_{in})$



Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch

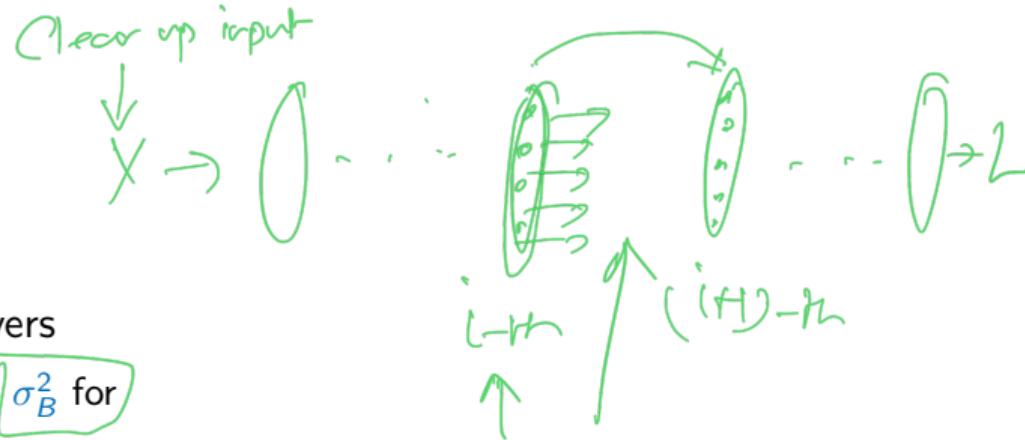
Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
 - Zero-centre and normalize each input

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
 - Zero-centre and normalize each input
 - Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$



Maybe
Clear Up
here too

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \text{Random Noise}$$

$$z_i = \lambda \cdot \hat{x}_i + \beta$$

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers

- Estimate mean μ_B and variance σ_B^2 for inputs across minibatch

- Zero-centre and normalize each input

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

- Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$
- Learn optimal scaling and shifting parameters for each layer



Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers

- Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
- Zero-centre and normalize each input

$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

- Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$
- Learn optimal scaling and shifting parameters for each layer

Batch normalization [Joffe, Szegedy 2015]

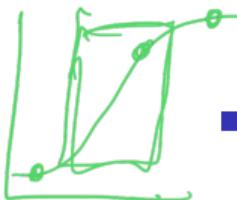
- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
 - Zero-centre and normalize each input
$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$
 - Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$
 - Learn optimal scaling and shifting parameters for each layer

- At input, BN layer avoids need for standardizing
- Difficulties
 - Mean and variance differ across minibatches
 - How to estimate parameters for entire dataset?
 - Practical solution: maintain a moving average of means and standard deviations for each layer

not much weaker
Good estimate for
curr minibatch

Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
 - Zero-centre and normalize each input
 - Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$
 - Learn optimal scaling and shifting parameters for each layer
- At input, BN layer avoids need for standardizing
- Difficulties
 - Mean and variance differ across minibatches
 - How to estimate parameters for entire dataset?
 - Practical solution: maintain a moving average of means and standard deviations for each layer
- Batch normalization greatly speeds up learning rate

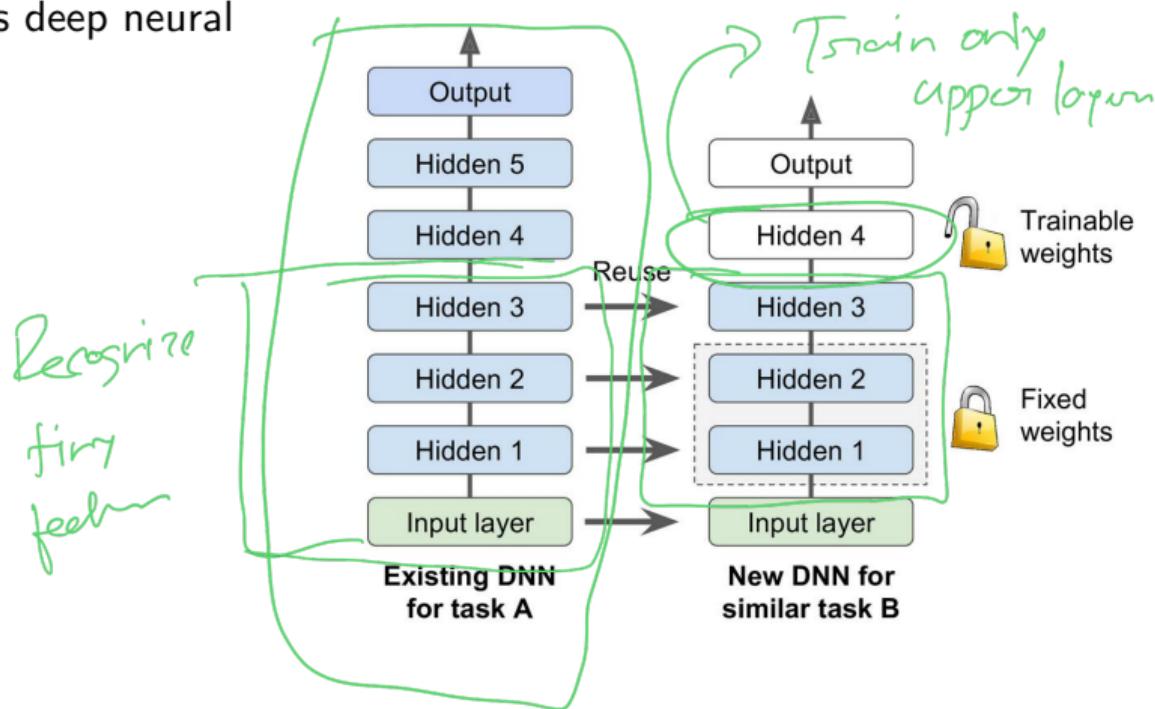


Batch normalization [Joffe, Szegedy 2015]

- Good activation function and initialization mitigates vanishing/exploding gradients
- May still recur during training
- Add batch normalization (BN) layers
 - Estimate mean μ_B and variance σ_B^2 for inputs across minibatch
 - Zero-centre and normalize each input
$$\hat{x}_i = \frac{x_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$
 - Scale and shift $z_i = \lambda \cdot \hat{x}_i + \beta$
 - Learn optimal scaling and shifting parameters for each layer
- At input, BN layer avoids need for standardizing
- Difficulties
 - Mean and variance differ across minibatches
 - How to estimate parameters for entire dataset?
 - Practical solution: maintain a moving average of means and standard deviations for each layer
- Batch normalization greatly speeds up learning rate
- Even works as a regularizer!

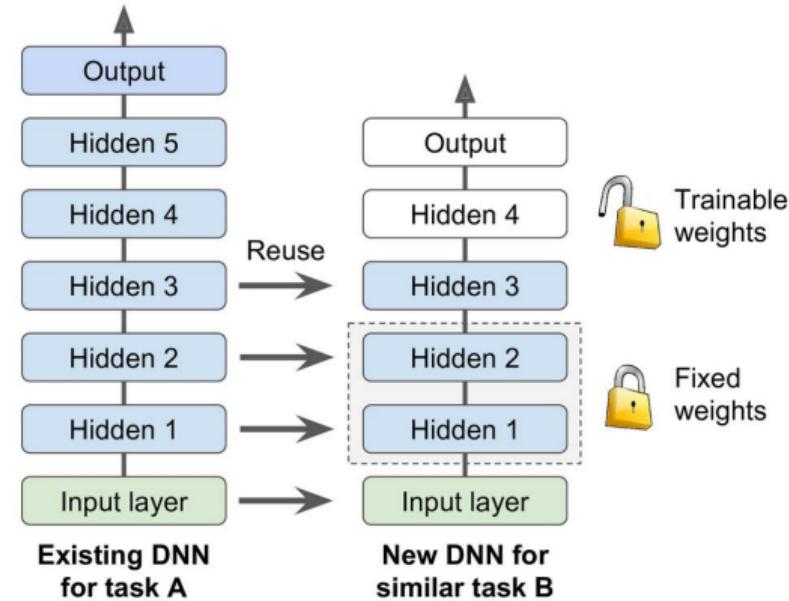
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)



