

“dislike” (activation near zero) certain linear regions of its input space. Hence, with an appropriate loss function on the neuron’s output, we can turn a single neuron into a linear classifier:

**Binary Softmax classifier.** For example, we can interpret  $\sigma(\sum_i w_i x_i + b)$  to be the probability of one of the classes  $P(y_i = 1 \mid x_i; w)$ . The probability of the other class would be  $P(y_i = 0 \mid x_i; w) = 1 - P(y_i = 1 \mid x_i; w)$ , since they must sum to one. With this interpretation, we can formulate the cross-entropy loss as we have seen in the Linear Classification section, and optimizing it would lead to a binary Softmax classifier (also known as *logistic regression*). Since the sigmoid function is restricted to be between 0-1, the predictions of this classifier are based on whether the output of the neuron is greater than 0.5.

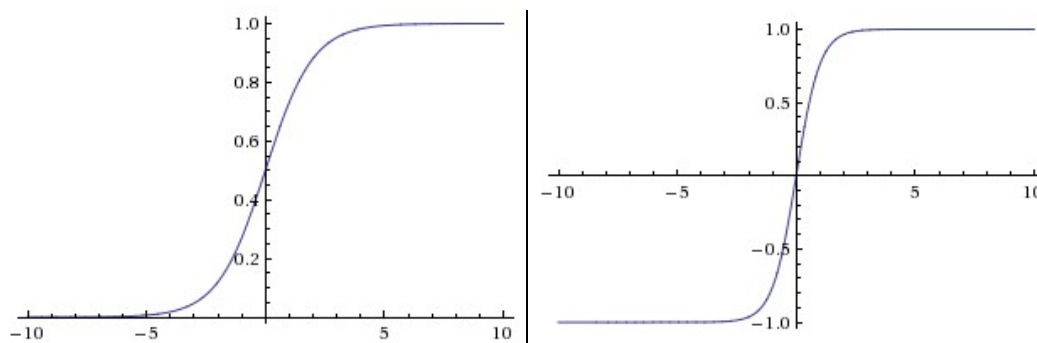
**Binary SVM classifier.** Alternatively, we could attach a max-margin hinge loss to the output of the neuron and train it to become a binary Support Vector Machine.

**Regularization interpretation.** The regularization loss in both SVM/Softmax cases could in this biological view be interpreted as *gradual forgetting*, since it would have the effect of driving all synaptic weights  $w$  towards zero after every parameter update.

*A single neuron can be used to implement a binary classifier (e.g. binary Softmax or binary SVM classifiers)*

## Commonly used activation functions

Every activation function (or *non-linearity*) takes a single number and performs a certain fixed mathematical operation on it. There are several activation functions you may encounter in practice:

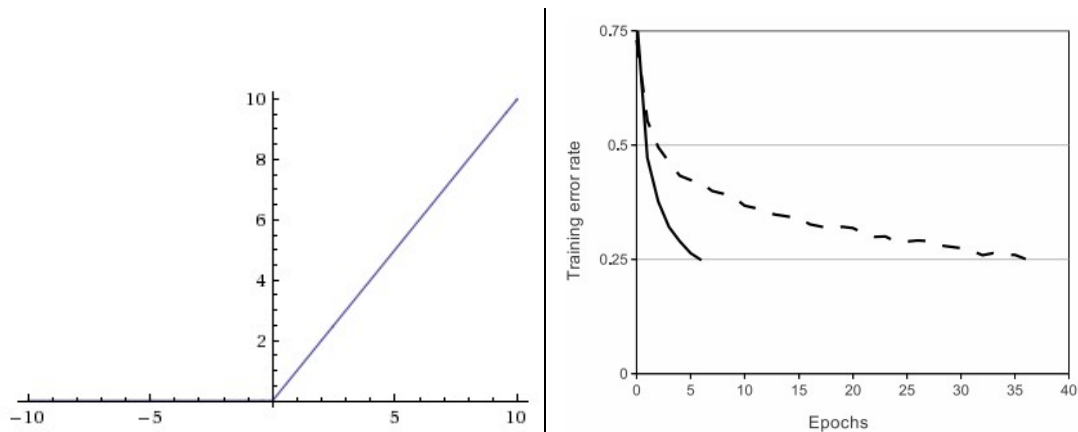


**Left:** Sigmoid non-linearity squashes real numbers to range between [0,1] **Right:** The tanh non-linearity squashes real numbers to range between [-1,1].

**Sigmoid.** The sigmoid non-linearity has the mathematical form  $\sigma(x) = 1/(1 + e^{-x})$  and is shown in the image above on the left. As alluded to in the previous section, it takes a real-valued number and “squashes” it into range between 0 and 1. In particular, large negative numbers become 0 and large positive numbers become 1. The sigmoid function has seen frequent use historically since it has a nice interpretation as the firing rate of a neuron: from not firing at all (0) to fully-saturated firing at an assumed maximum frequency (1). In practice, the sigmoid non-linearity has recently fallen out of favor and it is rarely ever used. It has two major drawbacks:

