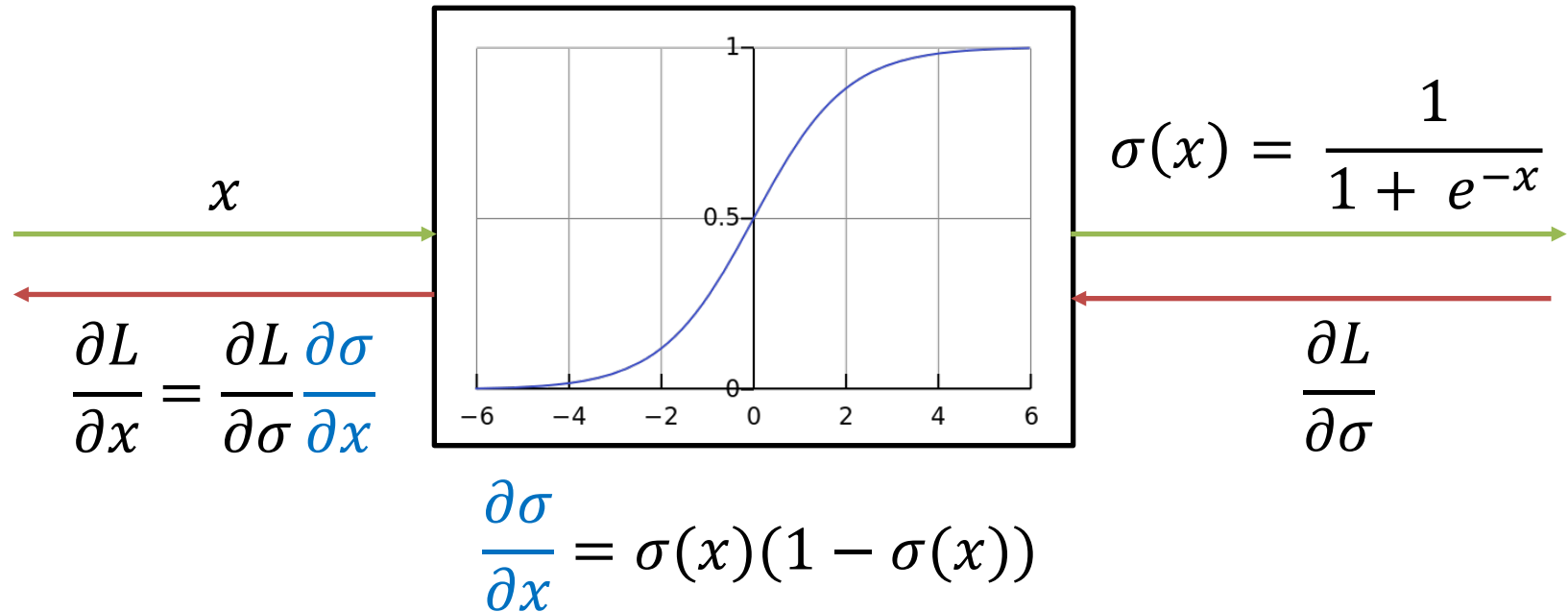
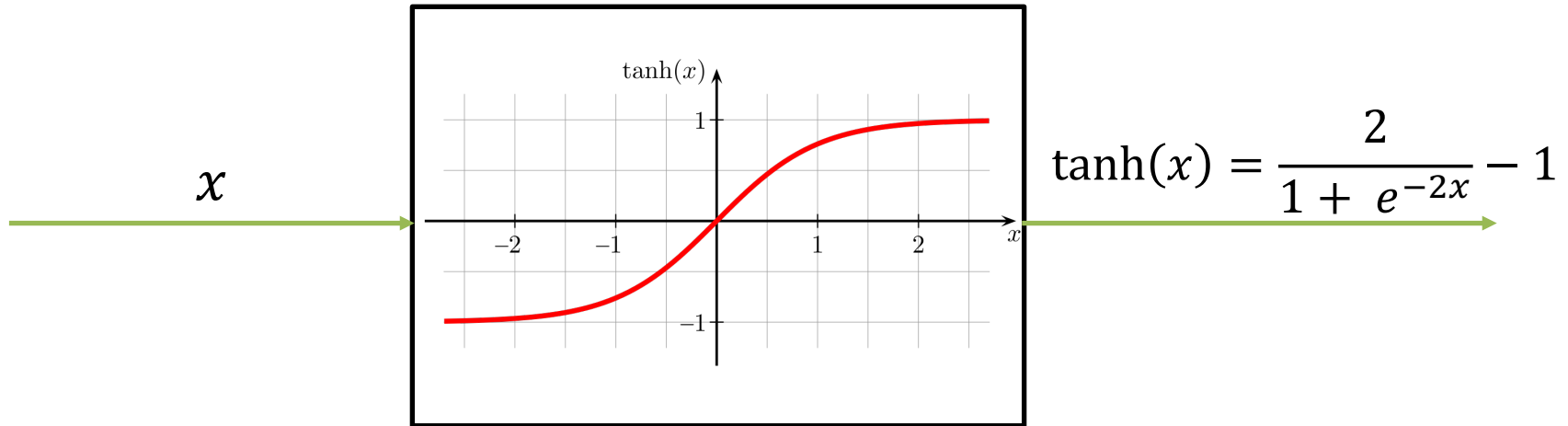


Sigmoid activation



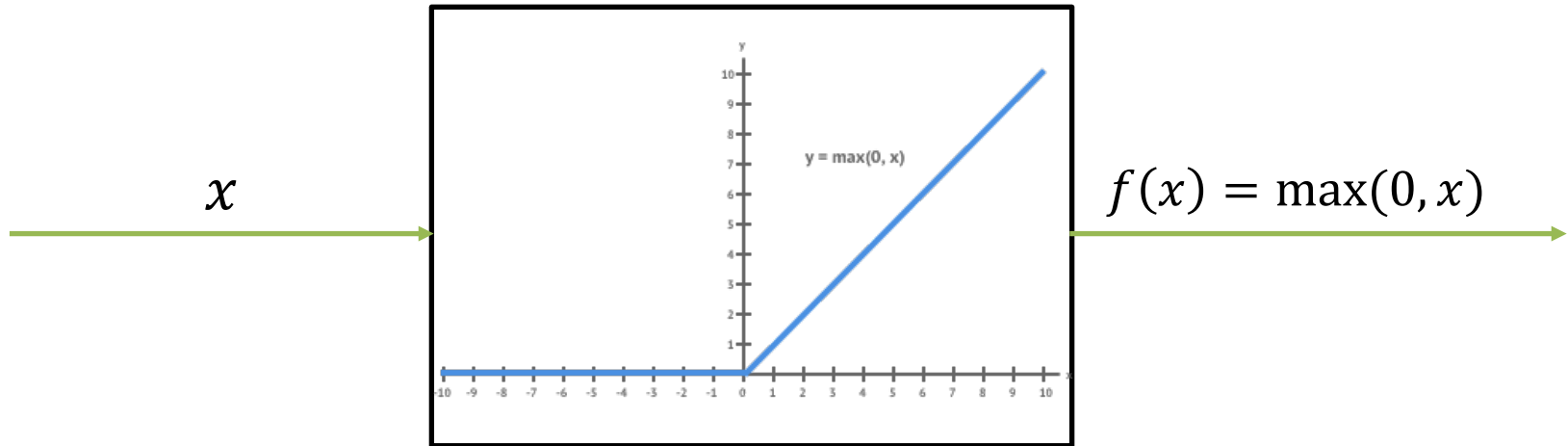
- Sigmoid neurons can saturate and lead to **vanishing gradients**.
- Not zero-centered. Why do we want it to be zero centered?
Neural networks like it when we have 0 mean and standard variance.
- e^x is computationally expensive.

Tanh activation



- Zero-centered.
- But still pretty much like sigmoid. Vanishing gradient problem.