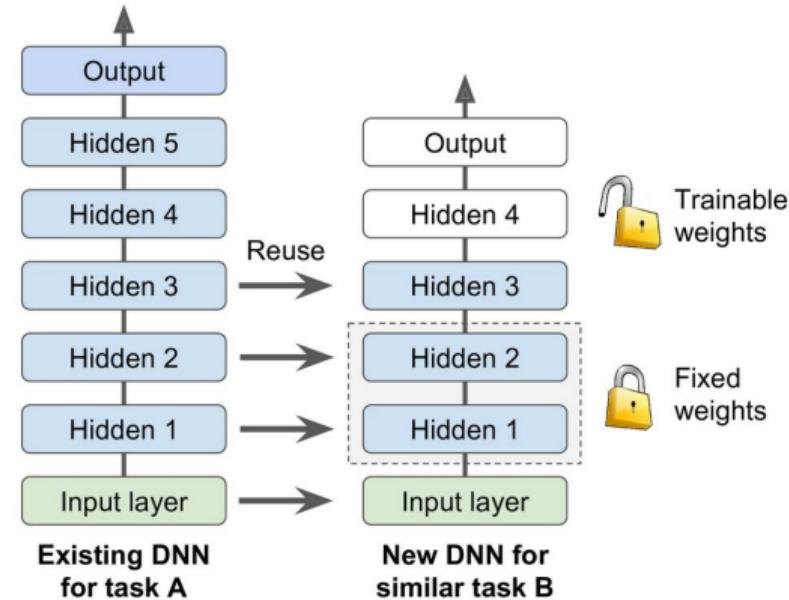
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)



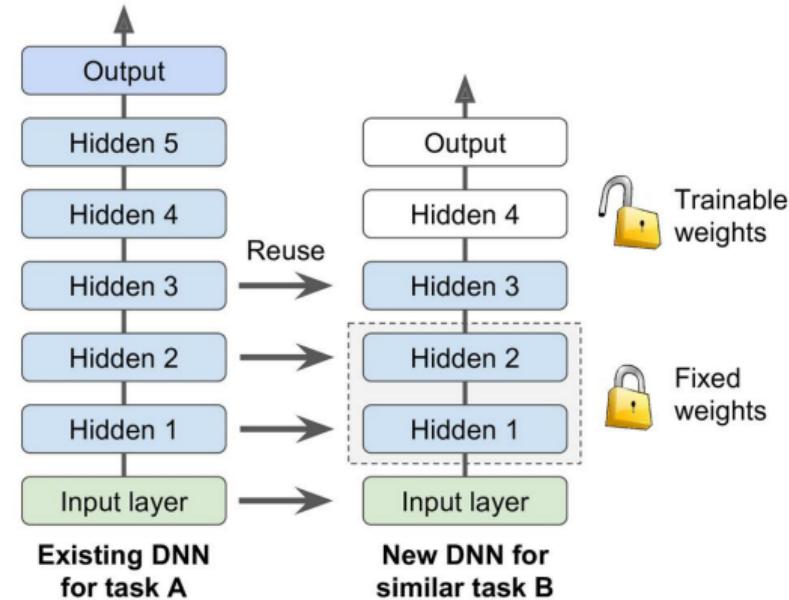
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)
- New DNN to classify types of vehicles



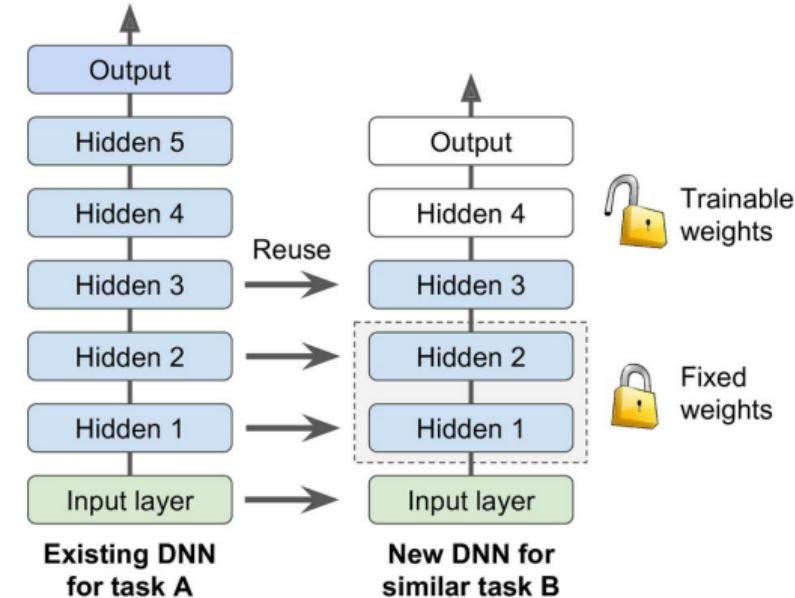
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)
- New DNN to classify types of vehicles
- Tasks similar, even overlapping