- *Sigmoids saturate and kill gradients.* A very undesirable property of the sigmoid neuron is that when the neuron’s activation saturates at either tail of 0 or 1, the gradient at these regions is almost zero. Recall that during backpropagation, this (local) gradient will be multiplied to the gradient of this gate’s output for the whole objective. Therefore, if the local gradient is very small, it will effectively “kill” the gradient and almost no signal will flow through the neuron to its weights and recursively to its data. Additionally, one must pay extra caution when initializing the weights of sigmoid neurons to prevent saturation. For example, if the initial weights are too large then most neurons would become saturated and the network will barely learn.
- *Sigmoid outputs are not zero-centered.* This is undesirable since neurons in later layers of processing in a Neural Network (more on this soon) would be receiving data that is not zero-centered. This has implications on the dynamics during gradient descent, because if the data coming into a neuron is always positive (e.g.  $x > 0$  elementwise in  $f = w^T x + b$ ), then the gradient on the weights  $w$  will during backpropagation become either all be positive, or all negative (depending on the gradient of the whole expression  $f$ ). This could introduce undesirable zig-zagging dynamics in the gradient updates for the weights. However, notice that once these gradients are added up across a batch of data the final update for the weights can have variable signs, somewhat mitigating this issue. Therefore, this is an inconvenience but it has less severe consequences compared to the saturated activation problem above.

**Tanh.** The tanh non-linearity is shown on the image above on the right. It squashes a real-valued number to the range  $[-1, 1]$ . Like the sigmoid neuron, its activations saturate, but unlike the sigmoid neuron its output is zero-centered. Therefore, in practice the *tanh non-linearity is always preferred to the sigmoid nonlinearity*. Also note that the tanh neuron is simply a scaled sigmoid neuron, in particular the following holds:  $\tanh(x) = 2\sigma(2x) - 1$ .



**Left:** Rectified Linear Unit (ReLU) activation function, which is zero when  $x < 0$  and then linear with slope 1 when  $x > 0$ . **Right:** A plot from [Krizhevsky et al. \(pdf\)](#) paper indicating the 6x improvement in convergence with the ReLU unit compared to the tanh unit.

**ReLU.** The Rectified Linear Unit has become very popular in the last few years. It computes the function  $f(x) = \max(0, x)$ . In other words, the activation is simply thresholded at zero (see image above on the left). There are several pros and cons to using the ReLUs:

- (+) It was found to greatly accelerate (e.g. a factor of 6 in [Krizhevsky et al.](#)) the convergence of stochastic gradient descent compared to the sigmoid/tanh functions. It is argued that this is due to its linear, non-saturating form.
- (+) Compared to tanh/sigmoid neurons that involve expensive operations (exponentials, etc.), the ReLU can be implemented by simply thresholding a matrix of activations at zero.
- (-) Unfortunately, ReLU units can be fragile during training and can “die”. For example, a large gradient flowing through a ReLU neuron could cause the weights to update in such a way that the neuron will never activate on any datapoint again. If this happens, then the gradient flowing through the unit will forever be zero from that point on. That is, the ReLU units can irreversibly die during training since they can get knocked off the data manifold. For example, you may find that as much as 40% of your network can be “dead” (i.e. neurons that never activate across the entire training dataset) if the learning rate is set too high. With a proper setting of the learning rate this is less frequently an issue.

**Leaky ReLU.** Leaky ReLUs are one attempt to fix the “dying ReLU” problem. Instead of the function being zero when  $x < 0$ , a leaky ReLU will instead have a small negative slope (of 0.01, or so). That is, the function computes  $f(x) = 1(x < 0)(\alpha x) + 1(x \geq 0)(x)$  where  $\alpha$  is a small constant. Some people report success with this form of activation function, but the results are not always consistent. The slope in the negative region can also be made into a parameter of each neuron, as seen in PReLU neurons, introduced in [Delving Deep into Rectifiers](#), by Kaiming He et al., 2015. However, the consistency of the benefit across tasks is presently unclear.

**Maxout.** Other types of units have been proposed that do not have the functional form  $f(w^T x + b)$  where a non-linearity is applied on the dot product between the weights and the data. One relatively popular choice is the Maxout neuron (introduced recently by [Goodfellow et al.](#)) that generalizes the ReLU and its leaky version. The Maxout neuron computes the function  $\max(w_1^T x + b_1, w_2^T x + b_2)$ . Notice that both ReLU and Leaky ReLU are a special case of this form (for example, for ReLU we have  $w_1, b_1 = 0$ ). The Maxout neuron therefore enjoys all the benefits of a ReLU unit (linear regime of operation, no saturation) and does not have its drawbacks (dying ReLU). However, unlike the ReLU neurons it doubles the number of parameters for every single neuron, leading to a high total number of parameters.

This concludes our discussion of the most common types of neurons and their activation functions. As a last comment, it is very rare to mix and match different types of neurons in the same network, even though there is no fundamental problem with doing so.

**TLDR:** “*What neuron type should I use?*” Use the ReLU non-linearity, be careful with your learning rates and possibly monitor the fraction of “dead” units in a network. If this concerns you, give Leaky ReLU or Maxout a try. Never use sigmoid. Try tanh, but expect it to work worse than ReLU/Maxout.

## Neural Network architectures

### Layer-wise organization

**Neural Networks as neurons in graphs.** Neural Networks are modeled as collections of neurons that are connected in an acyclic graph. In other words, the outputs of some neurons can become inputs to other neurons. Cycles are not allowed since that would imply an infinite loop in the forward pass of a network. Instead of an amorphous blobs of connected neurons, Neural Network models are often organized into distinct layers of neurons. For regular neural networks, the most common layer type is the **fully-connected layer** in which neurons between two adjacent layers are fully pairwise connected, but neurons within a single layer share no connections. Below are two example Neural Network topologies that use a stack of fully-connected layers: