

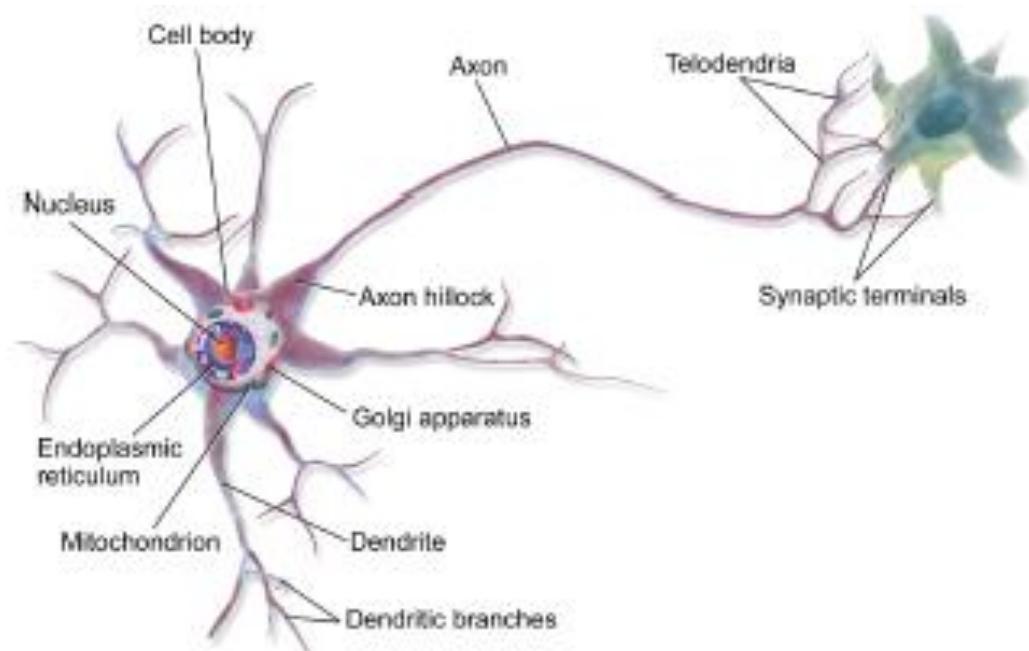
# *Deep Neural Networks* **(DNN)**



Saeed Sharifian

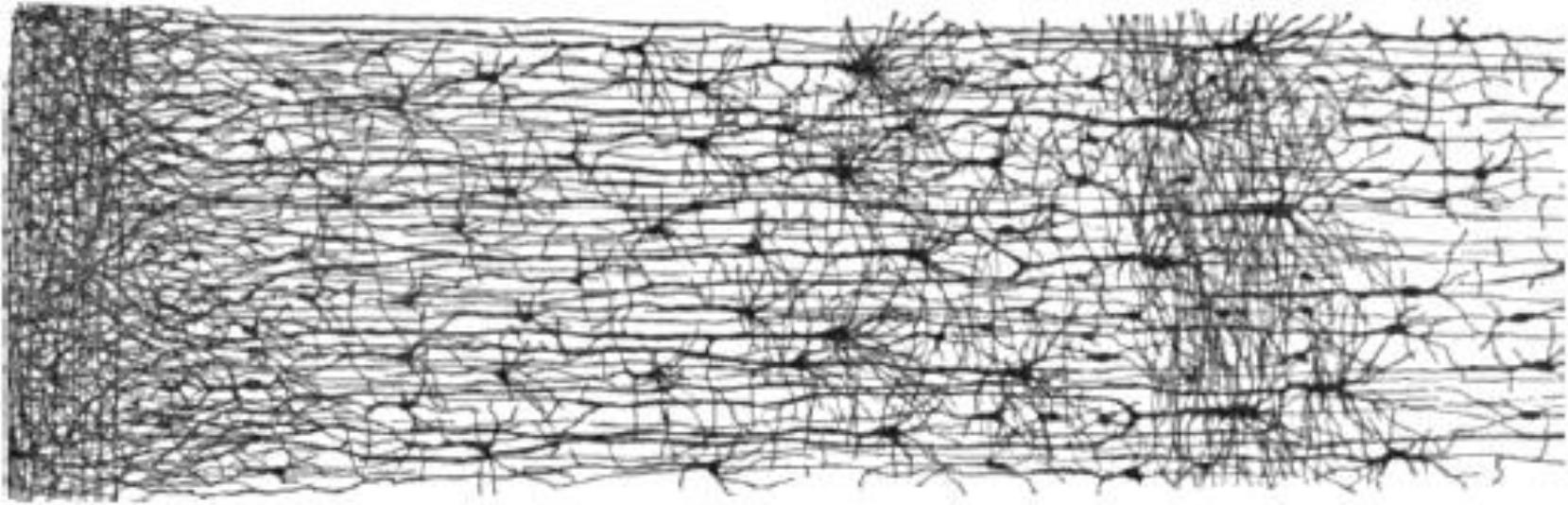
# Biological Neurons

- ▶ Brain architecture has inspired artificial neural networks.
- ▶ A biological neuron is composed of
  - Cell body, many dendrites (branching extensions), one axon (long extension), synapses
- ▶ Biological neurons receive signals from other neurons via these synapses.
- ▶ When a neuron receives a sufficient number of signals within a few milliseconds, it fires its own signals.



# Biological Neurons

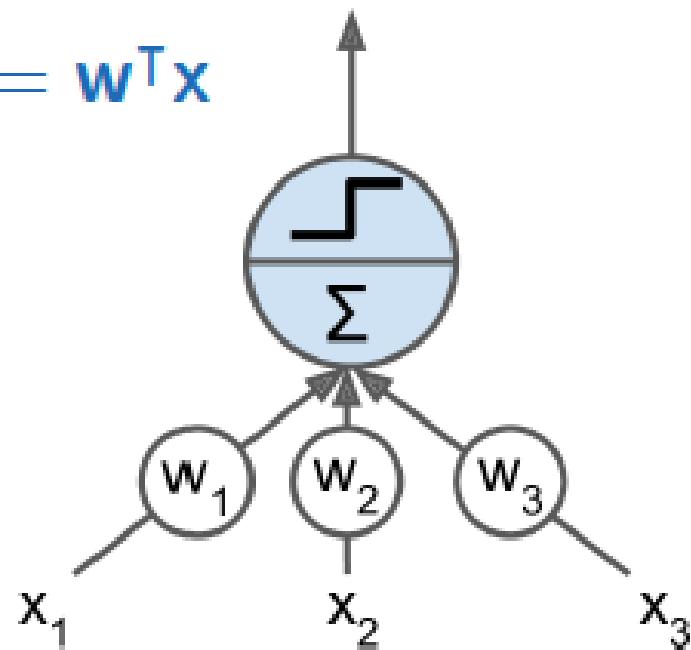
- ▶ Biological neurons are organized in a vast **network** of billions of neurons.
- ▶ Each neuron typically is **connected** to **thousands** of other neurons.



# The Linear Threshold Unit (LTU)

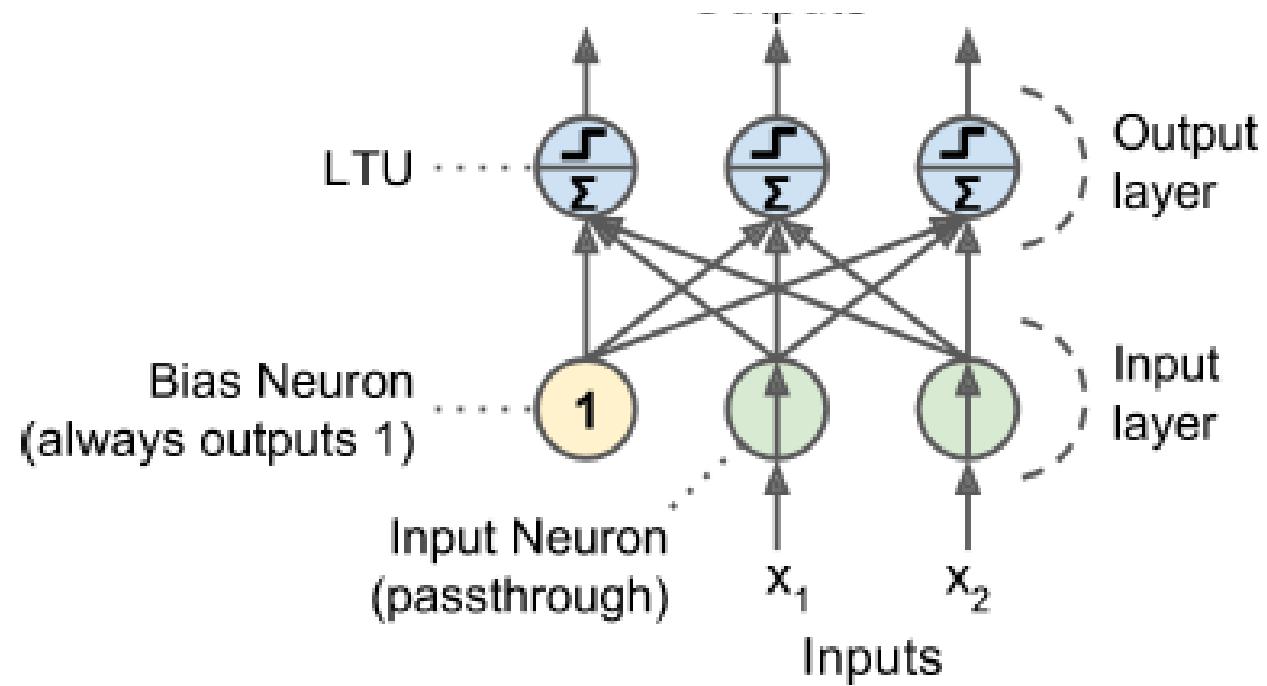
- ▶ Inputs of a LTU are **numbers**
- ▶ Each **input connection** is associated with a **weight**.
- ▶ Computes a **weighted sum** of its inputs and applies a **step function** to that **sum**.

- ▶  $z = w_1x_1 + w_2x_2 + \dots + w_nx_n = \mathbf{w}^T \mathbf{x}$
- ▶  $\hat{y} = \text{step}(z) = \text{step}(\mathbf{w}^T \mathbf{x})$



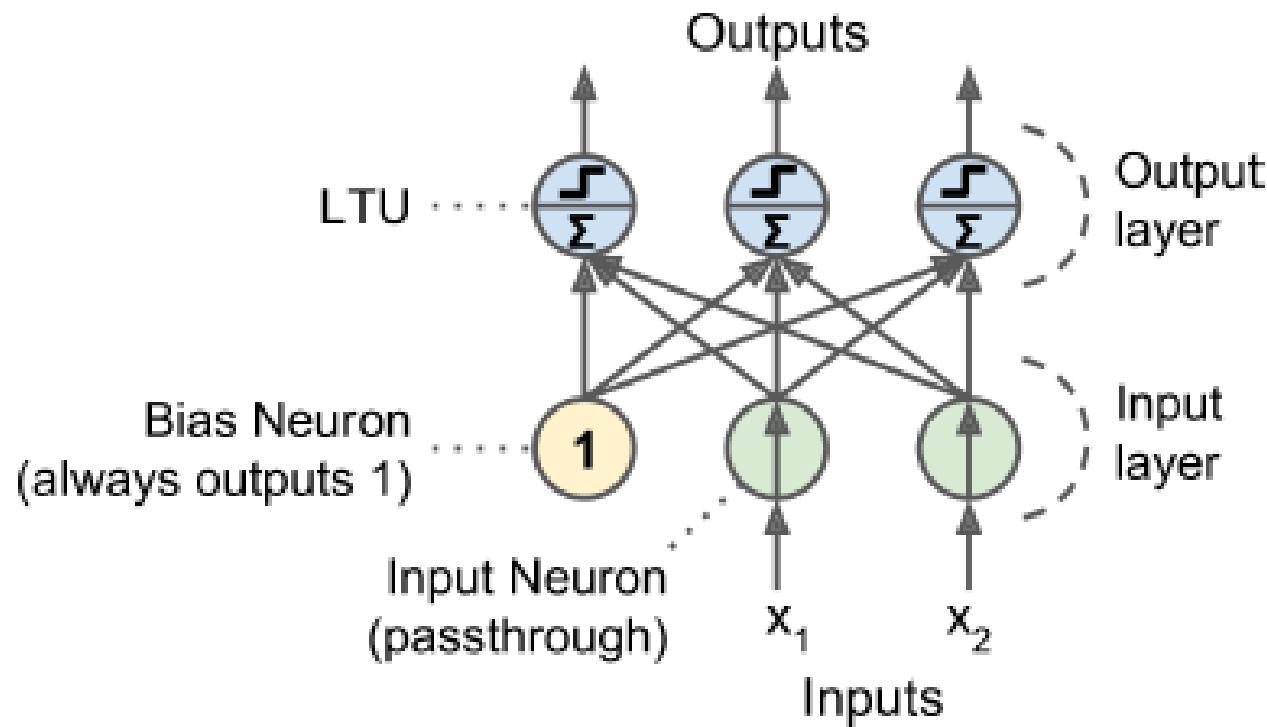
# The Perceptron

- ▶ The perceptron is a single layer of LTUs.
- ▶ The input neurons output whatever input they are fed.
- ▶ A bias neuron, which just outputs 1 all the time.
- ▶ If we use logistic function (sigmoid) instead of a step function, it computes a continuous output.



# How is a Perceptron Trained?

- ▶ The Perceptron training algorithm is inspired by Hebb's rule.
- ▶ When a biological neuron often triggers another neuron, the connection between these two neurons grows stronger.



# How is a Perceptron Trained?

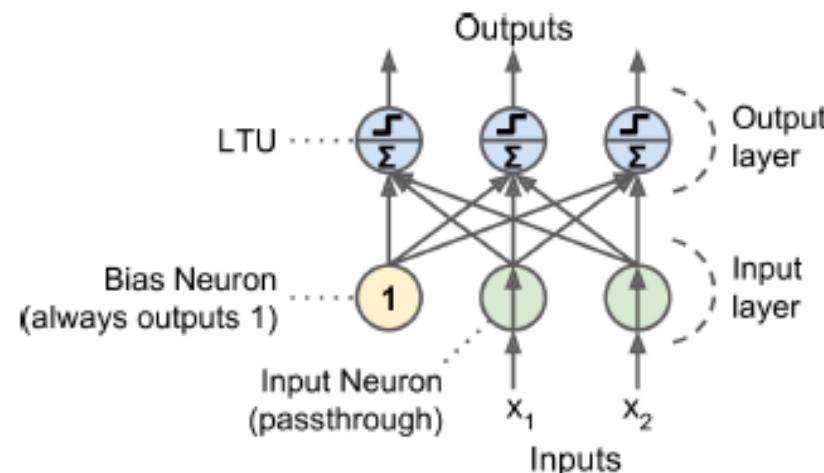
- ▶ Feed one training instance  $\mathbf{x}$  to each neuron  $j$  at a time and make its prediction  $\hat{y}_j$ .
- ▶ Update the connection weights.

$$\hat{y}_j = \sigma(\mathbf{w}_j^T \mathbf{x} + b)$$

$$J(\mathbf{w}_j) = \text{cross\_entropy}(y_j, \hat{y}_j)$$

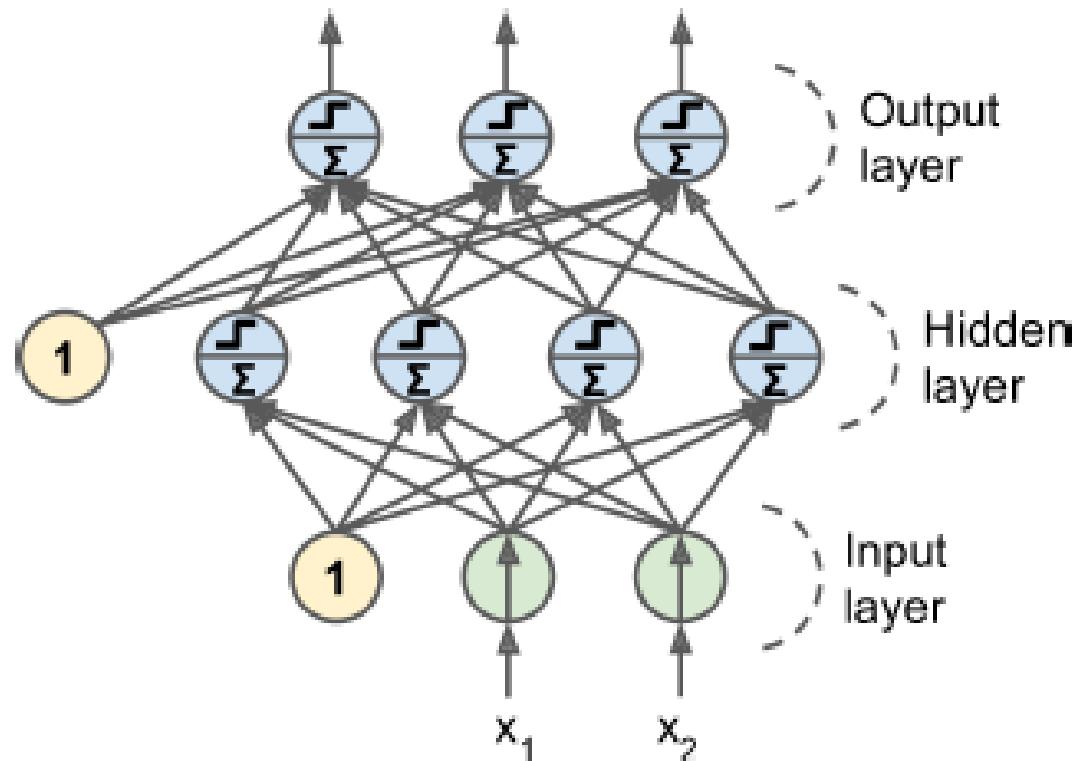
$$\mathbf{w}_{i,j}^{(\text{next})} = \mathbf{w}_{i,j} - \eta \frac{\partial J(\mathbf{w}_j)}{\mathbf{w}_i}$$

- ▶  $w_{i,j}$ : the weight between neurons  $i$  and  $j$ .
- ▶  $x_i$ : the  $i$ th input value.
- ▶  $\hat{y}_j$ : the  $j$ th predicted output value.
- ▶  $y_j$ : the  $j$ th true output value.
- ▶  $\eta$ : the learning rate.



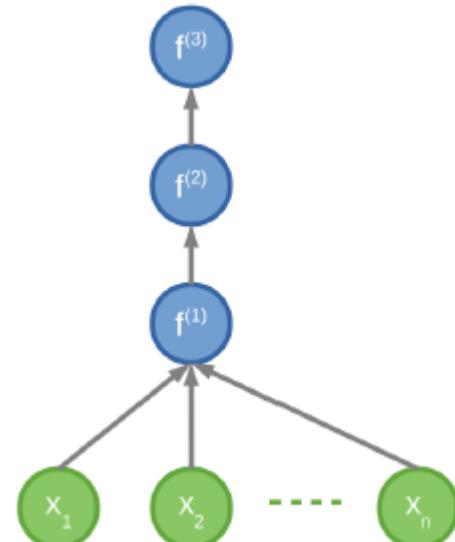
# Multi-Layer Perceptron (MLP)

- ▶ A feedforward neural network is composed of:
  - One input layer
  - One or more hidden layers
  - One final output layer
- ▶ Every layer except the output layer includes a bias neuron and is fully connected to the next layer.



# How Does it Work?

- ▶ The model is associated with a directed acyclic graph describing how the functions are composed together.
- ▶ E.g., assume a network with just a single neuron in each layer.
- ▶ Also assume we have three functions  $f^{(1)}$ ,  $f^{(2)}$ , and  $f^{(3)}$  connected in a chain:  $\hat{y} = f(x) = f^{(3)}(f^{(2)}(f^{(1)}(x)))$
- ▶  $f^{(1)}$  is called the first layer of the network.
- ▶  $f^{(2)}$  is called the second layer, and so on.
- ▶ The length of the chain gives the depth of the model



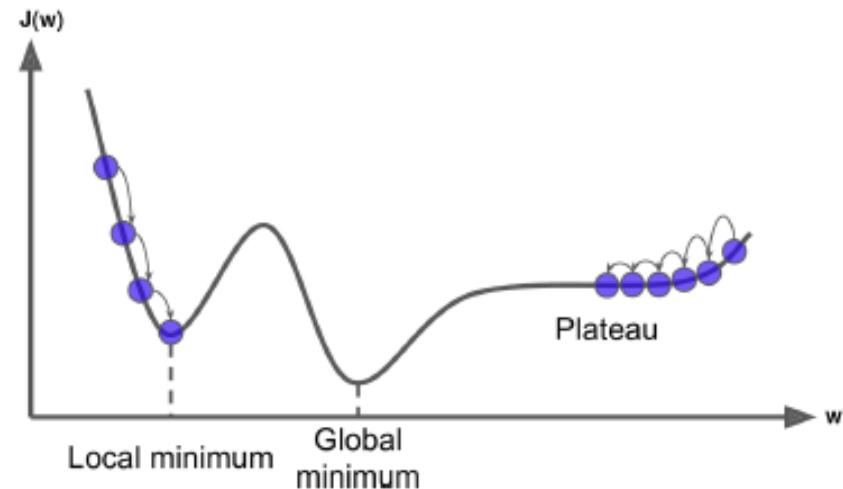
# How to Learn Model Parameters W?

- We use the **cross-entropy** (minimizing the negative log-likelihood) between the training data  $y$  and the model's predictions  $\hat{y}$  as the **cost function**.

$$\text{cost}(y, \hat{y}) = - \sum_i y_j \log(\hat{y}_j)$$

## Gradient-Based Learning

- The **most significant difference** between the **linear models** we have seen so far and **feedforward neural network**?
- The **non-linearity** of a neural network causes its **cost functions** to become **non-convex**.
- Linear models, with **convex cost function**, **guarantee** to find **global minimum**.
  - Convex optimization converges starting from **any initial parameters**.

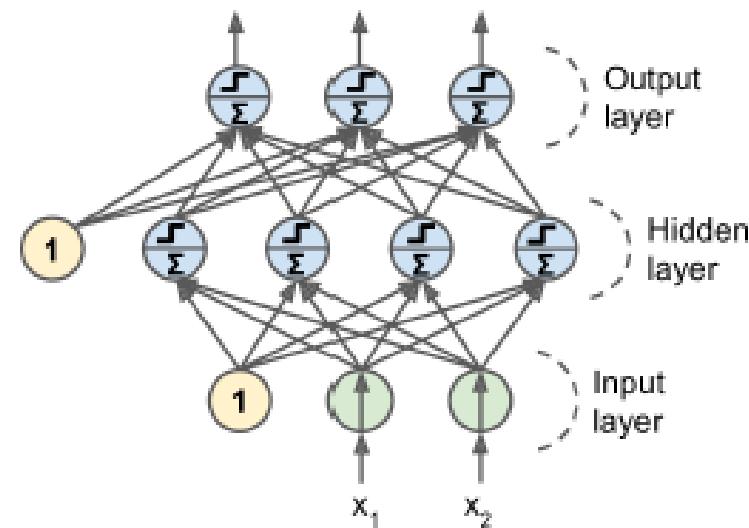


# Gradient-Based Learning

- ▶ Stochastic gradient descent applied to **non-convex cost functions** has no such convergence guarantee.
- ▶ It is **sensitive** to the values of the **initial parameters**.
- ▶ For feedforward neural networks, it is important to **initialize** all **weights** to small random values.
- ▶ The biases may be **initialized** to zero or to small positive values.

# Training Feedforward Neural Networks

- ▶ How to train a feedforward neural network?
- ▶ For each training instance  $\mathbf{x}^{(i)}$  the algorithm does the following steps:
  1. Forward pass: make a prediction (compute  $\hat{y}^{(i)} = f(\mathbf{x}^{(i)})$ ).
  2. Measure the error (compute  $\text{cost}(\hat{y}^{(i)}, y^{(i)})$ ).
  3. Backward pass: go through each layer in reverse to measure the error contribution from each connection.
  4. Tweak the connection weights to reduce the error (update  $\mathbf{W}$  and  $\mathbf{b}$ ).
- ▶ It's called the backpropagation training algorithm

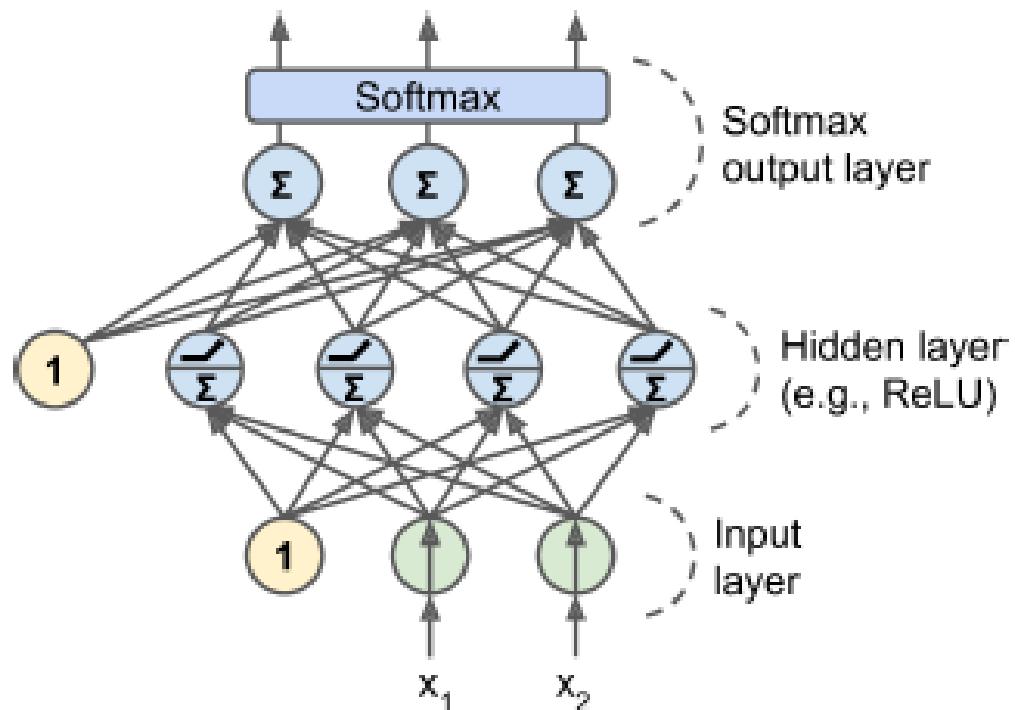


# Output Unit

- ▶ Linear units in neurons of the output layer.
  - ▶ Given  $\mathbf{h}$  as the output of neurons in the layer before the output layer.
  - ▶ Each neuron  $j$  in the output layer produces  $\hat{y}_j = \mathbf{w}_j^T \mathbf{h} + b_j$ .
  - ▶ Minimizing the cross-entropy is then equivalent to minimizing the mean squared error.
- 
- ▶ Sigmoid units in neurons of the output layer (binomial classification).
  - ▶ Given  $\mathbf{h}$  as the output of neurons in the layer before the output layer.
  - ▶ Each neuron  $j$  in the output layer produces  $\hat{y}_j = \sigma(\mathbf{w}_j^T \mathbf{h} + b_j)$ .
  - ▶ Minimizing the cross-entropy.

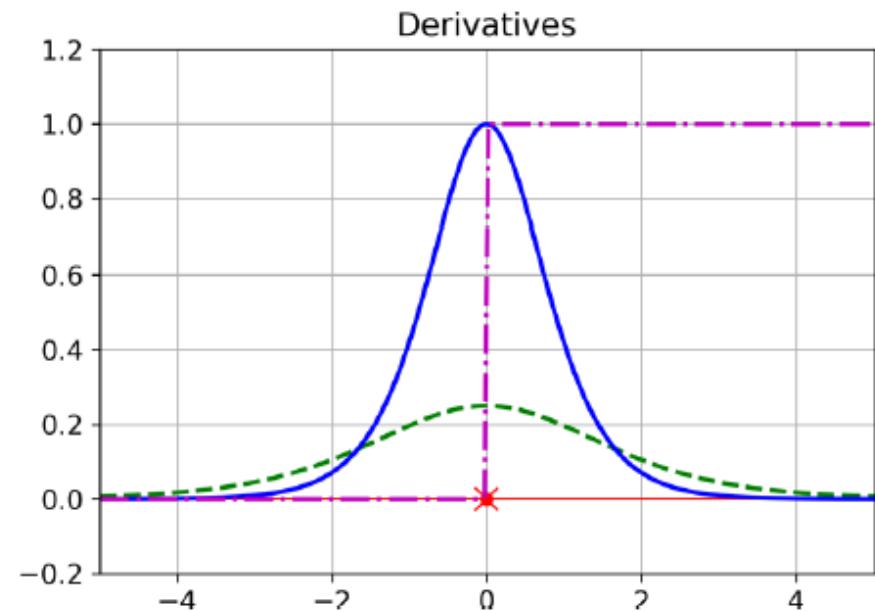
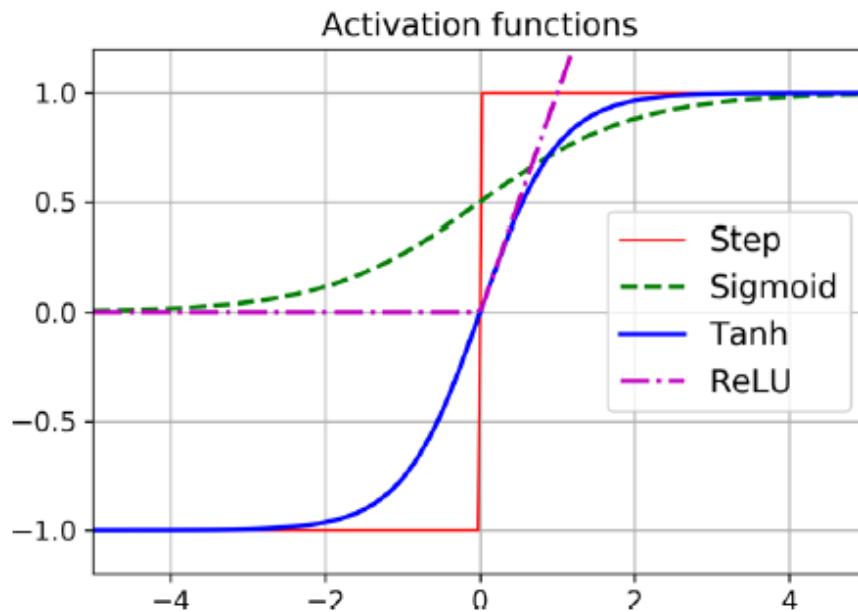
# Output Unit

- ▶ Softmax units in neurons of the output layer (multinomial classification).
- ▶ Given  $\mathbf{h}$  as the output of neurons in the layer before the output layer.
- ▶ Each neuron  $j$  in the output layer produces  $\hat{y}_j = \text{softmax}(\mathbf{w}_j^T \mathbf{h} + b_j)$ .
- ▶ Minimizing the cross-entropy.



# Hidden Units

- In order for the backpropagation algorithm to work properly, we need to replace the step function with other activation functions. Why?
- Alternative activation functions:
  - Logistic function (sigmoid):  $\sigma(z) = \frac{1}{1+e^{-z}}$
  - Hyperbolic tangent function:  $\tanh(z) = 2\sigma(2z) - 1$
  - Rectified linear units (ReLUs):  $\text{ReLU}(z) = \max(0, z)$



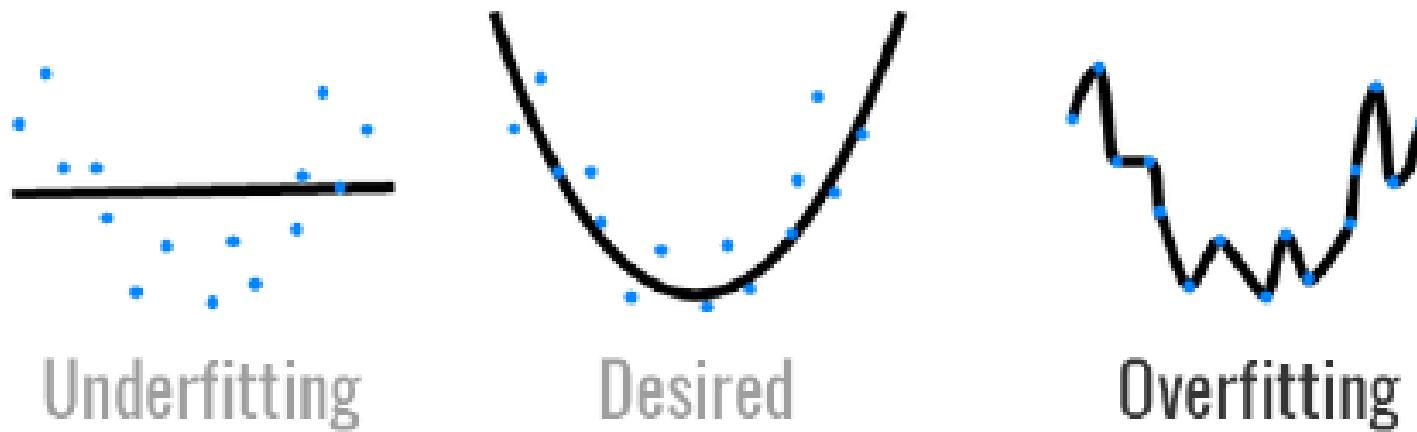
# Challenges of Training Deep NN

- ▶ Overfitting: risk of overfitting a model with large number of parameters.
- ▶ Vanishing/exploding gradients: hard to train lower layers.
- ▶ Training speed: slow training with large networks.



# High Degree of Freedom and Overfitting

- ▶ With large number of parameters, a network has a **high** degree of freedom.
- ▶ It can fit a huge variety of **complex datasets**.
- ▶ This **flexibility** also means that it is prone to overfitting on training set.
- ▶ **Regularization**: a way to **reduce** the risk of overfitting.
- ▶ It **reduces** the degree of freedom a model.

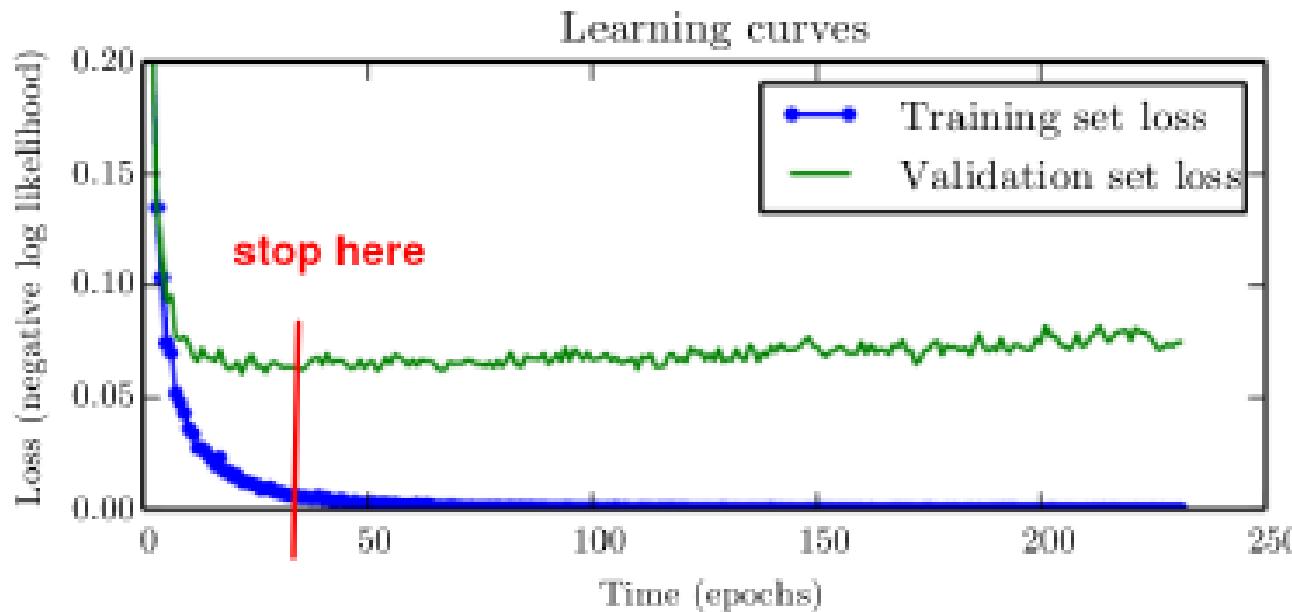


# Avoiding Overfitting Through Regularization

- ▶ Early stopping
- ▶  $\ell_1$  and  $\ell_2$  regularization
- ▶ Max-norm regularization
- ▶ Dropout
- ▶ Data augmentation

# Early Stopping

- ▶ As the training steps go by, its prediction error on the training/validation set naturally goes down.
- ▶ After a while the validation error stops decreasing and starts to go back up.
  - The model has started to overfit the training data.
- ▶ In the early stopping, we stop training when the validation error reaches a minimum.



# L1 and L2 Regularization

- ▶ Penalize large values of weights  $w_j$ .  $\tilde{J}(w) = J(w) + \lambda R(w)$
- ▶ Two questions:
  1. How should we define  $R(w)$ ?
  2. How do we determine  $\lambda$ ?
- ▶ L1 regression:  $R(w) = \lambda \sum_{i=1}^n |w_i|$  is added to the cost function.  
$$\tilde{J}(w) = J(w) + \lambda \sum_{i=1}^n |w_i|$$
- ▶ L2 regression:  $R(w) = \lambda \sum_{i=1}^n w_i^2$  is added to the cost function.  
$$\tilde{J}(w) = J(w) + \lambda \sum_{i=1}^n w_i^2$$

# Max-Norm Regularization

- ▶ Max-norm regularization: constrains the weights  $w_j$  of the incoming connections for each neuron  $j$ .

- Prevents them from getting too large.

- ▶ After each training step, clip  $w_j$  as below:

$$w_j \leftarrow w_j \frac{r}{\|w_j\|_2}$$

- ▶ We have  $\|w_j\|_2 \leq r$ .

- $r$  is the max-norm hyperparameter

- $\|w_j\|_2 = (\sum_i w_{i,j}^2)^{\frac{1}{2}} = \sqrt{w_{1,j}^2 + w_{2,j}^2 + \dots + w_{n,j}^2}$

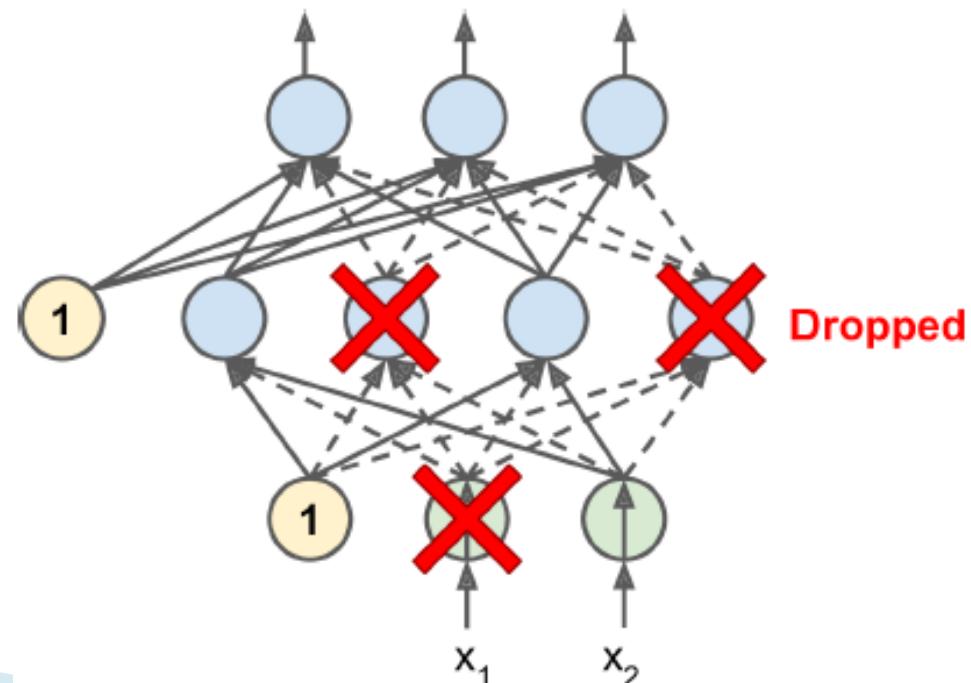
# Dropout

- ▶ Would a **company** perform better if its employees were told to toss a coin every morning to decide **whether or not to go to work?**



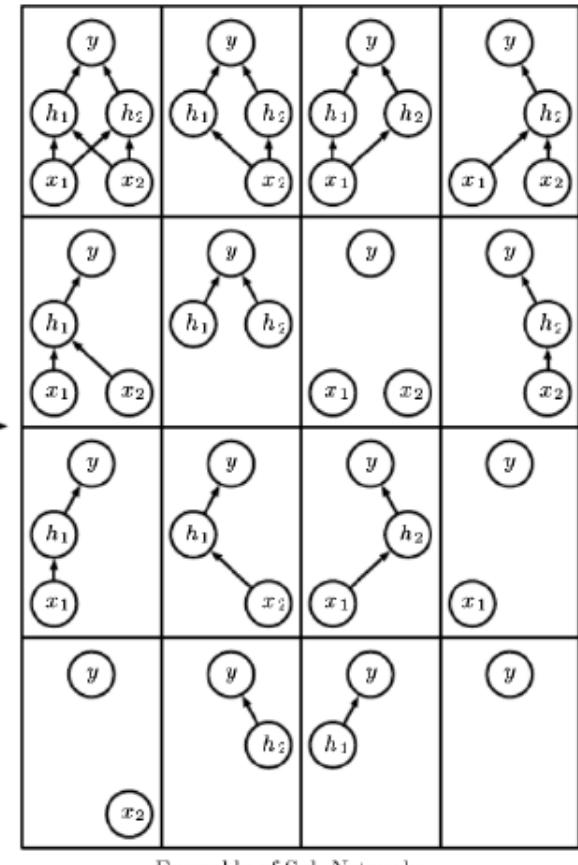
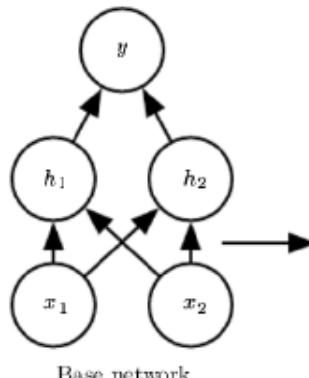
# Dropout

- ▶ At each **training step**, each neuron drops out temporarily with a **probability p**.
  - The **hyperparameter p** is called the **dropout rate**.
  - A neuron will be **entirely ignored** during **this training step**.
  - It may be **active** during the **next step**.
  - Exclude the **output neurons**.
- ▶ After training, neurons **don't get dropped** anymore.



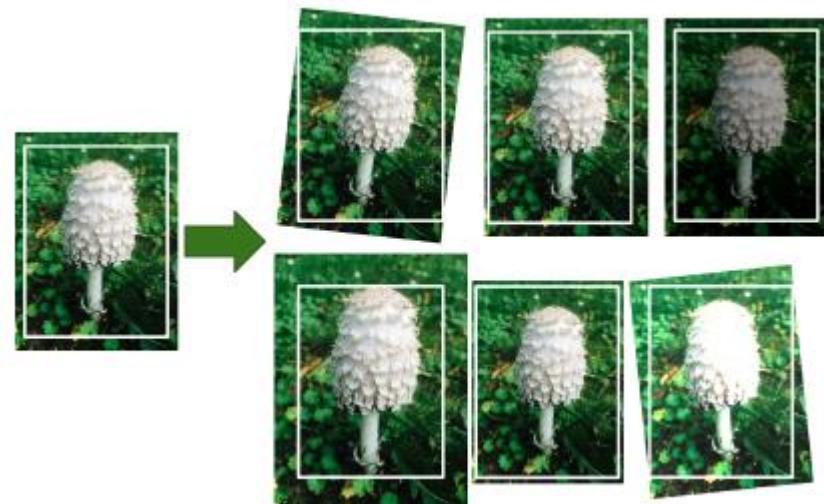
# Dropout

- ▶ Each neuron can be either present or absent.
- ▶  $2^N$  possible networks, where  $N$  is the total number of droppable neurons.
  - $N = 4$  in this figure.



# Data Augmentation

- ▶ One way to make a model **generalize better** is to train it on more data.
- ▶ This will **reduce overfitting**.
- ▶ Create **fake data** and add it to the **training set**.
  - E.g., in an **image classification** we can slightly shift, rotate and resize an image.
  - Add the resulting **pictures** to the training set.



# Vanishing/Exploding Gradients Problem

- ▶ The backpropagation goes from **output** to **input** layer, and propagates the error gradient on the way.

$$\mathbf{w}^{(\text{next})} = \mathbf{w} - \eta \frac{\partial J(\mathbf{w})}{\partial \mathbf{w}}$$

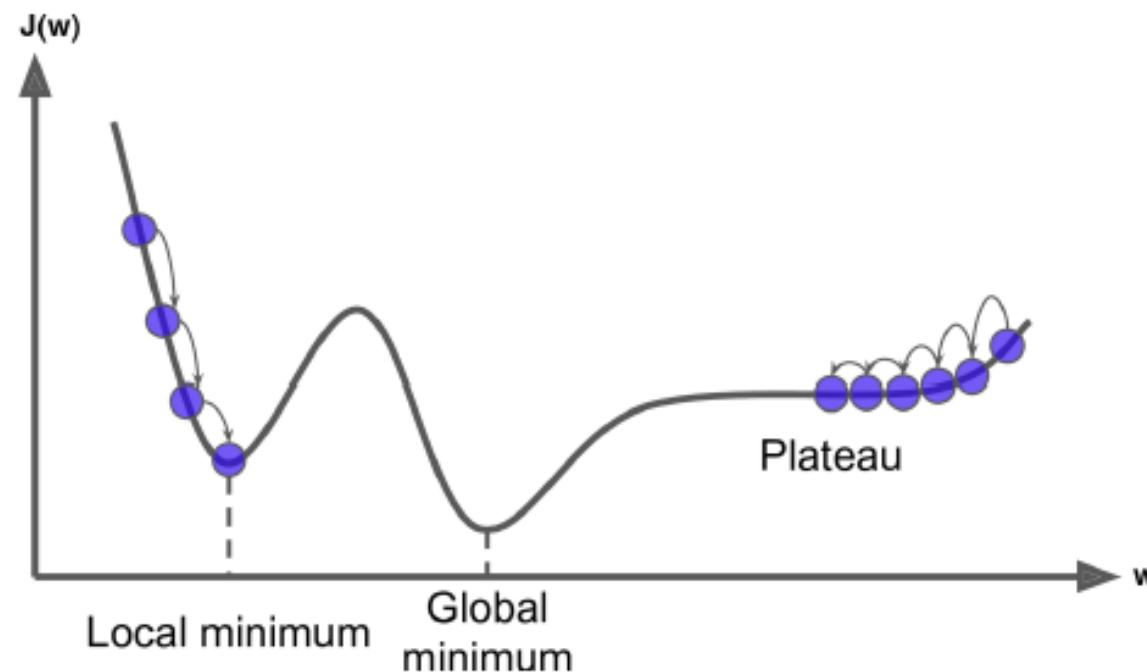
- ▶ Gradients often get **smaller and smaller** as the algorithm progresses **down to the lower layers**.
- ▶ As a result, the gradient descent update leaves the lower layer connection weights virtually **unchanged**.
- ▶ This is called the **vanishing gradients** problem.

# Overcoming the Vanishing Gradient

- ▶ Parameter initialization strategies
- ▶ Nonsaturating activation function
- ▶ Batch normalization
- ▶ Gradient clipping

# Parameter Initialization Strategies

- ▶ The non-linearity of a neural network causes the cost functions to become **non-convex**.
- ▶ The stochastic gradient descent on **non-convex** cost functions performs is **sensitive** to the values of the **initial parameters**.
- ▶ Designing initialization strategies is a **difficult task**.



# Parameter Initialization Strategies

- ▶ The initial parameters need to **break symmetry** between **different units**.
- ▶ Two **hidden units** with the same activation function connected to the **same inputs**, must have different initial parameters.
  - The goal of having each unit compute a different function.
- ▶ It motivates **random initialization** of the parameters.
  - Typically, we set the **biases** to **constants**, and initialize only the **weights** **randomly**.
- ▶ We need the signals to flow properly in **both directions**.
- ▶ The **Xavier initialization** proposed that:
  - The **variance** of the **outputs** of each layer to be **equal** to the **variance** of its **inputs**.
  - The **gradients** to have **equal variance** **before and after** flowing through a layer in the **reverse direction**.

# Parameter Initialization Strategies

- Based on the Xavier initialization, the weights are initialized using normal distribution with mean 0 and the following standard deviation.

- For the sigmoid activation function:

$$\sigma = \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$$

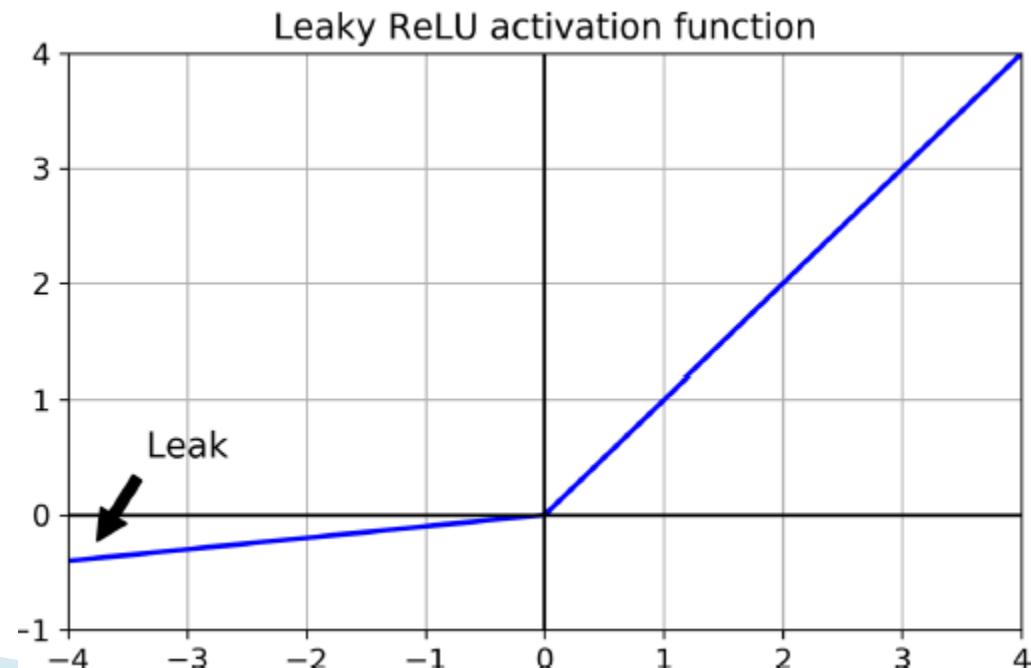
- For the ReLU activation function:

$$\sigma = \sqrt{2} \sqrt{\frac{2}{n_{\text{inputs}} + n_{\text{outputs}}}}$$

- $n_{\text{inputs}}$  and  $n_{\text{outputs}}$  are the number of input and output connections for the layer whose weights are being initialized.

# No saturating Activation Functions

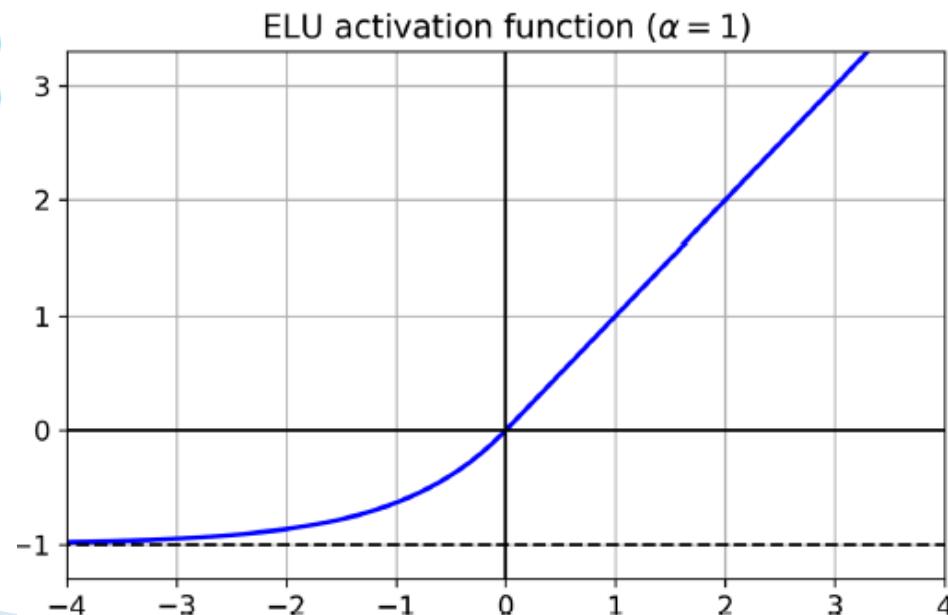
- ▶  $\text{ReLU}(z) = \max(0, z)$
- ▶ The **dying ReLUs** problem.
  - During **training**, some neurons **stop outputting anything other than 0**.
  - E.g., when the **weighted sum of the neuron's inputs** is negative, it starts outputting 0.
- ▶ Use **leaky ReLU** instead:  $\text{LeakyReLU}_\alpha(z) = \max(\alpha z, z)$ .
  - $\alpha$  is the **slope** of the function for  $z < 0$ .



# No saturating Activation Functions

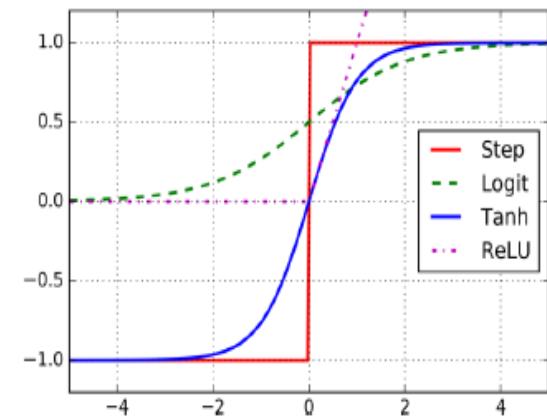
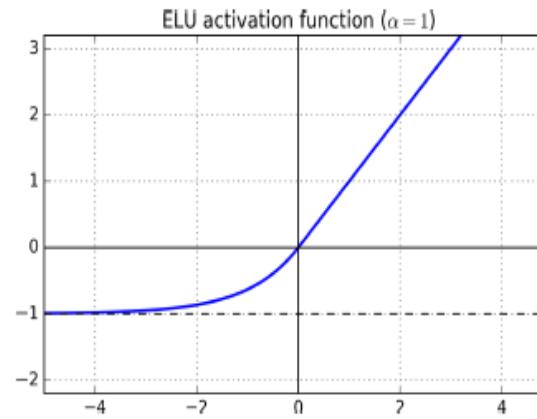
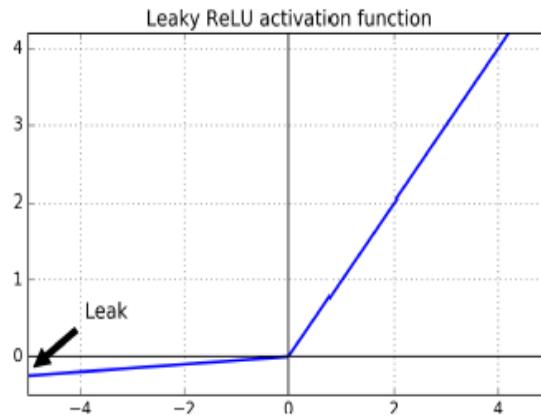
- ▶ Randomized Leaky ReLU (RReLU)
  - $\alpha$  is picked randomly during training, and it is fixed during testing.
- ▶ Parametric Leaky ReLU (PReLU)
  - Learn  $\alpha$  during training (instead of being a hyperparameter).
- ▶ Exponential Linear Unit (ELU)

$$\text{ELU}_\alpha(z) = \begin{cases} \alpha(\exp(z) - 1) & \text{if } z < 0 \\ z & \text{if } z \geq 0 \end{cases}$$



# No saturating Activation Functions

- ▶ Which activation function should we use?
- ▶ In general logistic < tanh < ReLU < leaky ReLU (and its variants) < ELU
- ▶ If you care about runtime performance, then leaky ReLUs works better than ELUs.



# Batch Normalization

- ▶ The gradient tells how to update each parameter, under the assumption that the other layers do not change.
  - In practice, we update all of the layers simultaneously.
  - However, unexpected results can happen.
- ▶ Batch normalization makes the learning of layers in the network more independent of each other.
  - It is a technique to address the problem that the distribution of each layer's inputs changes during training, as the parameters of the previous layers change.
- ▶ The technique consists of adding an operation in the model just before the activation function of each layer.

# Batch Normalization

- ▶ It's zero-centering and normalizing the inputs, then scaling and shifting the result.
  - Estimates the inputs' mean and standard deviation of the current mini-batch.

$$\mu_B = \frac{1}{m_B} \sum_{i=1}^{m_B} x^{(i)}$$

$$\sigma_B^2 = \frac{1}{m_B} \sum_{i=1}^{m_B} (x^{(i)} - \mu_B)^2$$

- ▶  $\mu_B$ : the empirical mean, evaluated over the whole mini-batch B.
- ▶  $\sigma_B$ : the empirical standard deviation, also evaluated over the whole mini-batch.
- ▶  $m_B$ : the number of instances in the mini-batch.

# Batch Normalization

$$\hat{x}^{(i)} = \frac{x^{(i)} - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}}$$

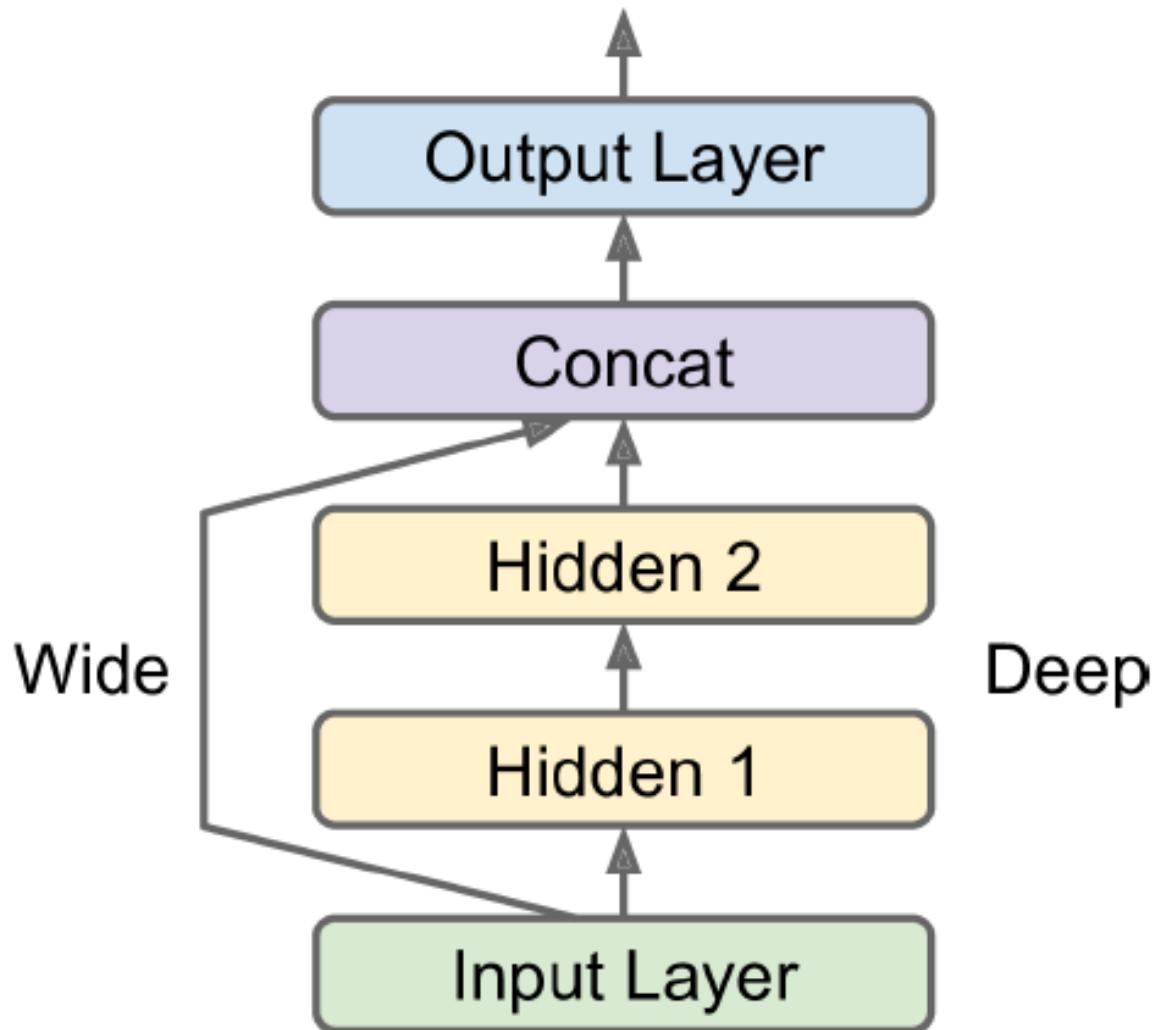
$$z^{(i)} = \gamma \hat{x}^{(i)} + \beta$$

- ▶  $\hat{x}^{(i)}$ : the zero-centered and normalized input.
- ▶  $\gamma$ : the scaling parameter for the layer.
- ▶  $\beta$ : the shifting parameter (offset) for the layer.
- ▶  $\epsilon$ : a tiny number to avoid division by zero.
- ▶  $z^{(i)}$ : the output of the BN operation, which is a scaled and shifted version of the inputs.

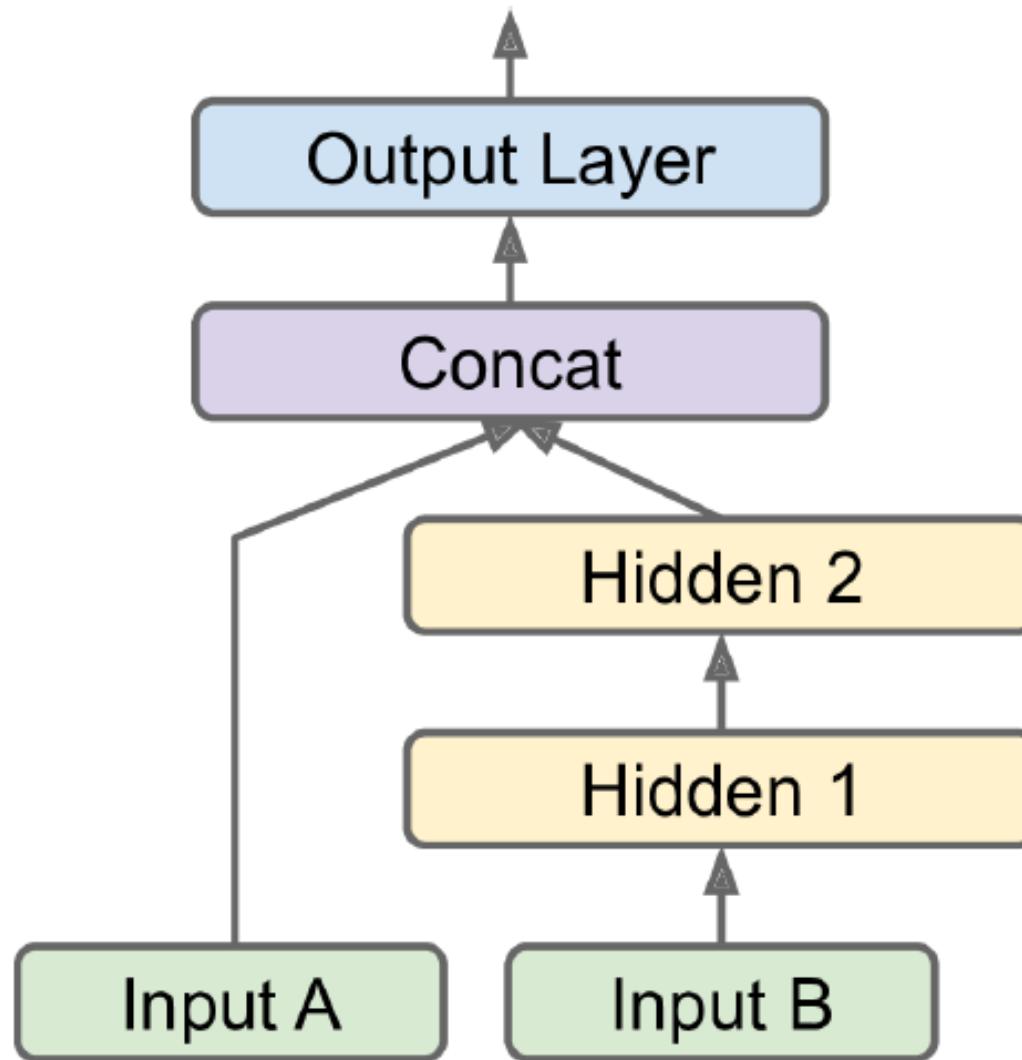
# Gradient Clipping

- ▶ Gradient clipping: clip the gradients during backpropagation so that they never exceed some threshold.
- ▶ In TensorFlow, the optimizer's `minimize()` function takes care of:
  1. Compute the gradients with `compute_gradients()`
  2. Apply the processed gradients with `apply_gradients()`
- ▶ To enable the gradient clipping, you must instead of calling `minimize()`, call:
  1. Compute the gradients with `compute_gradients()`
  2. Process the gradients as you wish.
  3. Apply the processed gradients with `apply_gradients()`

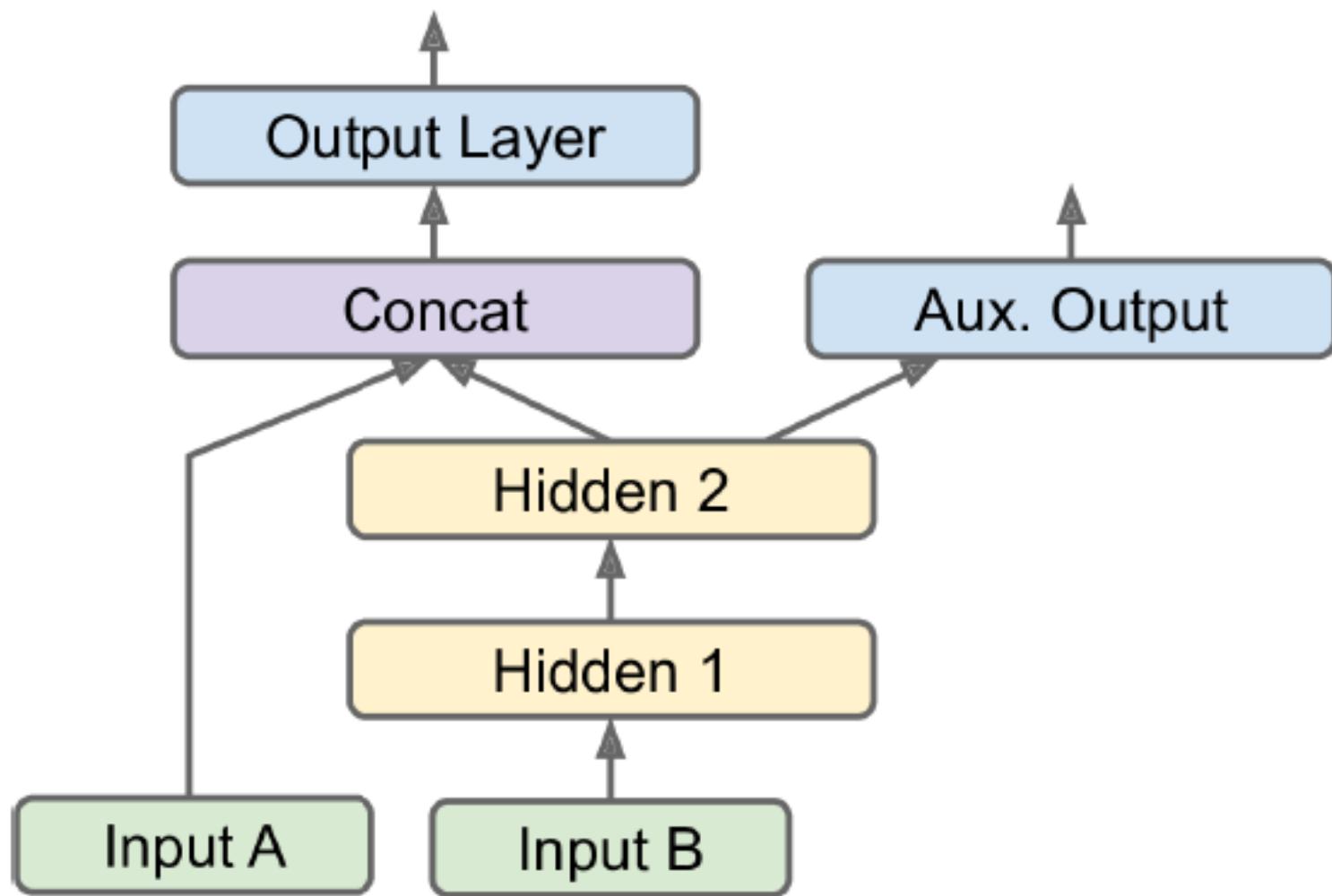
# Building Complex Models



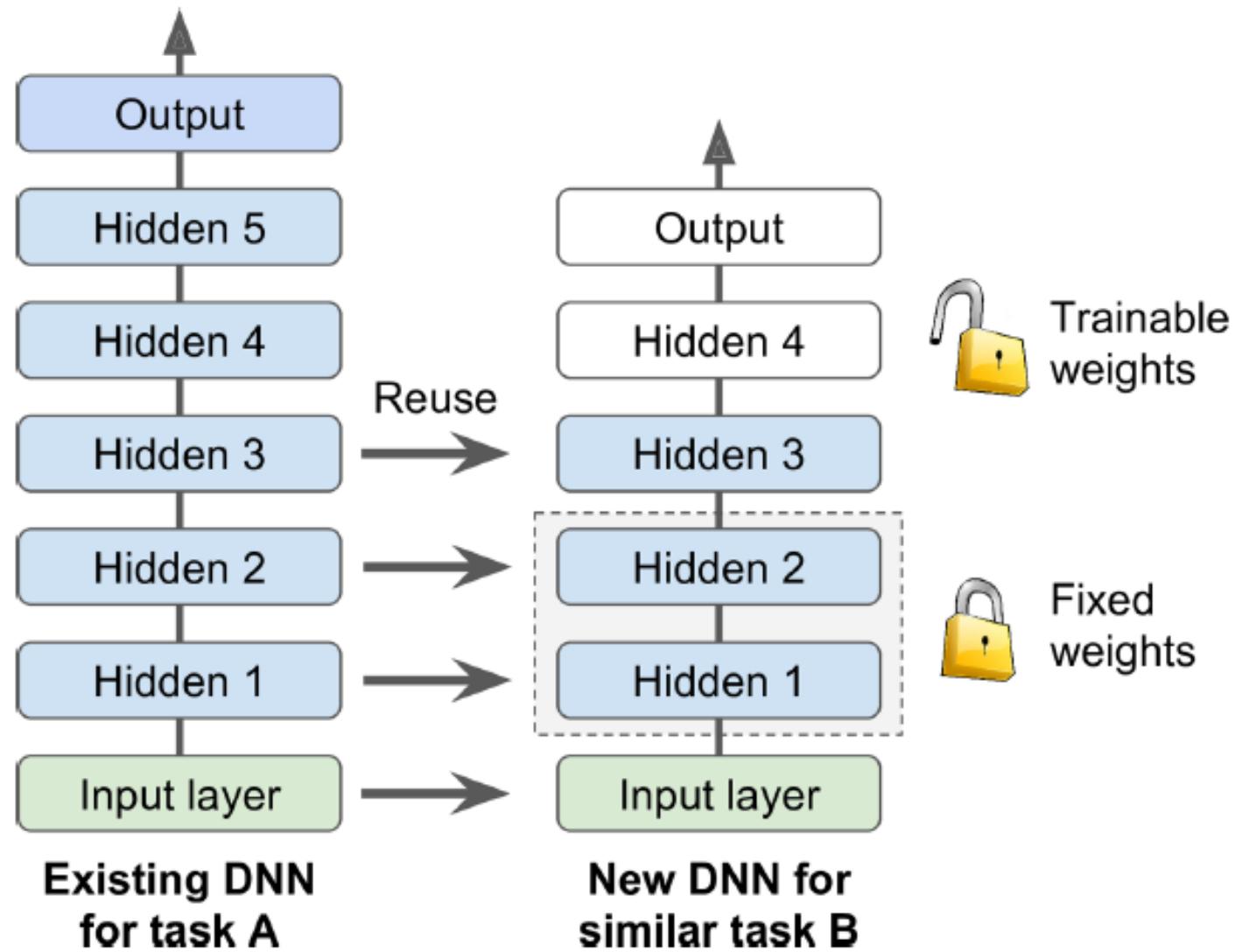
# Handling multiple inputs



# Handling multiple inputs/outputs



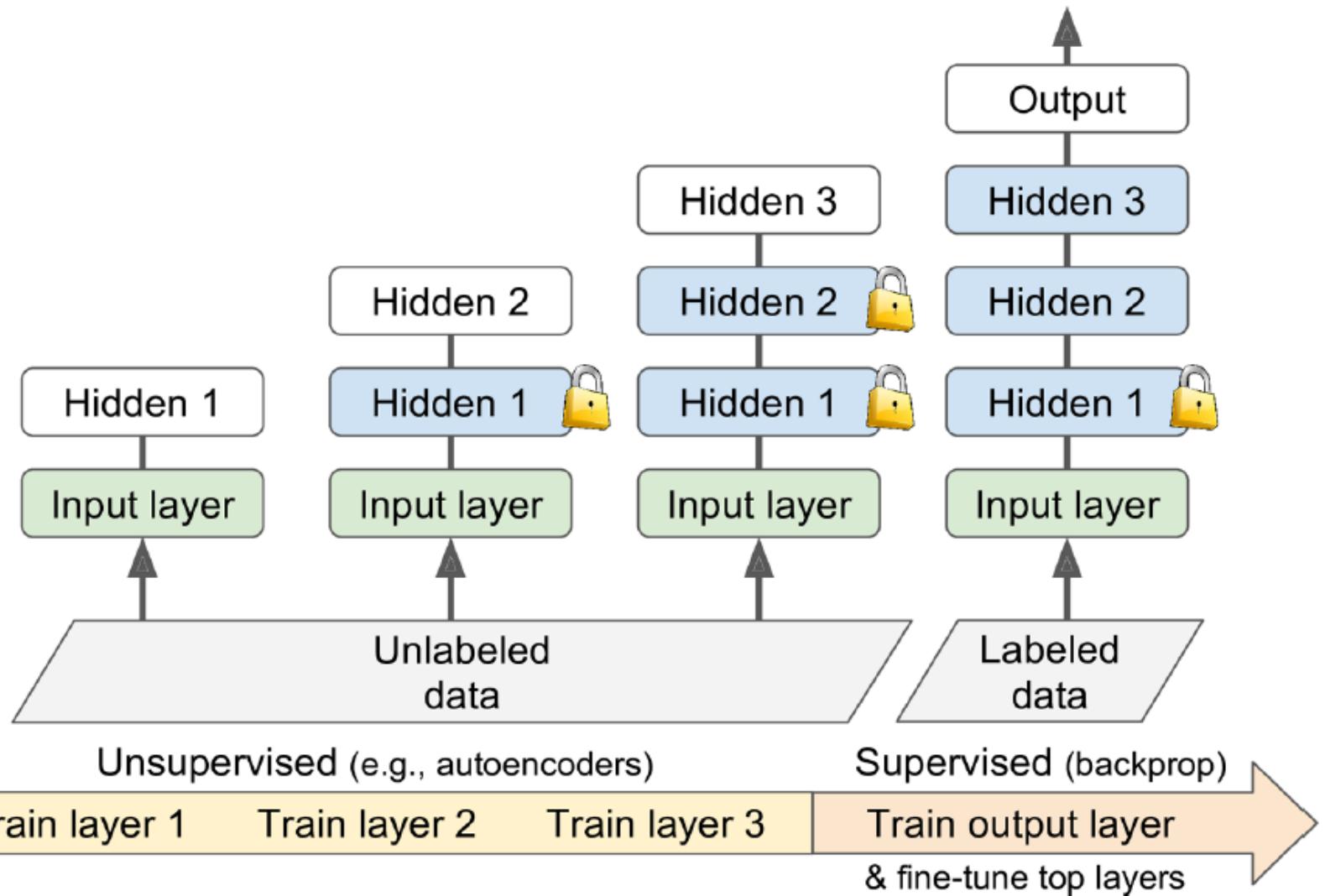
# Reusing Pretrained Layers



# Transfer Learning

	<b>very similar dataset</b>	<b>very different dataset</b>
<b>very little data</b>	Use Linear Classifier on top layer	You're in trouble... Try linear classifier from different stages
<b>quite a lot of data</b>	Finetune a few layers	Finetune a larger number of layers

# unsupervised training

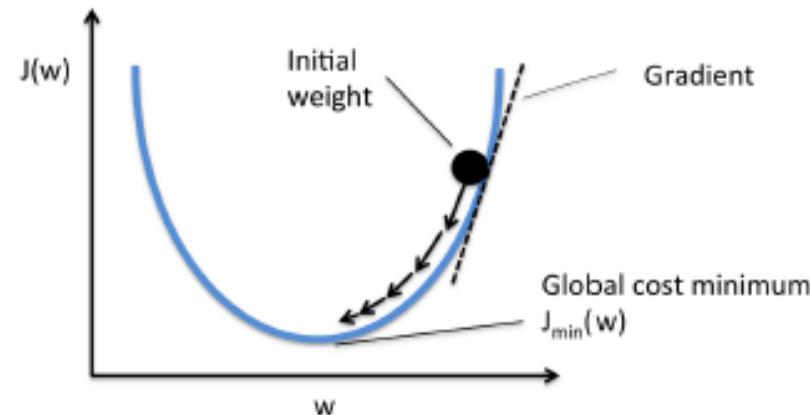


# Optimization Algorithms

- ▶ Momentum
- ▶ Nesterov momentum
- ▶ AdaGrad
- ▶ RMSProp
- ▶ Adam Optimization

# Momentum

- ▶ Momentum is a concept from physics: an object in motion will have a tendency to keep moving.
- ▶ It measures the resistance to change in motion.
  - The higher momentum an object has, the harder it is to stop it.
- ▶ This is the very simple idea behind momentum optimization.
- ▶ We can see the change in the parameters  $w$  as motion:  $w_i^{(\text{next})} = w_i - \eta \frac{\partial J(w)}{\partial w_i}$
- ▶ We can thus use the concept of momentum to give the update process a tendency to keep moving in the same direction.
- ▶ It can help to escape from bad local minima pits.



# Momentum

- ▶ Momentum optimization cares about what previous gradients were.
- ▶ At each iteration, it adds the local gradient to the momentum vector  $\mathbf{m}$ .

$$m_i = \beta m_i + \eta \frac{\partial J(\mathbf{w})}{\partial w_i}$$

- ▶  $\beta$  is called momentum, and it is between 0 and 1.
- ▶ Updates the weights by subtracting this momentum vector.

$$w_i^{(\text{next})} = w_i - m_i$$

# Nesterov Momentum

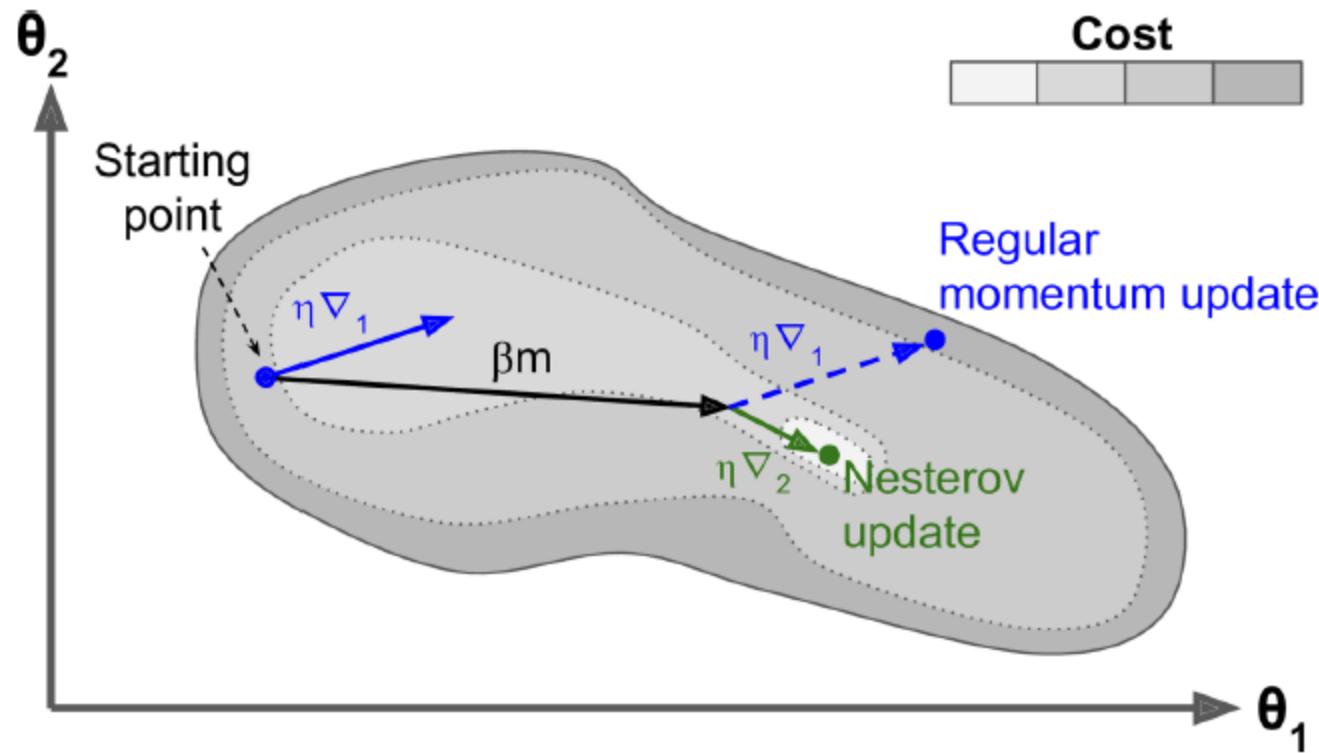
- ▶ Nesterov Momentum is a small variant to Momentum optimization.
- ▶ Faster than vanilla Momentum optimization.
- ▶ Measure the gradient of the cost function slightly ahead in the direction of the momentum (not at the local position).

$$m_i = \beta m_i + \eta \frac{\partial J(w + \beta m)}{\partial w_i}$$

$$w_i^{(\text{next})} = w_i - m_i$$

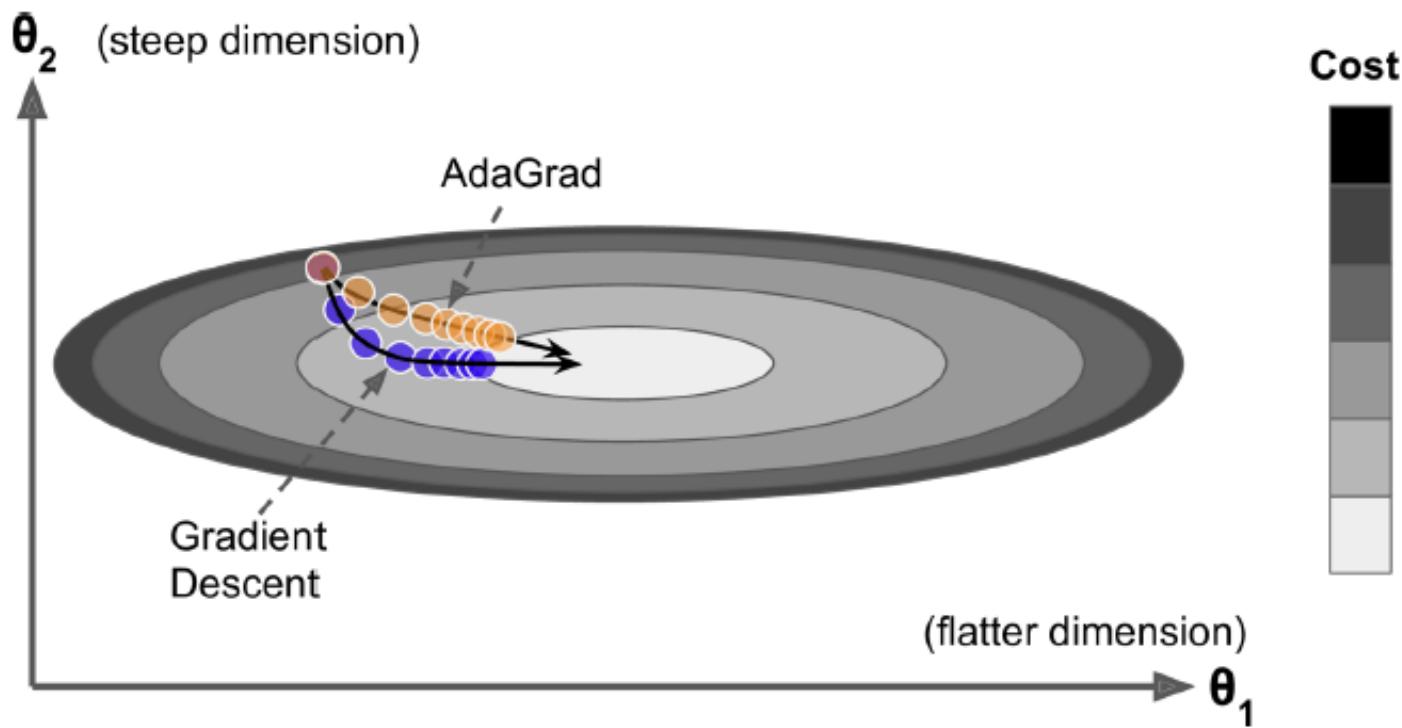
# Nesterov Momentum

- ▶  $\nabla_1$  represents the gradient of the cost function measured at the starting point  $w$ , and  $\nabla_2$  represents the gradient at the point located at  $w + \beta m$ .



# Ada Grad

- ▶ AdaGrad keeps track of a learning rate for each parameter.
- ▶ Adapts the learning rate over time (adaptive learning rate).



# Ada Grad

- ▶ For each feature  $w_i$ , we do the following steps:

$$s_i = s_i + \left( \frac{\partial J(\mathbf{w})}{\partial w_i} \right)^2$$

$$w_i^{(\text{next})} = w_i - \frac{\eta}{\sqrt{s_i + \epsilon}} \frac{\partial J(\mathbf{w})}{\partial w_i}$$

- ▶ Parameters with large partial derivative of the cost have a rapid decrease in their learning rate.
- ▶ Parameters with small partial derivatives have a small decrease in their learning rate.

# RMS Prop

- ▶ AdaGrad often stops too early when training neural networks.
- ▶ The learning rate gets scaled down so much that the algorithm ends up stopping entirely before reaching the global optimum.
- ▶ The RMSProp fixed the AdaGrad problem.
- ▶ It is like the AdaGrad problem, but accumulates only the gradients from the most recent iterations (not from the beginning of training).
- ▶ For each feature  $w_i$ , we do the following steps:

$$s_i = \beta s_i + (1 - \beta) \left( \frac{\partial J(w)}{\partial w_i} \right)^2$$

$$w_i^{(\text{next})} = w_i - \frac{\eta}{\sqrt{s_i + \epsilon}} \frac{\partial J(w)}{\partial w_i}$$

# Adam and Nadam Optimization

- ▶ Adam (Adaptive moment estimation) combines the ideas of Momentum optimization and RMSProp.
- ▶ Like Momentum optimization, it keeps track of an exponentially decaying average of past gradients.
- ▶ Like RMSProp, it keeps track of an exponentially decaying average of past squared gradients.

# Adam and Nadam Optimization

$$1. \quad \mathbf{m}^{(\text{next})} = \beta_1 \mathbf{m} + (1 - \beta_1) \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$2. \quad \mathbf{s}^{(\text{next})} = \beta_2 \mathbf{s} + (1 - \beta_2) \nabla_{\mathbf{w}} J(\mathbf{w}) \otimes \nabla_{\mathbf{w}} J(\mathbf{w})$$

$$3. \quad \mathbf{m}^{(\text{next})} = \frac{\mathbf{m}}{1 - \beta_1^T}$$

$$4. \quad \mathbf{s}^{(\text{next})} = \frac{\mathbf{s}}{1 - \beta_2^T}$$

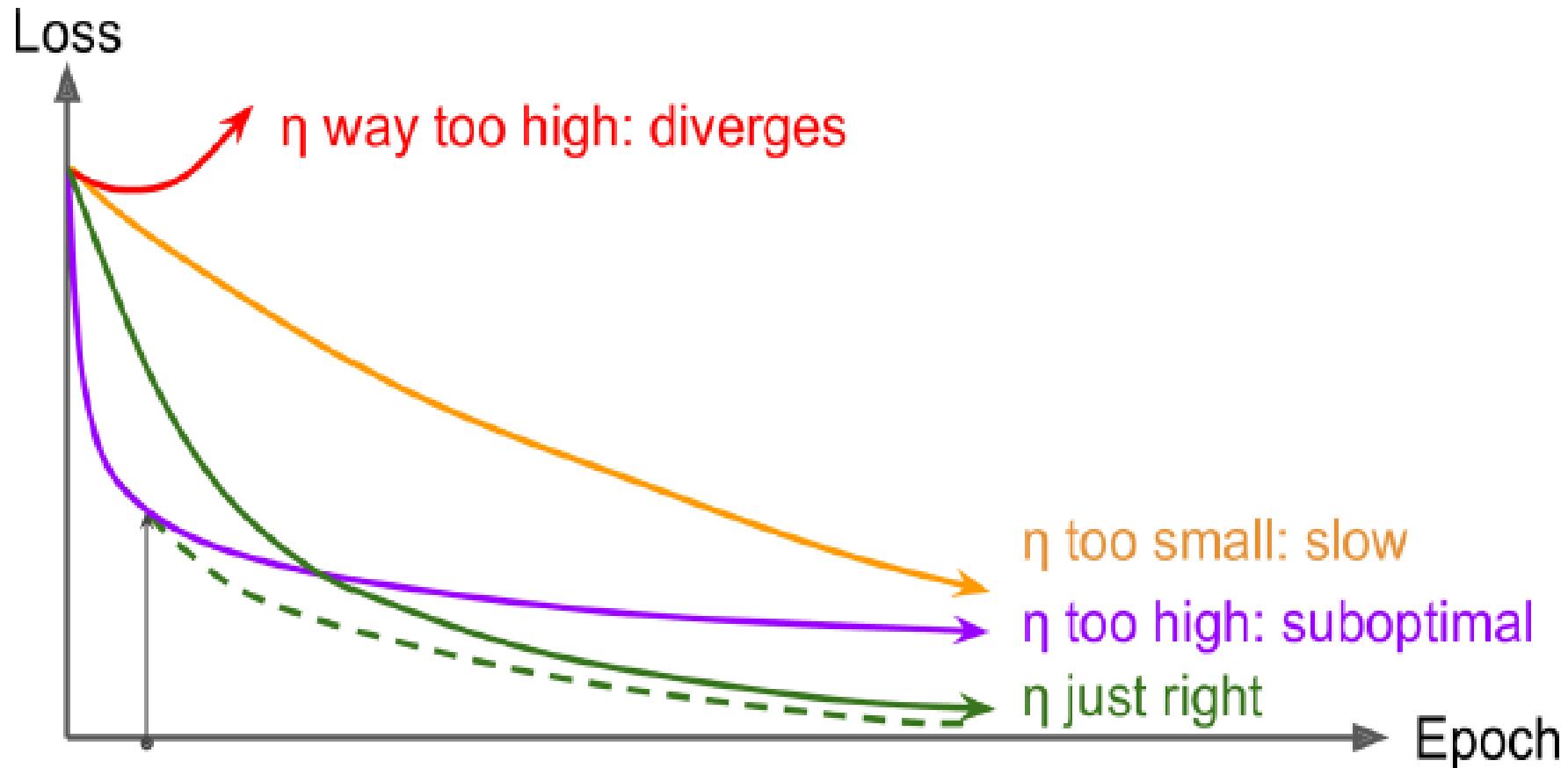
$$5. \quad \mathbf{w}^{(\text{next})} = \mathbf{w} - \eta \mathbf{m} \oslash \sqrt{\mathbf{s} + \epsilon}$$

- ▶  $\otimes$  and  $\oslash$  represents the represents the element-wise multiplication and division.
- ▶ Steps 1, 2, and 5: similar to both Momentum optimization and RMSProp.
- ▶ Steps 3 and 4: since  $\mathbf{m}$  and  $\mathbf{s}$  are initialized at 0, they will be biased toward 0 at the beginning of training, so these two steps will help boost  $\mathbf{m}$  and  $\mathbf{s}$  at the beginning of training.