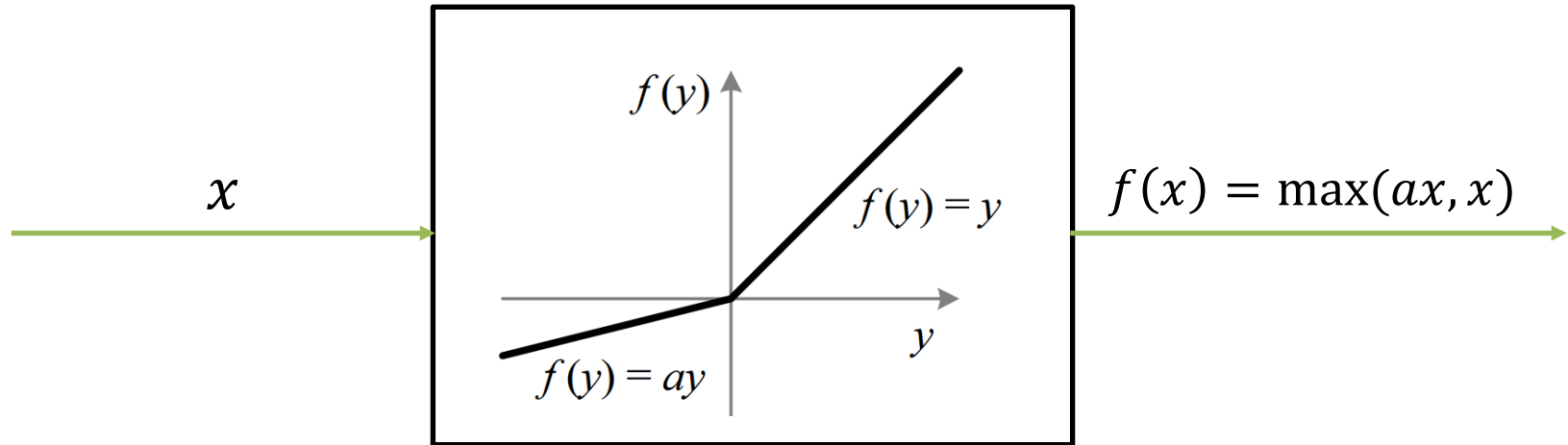
ReLU activation



- Fast to compute.
- Gradients do not vanish for $x > 0$.
- Provides faster convergence in practice!
- Not zero-centered.
- Can die: if not activated, never updates!

If x always less than 0,
then always 0 gradient.
This is the problem of
the dying neuron.

Leaky ReLU activation



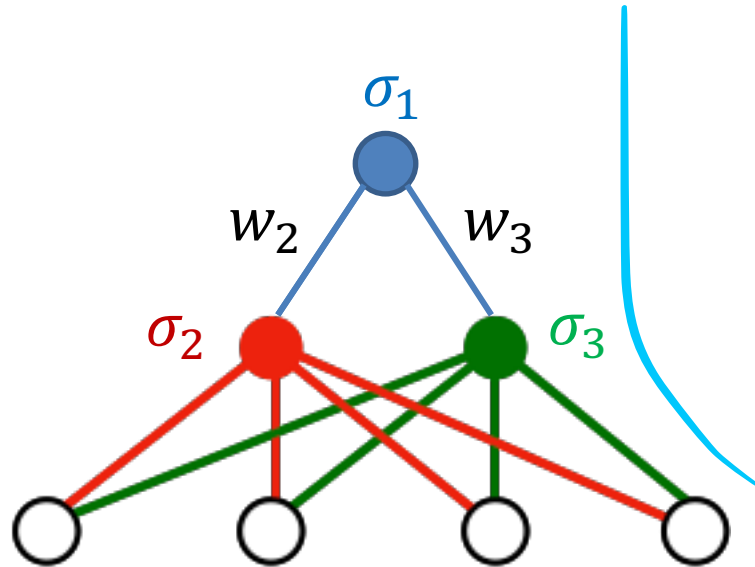
- Will not die!
- $a \neq 1$

This time, a leaky relu
will make $x < 0$ gradients
not 0.

Weights initializations

Maybe start with all zeros?

If we do this, then:



Inputs:

$$\frac{\partial L}{\partial w_2} = \frac{\partial L}{\partial \sigma_1} \sigma_1 (1 - \sigma_1) \sigma_2$$

$$\frac{\partial L}{\partial w_3} = \frac{\partial L}{\partial \sigma_1} \sigma_1 (1 - \sigma_1) \sigma_3$$

σ_2 and σ_3 will always get the same updates!

we will have the same neurons and not learn complex representations!

- Need to break symmetry!
- Maybe start with small random numbers then?
- But how small? $0.03 \cdot \mathcal{N}(0,1)$?

Weights initializations

- Linear models work best when inputs are normalized.
- Neuron is a linear combination of inputs + activation.
- Neuron output will be used by consecutive layers.

It would be great if we could normalize the output of the neurons.

Weights initializations

Let's look at the neuron output **before activation**: $\sum_{i=1}^n x_i w_i$.

E being expectation.

If $E(x_i) = E(w_i) = 0$ and we generate weights independently from inputs, then $E(\sum_{i=1}^n x_i w_i) = 0$.

But variance **can grow** with consecutive layers.

Empirically this hurts convergence for deep networks!

When gradients of different outputs become of different scale (like 1 to 100) gradient descent methods slow down drastically, this is called ill-conditioning!

With variance growing, we may overfit the model, and have poor testing accuracy.

Weights initializations

NOTE: for two random variables X and Y, $\text{Var}[XY] = (E[Y])^2 \text{Var}[X] + (E[X])^2 \text{Var}[Y] + \text{Var}[X] \cdot \text{Var}[Y]$

Let's look at the variance of $\sum_{i=1}^n x_i w_i$:

$$\text{Var}(\sum_{i=1}^n x_i w_i) = \text{i.i.d. } w_i \text{ and mostly uncorrelated } x_i$$

$$= \sum_{i=1}^n \text{Var}(x_i w_i) = \text{independent factors } w_i \text{ and } x_i$$

$$= \sum_{i=1}^n \left(\begin{array}{l} [E(x_i)]^2 \text{Var}(w_i) \\ + [E(w_i)]^2 \text{Var}(x_i) \\ + \text{Var}(x_i) \text{Var}(w_i) \end{array} \right) = \text{ } w_i \text{ and } x_i \text{ have 0 mean}$$

thus,...

$$= \sum_{i=1}^n \text{Var}(x_i) \text{Var}(w_i) = \text{Var}(x) [n \text{Var}(w)]$$

↑
We want this to be 1

Why 1? because if greater than 1, then variance of consecutive layers will grow.

Weights initializations

Our goal is for $n \text{Var}(w)$ to be 1.

- Let's use the fact that $\text{Var}(aw) = a^2 \text{Var}(w)$.
 - For $[n \text{Var}(aw)]$ to be 1 we need to multiply $\mathcal{N}(0,1)$ weights ($\text{Var}(w) = 1$) by $a = 1/\sqrt{n}$.
 - the weights are drawn from standard normal distribution, mean of 0 and stdev 1
 - Working out:
Want $n\text{Var}(aw) = 1$
 $a^2 n \text{Var}(w) = 1$
 $a = 1 / \text{sqrt}(n)$
 - Xavier initialization (Glorot et al.) multiplies weights by $\sqrt{2}/\sqrt{n_{in} + n_{out}}$.
 - Initialization for ReLU neurons (He et al.) uses multiplication by $\sqrt{2}/\sqrt{n_{in}}$.
- THE IDEA is to multiply initial weights with some constant. This limits the growth of variance initially.

Batch normalization

- We know how to initialize our network to constrain variance.
- But what if it grows during backpropagation?
- Batch normalization controls mean and variance of outputs **before activations**.

Weights
initialization may
not help!

Batch normalization

- Let's normalize h_i – neuron output before activation:

We normalize the neuron
so that the variance won't
grow.

$$h_i = \gamma_i \frac{h_i - \mu_i}{\sqrt{\sigma_i^2}} + \beta_i$$

→ 0 mean, unit variance

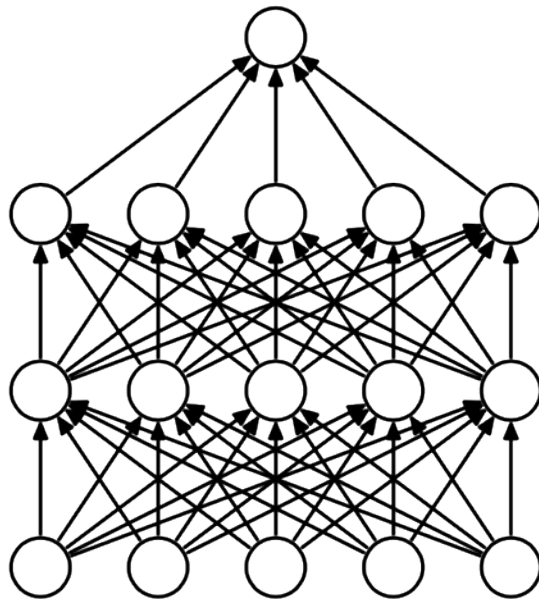
- Where do μ_i and σ_i^2 come from? We can estimate them having a **current training batch!**
- During testing we will use an **exponential moving average** over train batches:

$$0 < \alpha < 1$$
$$\mu_i = \alpha \cdot \mathbf{mean}_{\text{batch}} + (1 - \alpha) \cdot \mu_i$$
$$\sigma_i^2 = \alpha \cdot \mathbf{variance}_{\text{batch}} + (1 - \alpha) \cdot \sigma_i^2$$

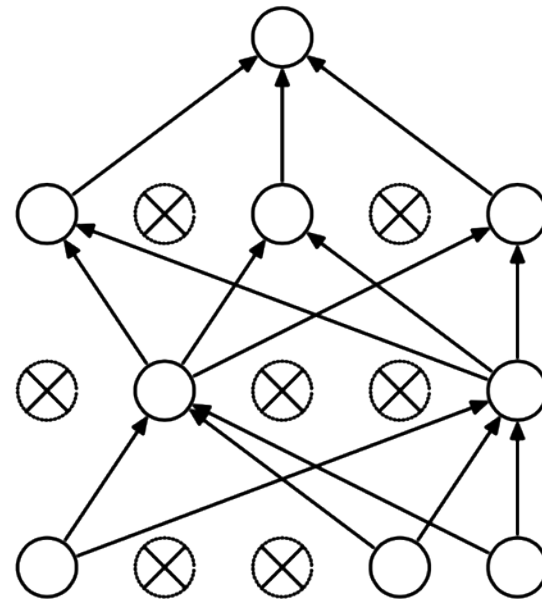
- What about γ_i and β_i ? Normalization is a differentiable operation and we can apply backpropagation!

Dropout

- Regularization technique to reduce overfitting.
- We keep neurons active (non-zero) with probability p .
- This way we sample the network during training and change only a subset of its parameters on every iteration.



(a) Standard Neural Net

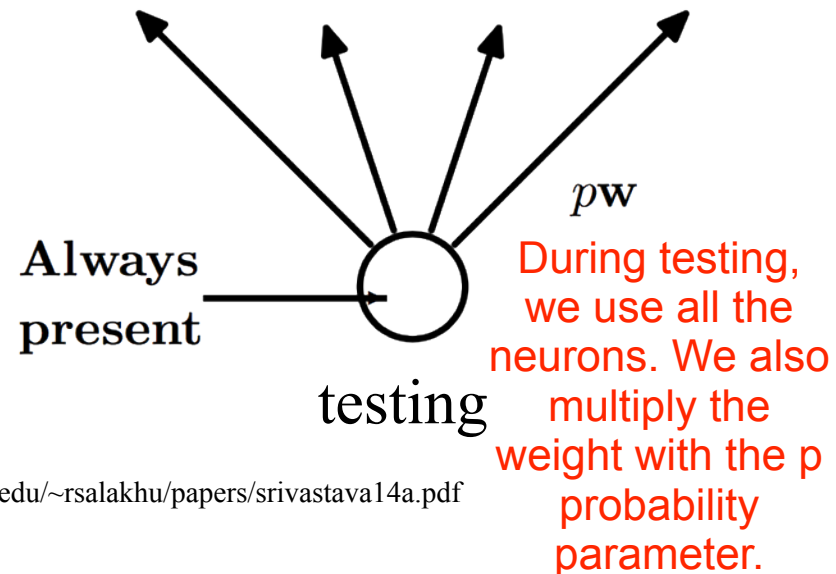
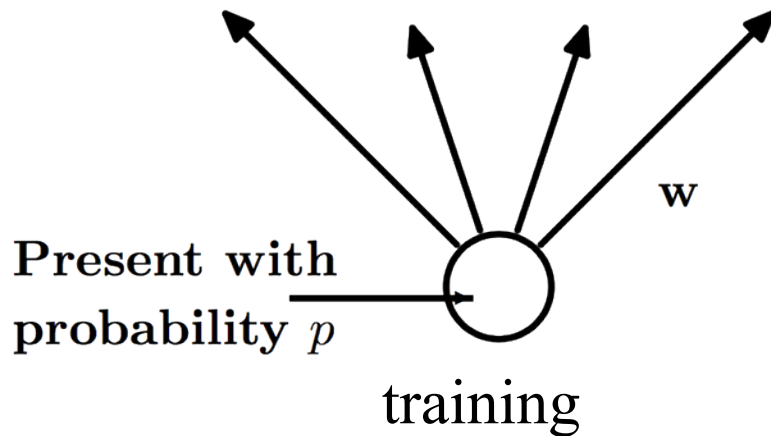


(b) After applying dropout.

Dropout

- During testing all neurons are present but their outputs are multiplied by p to maintain the scale of inputs:

Expected input weight: $p \cdot w + (1 - p) \cdot 0$



Nitish Srivastava, <http://www.cs.toronto.edu/~rsalakhu/papers/srivastava14a.pdf>

- The authors of dropout say it's similar to having an ensemble of exponentially large number of smaller networks.

Data augmentation

- Modern CNN's have millions of parameters!
- But datasets are not that huge!
- We can generate new examples applying distortions: flips, rotations, color shifts, scaling, etc.
- Remember: CNN's are invariant to translation

So we don't need to move images around to augment input images!



Takeaways

- Use ReLU activation
- Use He et al. initialization
- Try to add BN or dropout
- Try to augment your training data
- In the next video you will learn how modern convolutional networks look like