

Lecture 6

Training Neural Networks (I)



Where we are now...

Mini-batch SGD

Loop:

1. **Sample** a batch of data
2. **Forward** prop it through the graph
(network), get loss
3. **Backprop** to calculate the gradients
4. **Update** the parameters using the gradient

Next: Training Neural Networks

1. One time setup

Activation functions, data preprocessing, weight initialization, regularization

Today

2. Training dynamics

Learning rate schedules; large-batch training;
hyperparameter optimization

Next time

3. After training

Model ensembles, transfer learning

Activation Functions

Where are Activation Functions?

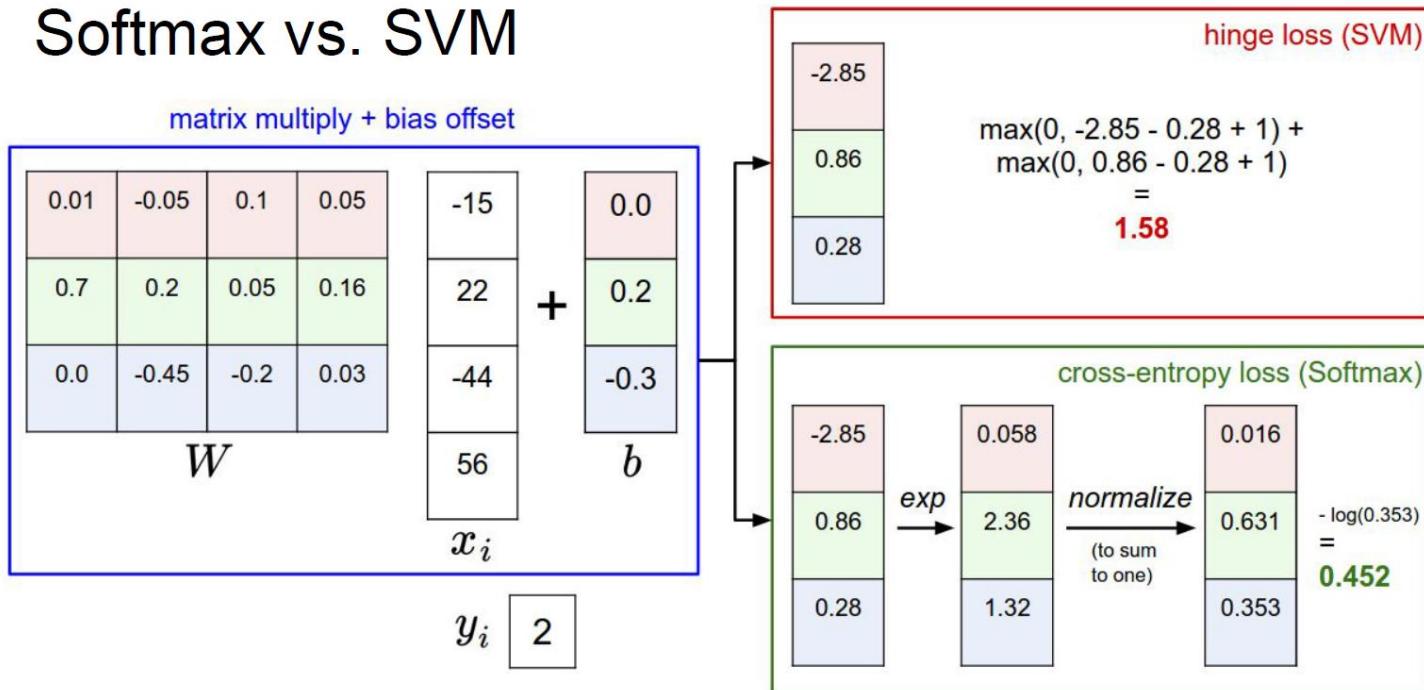
(Before) Linear score function: $f = Wx$

(Now) 2-layer Neural Network $f = W_2 \max(0, W_1 x)$
or 3-layer Neural Network

$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$

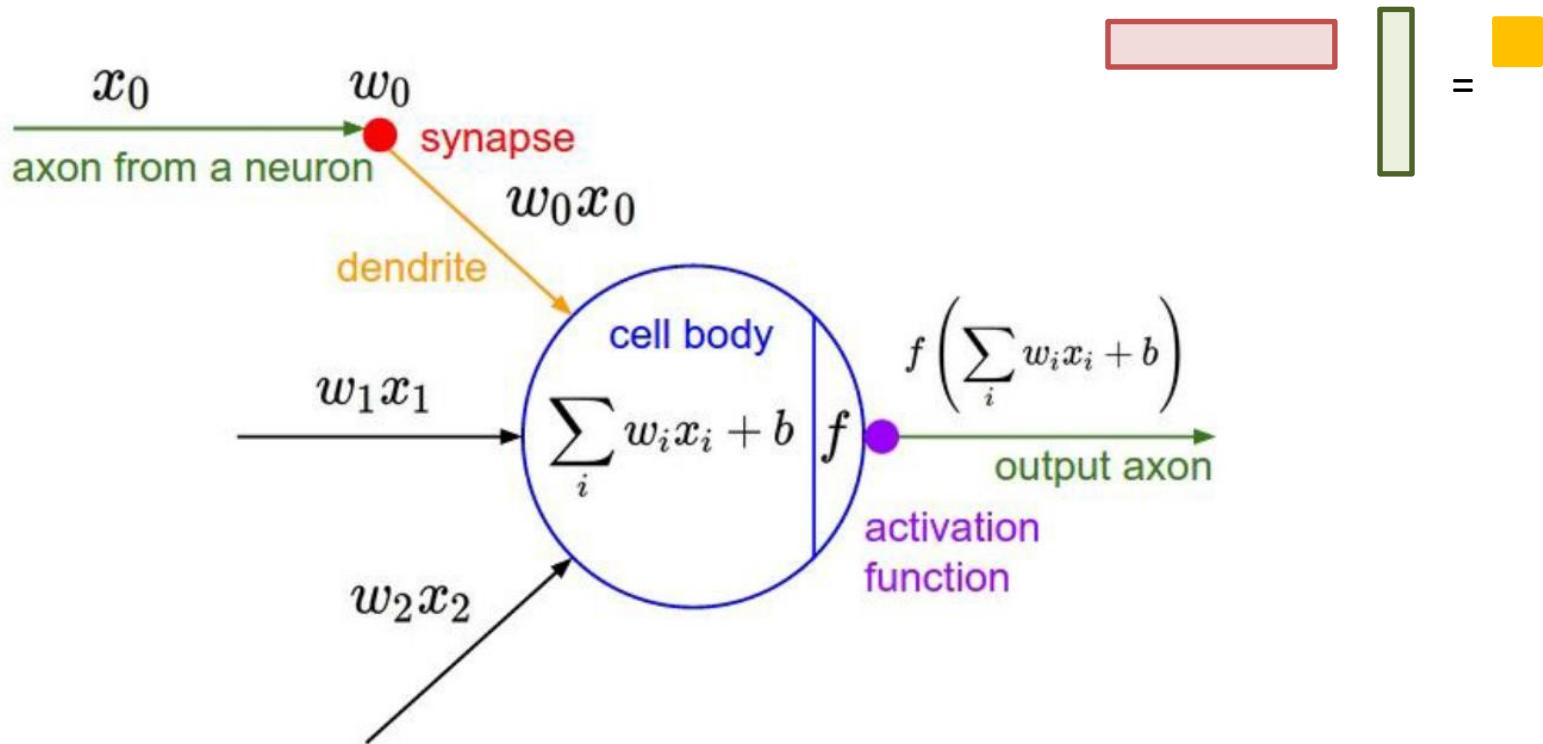
Don't confuse with **Loss Function**

Softmax vs. SVM



Activation Functions

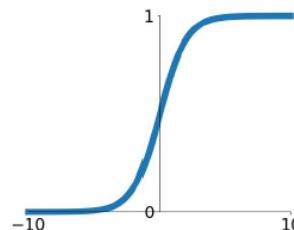
$$f = W_3 \max(0, W_2 \max(0, W_1 x))$$



Activation Functions

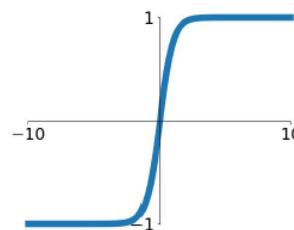
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