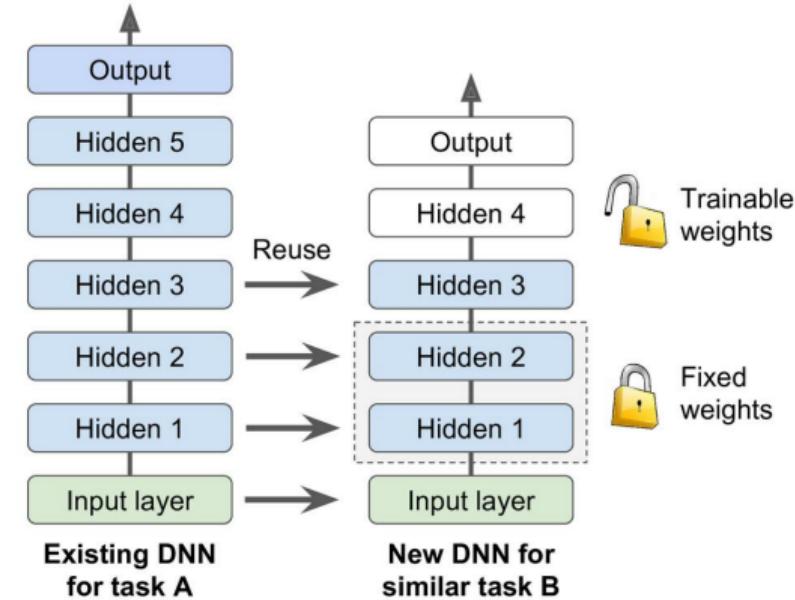
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)
- New DNN to classify types of vehicles
- Tasks similar, even overlapping
- Lower layers identify basic features, upper layers combine them to classify



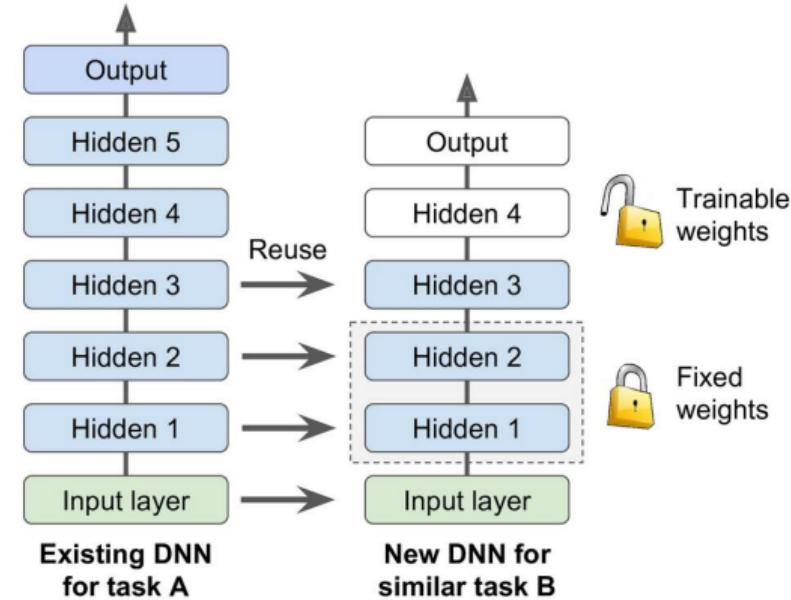
Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)
- New DNN to classify types of vehicles
- Tasks similar, even overlapping
- Lower layers identify basic features, upper layers combine them to classify
- Freeze weights of lower layers, re-learn upper layers



Transfer learning

- Reuse trained layers across deep neural networks (DNNs)
- Old DNN trained on images of daily objects (animals, plants, vehicles, ...)
- New DNN to classify types of vehicles
- Tasks similar, even overlapping
- Lower layers identify basic features, upper layers combine them to classify
- Freeze weights of lower layers, re-learn upper layers
- Unfreeze in stages to determine how much to reuse



Still to come

- Optimizing rate of updates in backpropagation

$$w_i \leftarrow w_i - \alpha \frac{\partial L}{\partial w_i}$$

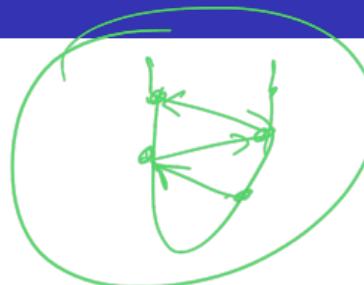


Very important value

$-\alpha$ too small \rightarrow long time to train

$-\alpha$ too large \rightarrow unable to properly minimize loss

- Want α large at start of training
- Small α toward end of training.



Still to come

- Optimizing rate of updates in backpropagation
- How problematic are local minima?

Still to come

- Optimizing rate of updates in backpropagation
- How problematic are local minima?
- Identifying and dealing with unstable gradients