# Optimizer comparison

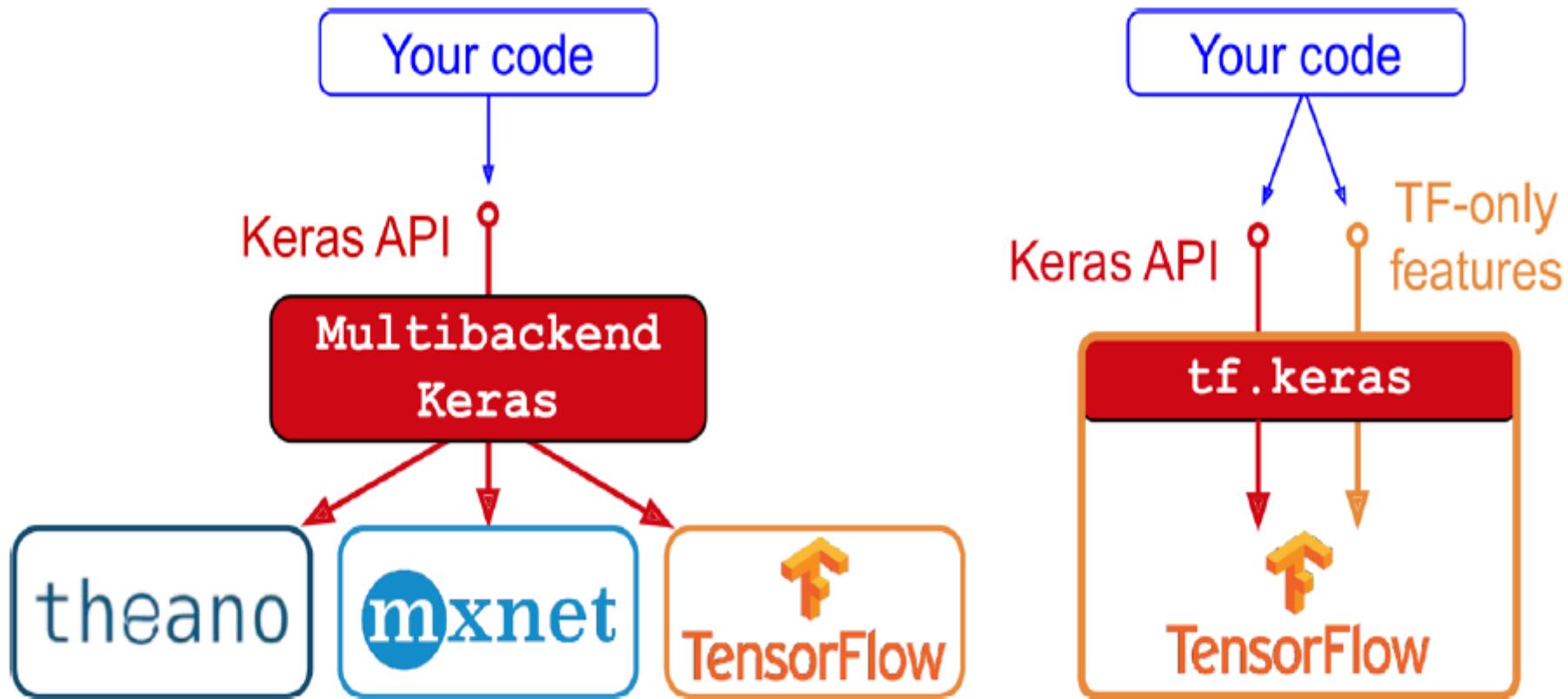
Class	Convergence speed	Convergence quality
SGD	*	***
SGD(momentum=...)	**	***
SGD(momentum=..., nesterov=True)	**	***
Adagrad	***	* (stops too early)
RMSprop	***	** or ***
Adam	***	** or ***
Nadam	***	** or ***
AdaMax	***	** or ***

# Learning Rate Scheduling

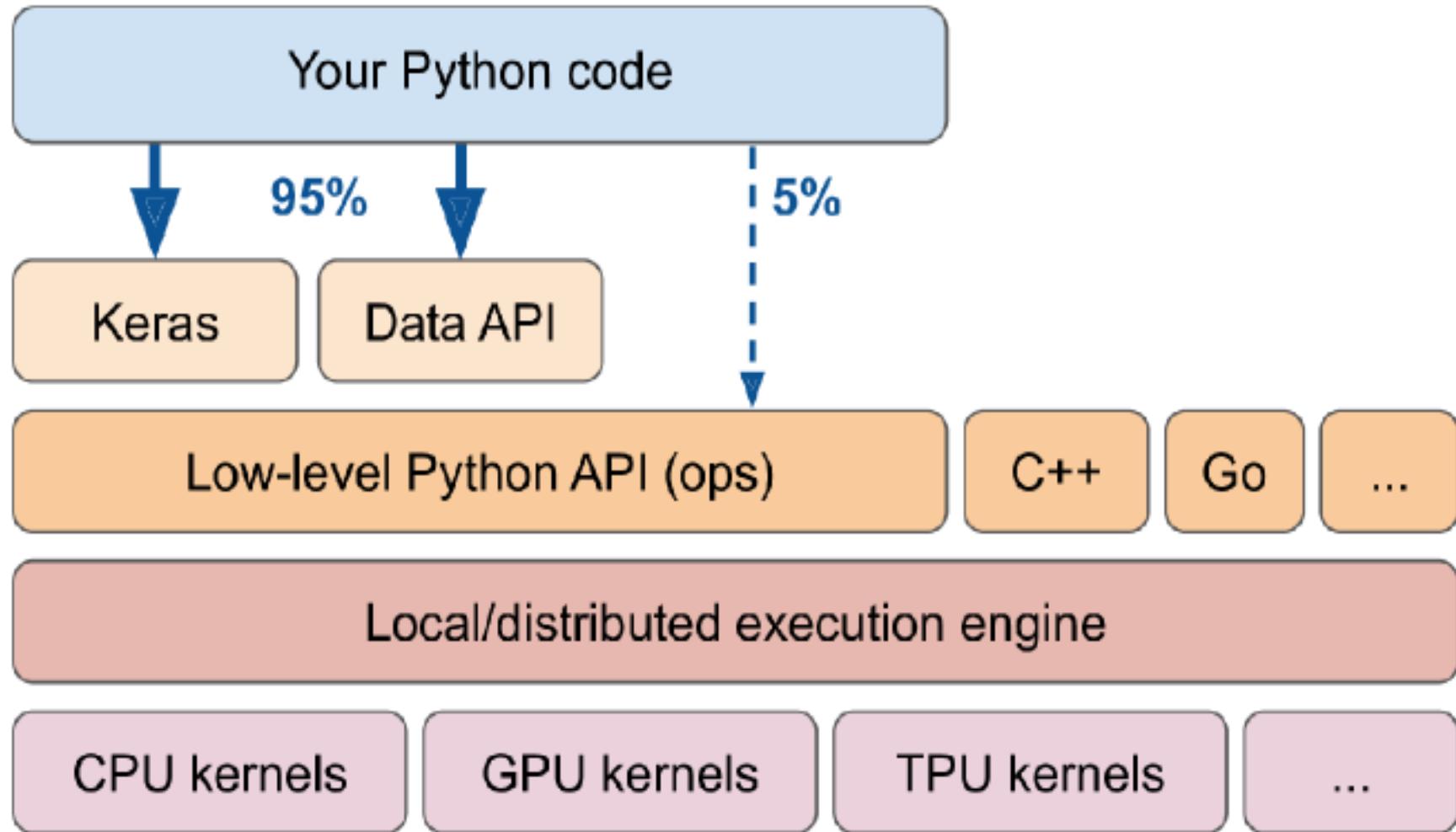


Start with a high learning rate then reduce it: perfect!

# *TensorFlow's* tf.keras

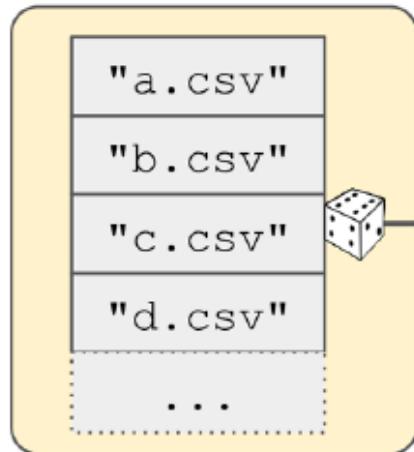


# *TensorFlow's architecture*

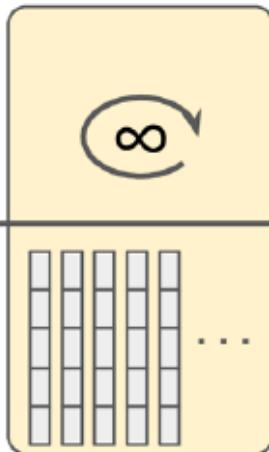


# Loading and preprocessing data

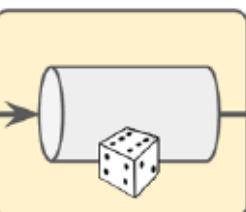
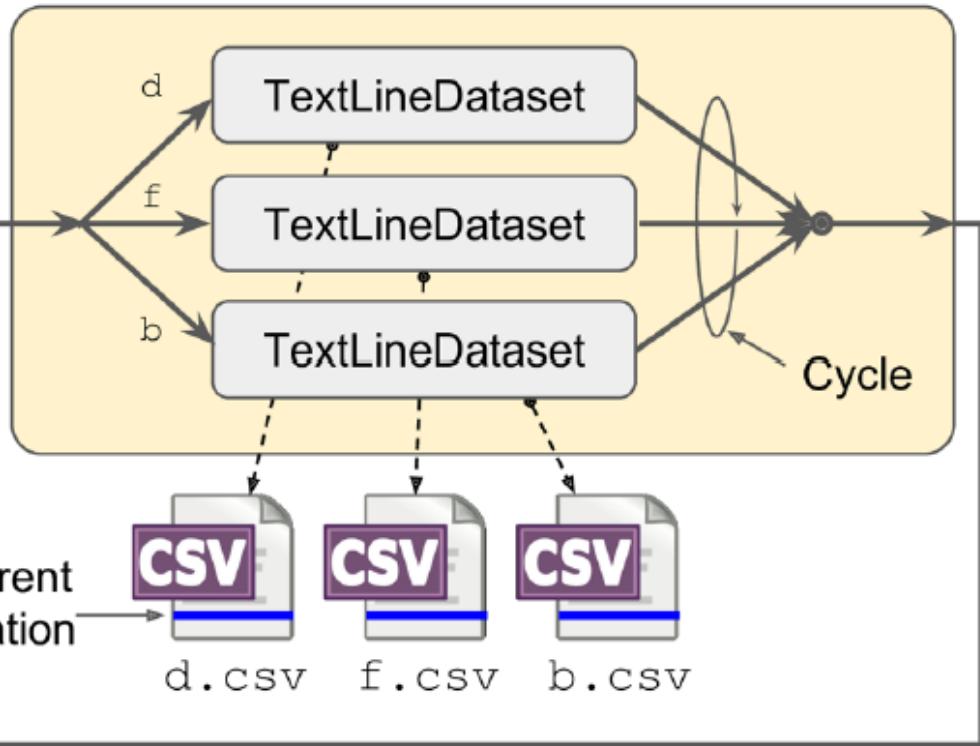
`list_files()`



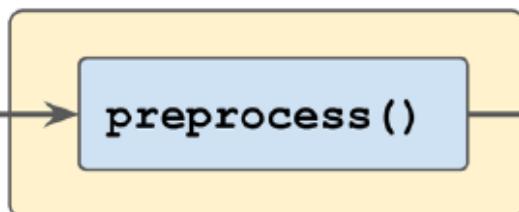
`repeat()`



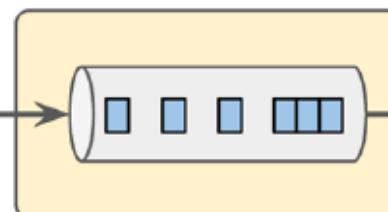
`interleave()`



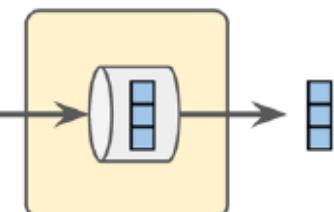
`shuffle()`



`map()`



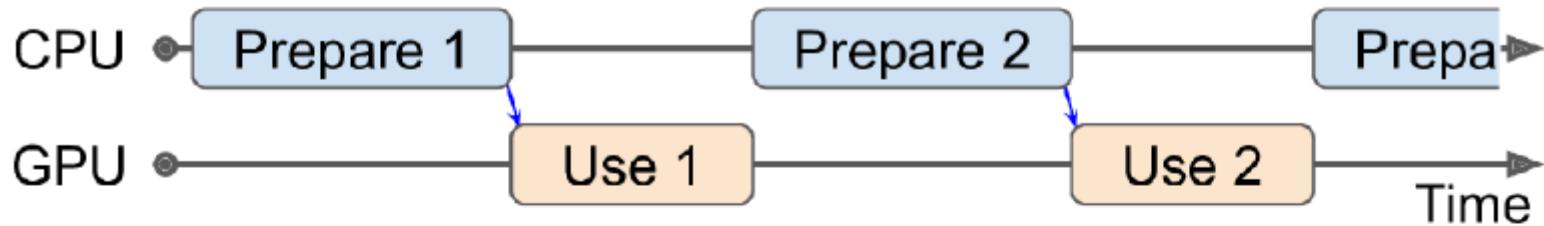
`batch()`



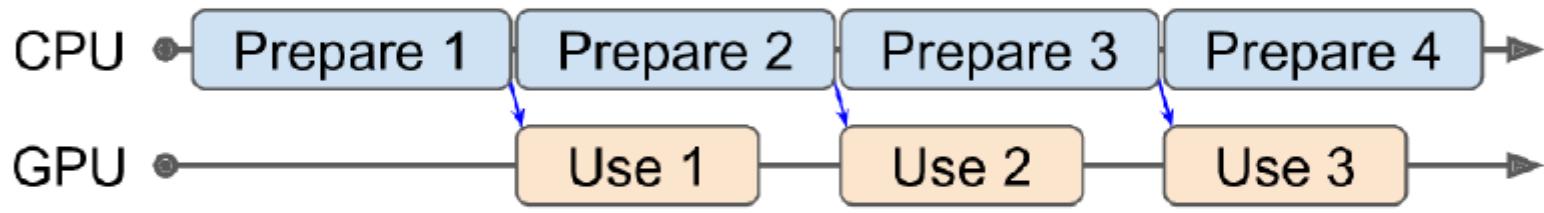
`prefetch()`

# Prefetching

Without prefetching



With prefetching



With prefetching + multithreaded loading & preprocessing