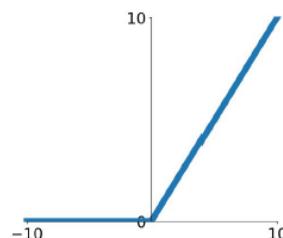
tanh

$$\tanh(x)$$



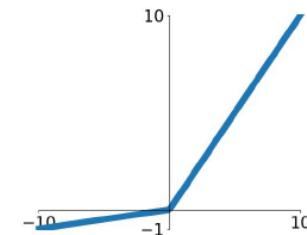
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(0.1x, x)$$

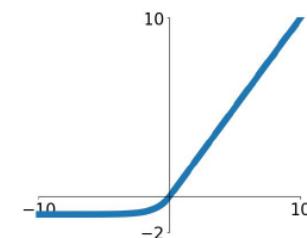


Maxout

$$\max(w_1^T x + b_1, w_2^T x + b_2)$$

ELU

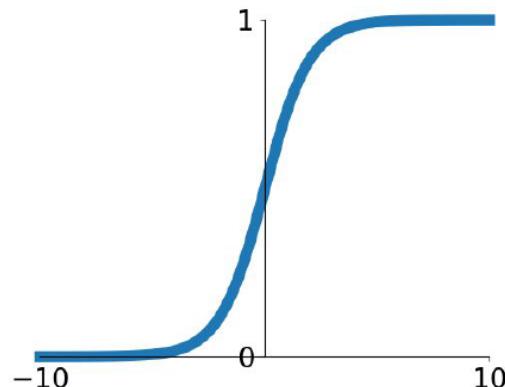
$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Activation Functions

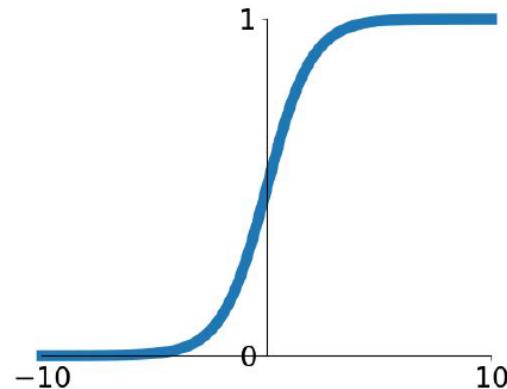
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron



Sigmoid

Activation Functions



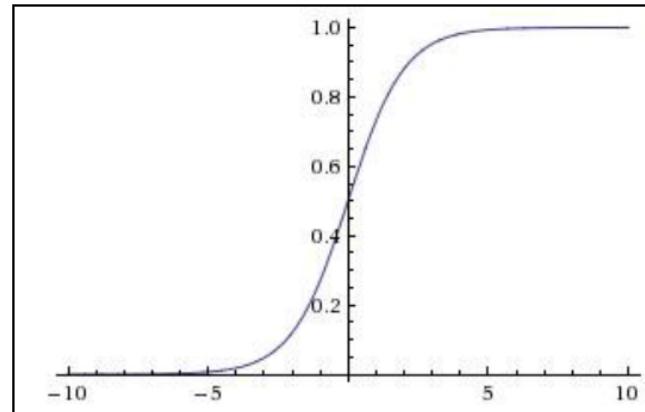
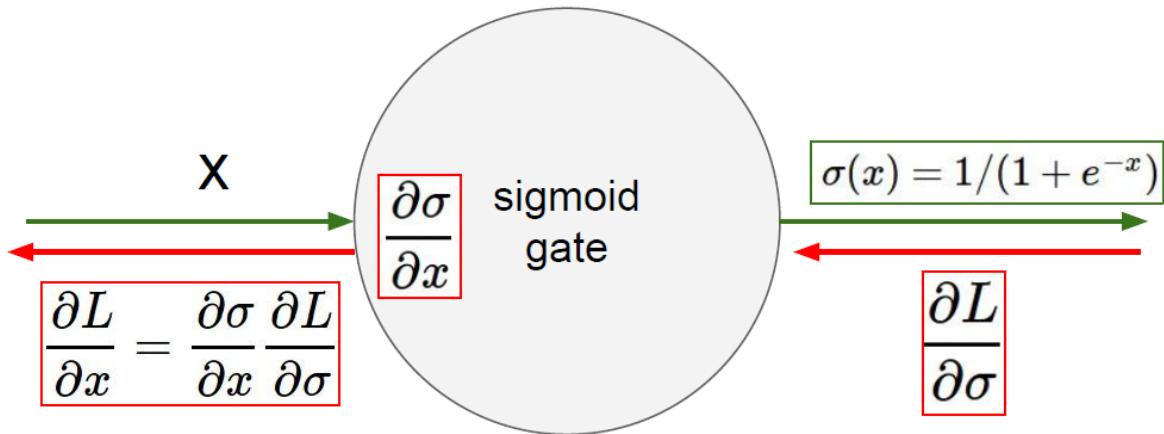
Sigmoid

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

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients

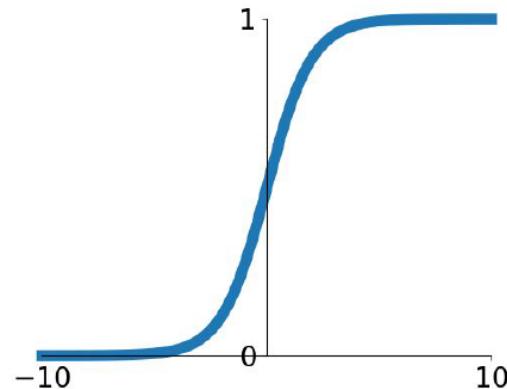


What happens when $x = -10$?

What happens when $x = 0$?

What happens when $x = 10$?

Activation Functions



Sigmoid

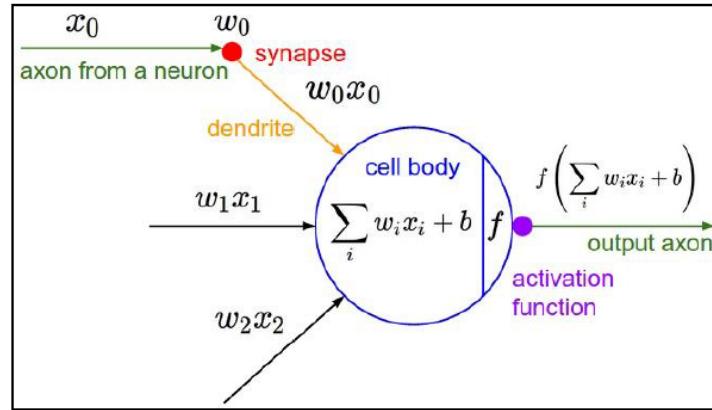
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered

Consider what happens when the input to a neuron (x) is always positive:



$$f \left(\sum_i w_i x_i + b \right) = Z$$

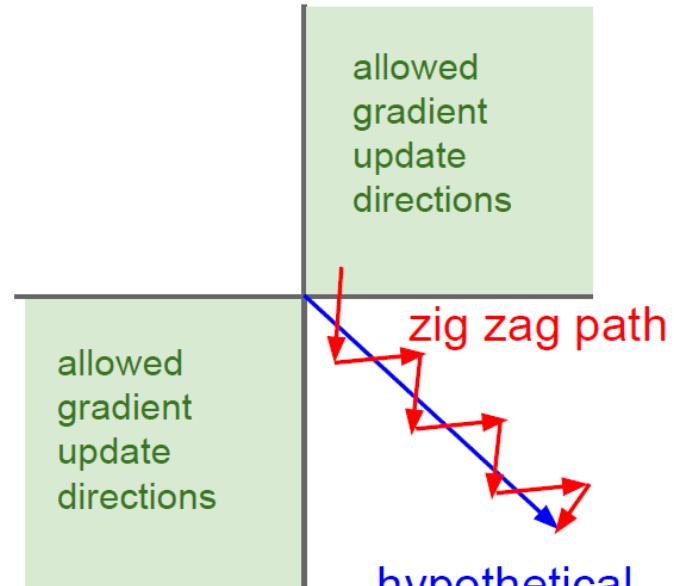
What can we say about the gradients on w ?

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_1} = x_1$$

$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_2} = x_2$$

Consider what happens when the input to a neuron is always positive...

$$\frac{\partial f}{\partial w_1} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_1} = x_1$$
$$\frac{\partial f}{\partial w_2} = \frac{\partial f}{\partial z} \cdot \frac{\partial z}{\partial w_2} = x_2$$

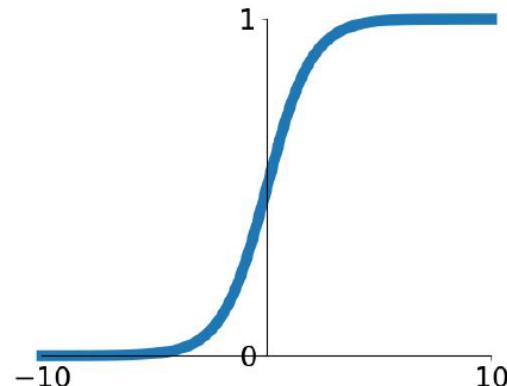


What can we say about the gradients on w ?

Always all positive or all negative :(

(this is also why you want zero-mean data!)

Activation Functions



Sigmoid

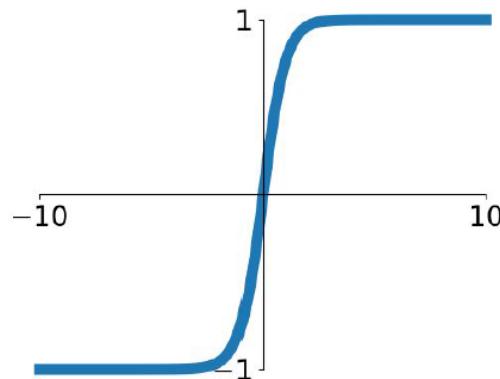
$$\sigma(x) = 1/(1 + e^{-x})$$

- Squashes numbers to range [0,1]
- Historically popular since they have nice interpretation as a saturating “firing rate” of a neuron

3 problems:

1. Saturated neurons “kill” the gradients
2. Sigmoid outputs are not zero-centered
3. $\exp()$ is a bit compute expensive

Activation Functions

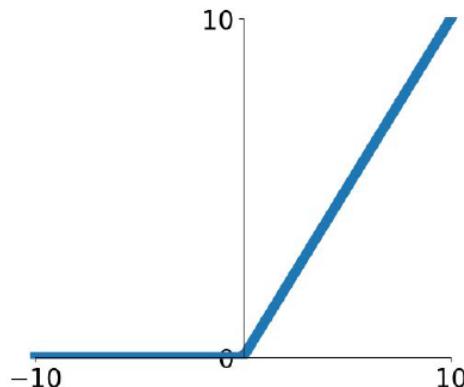


$\tanh(x)$

- Squashes numbers to range [-1,1]
- zero centered (nice)
- still kills gradients when saturated :(

[LeCun et al., 1991]

Activation Functions

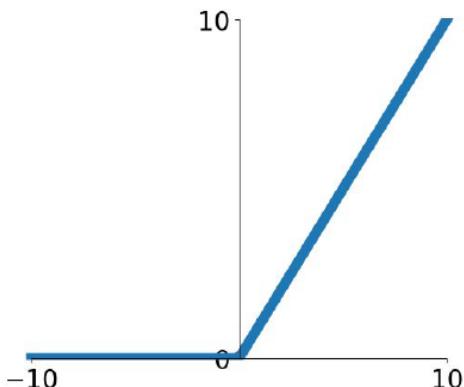


- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

ReLU
(Rectified Linear Unit)

[Krizhevsky et al., 2012]

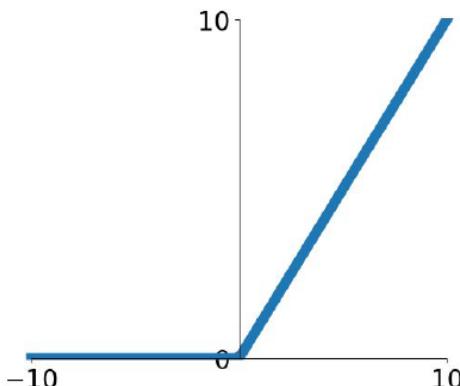
Activation Functions



ReLU
(Rectified Linear Unit)

- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid
- Not zero-centered output

Activation Functions

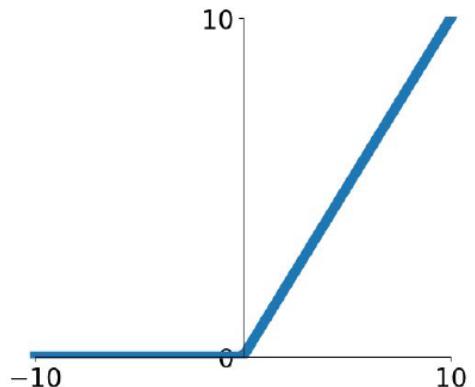
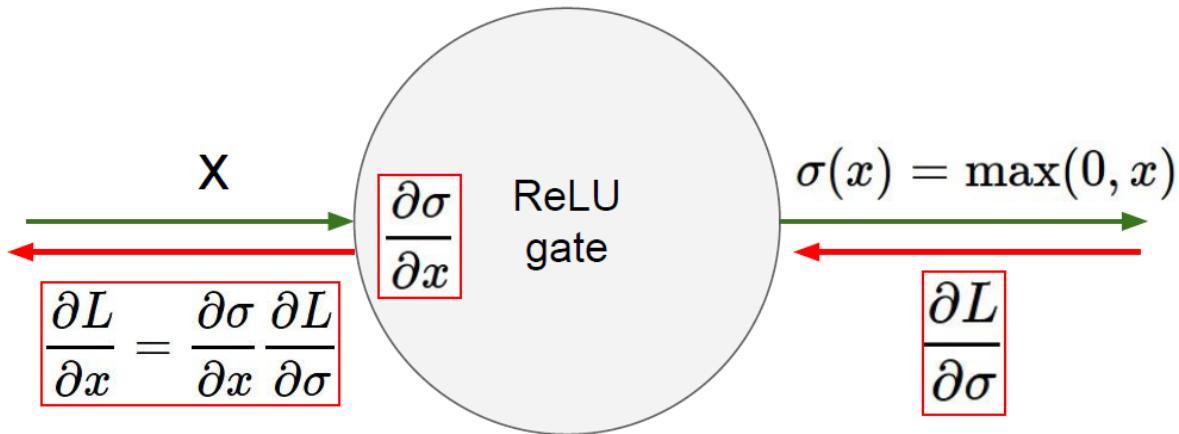


ReLU
(Rectified Linear Unit)

- Computes $f(x) = \max(0, x)$
- Does not saturate (in +region)
- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)
- Actually more biologically plausible than sigmoid

- Not zero-centered output
- An annoyance:

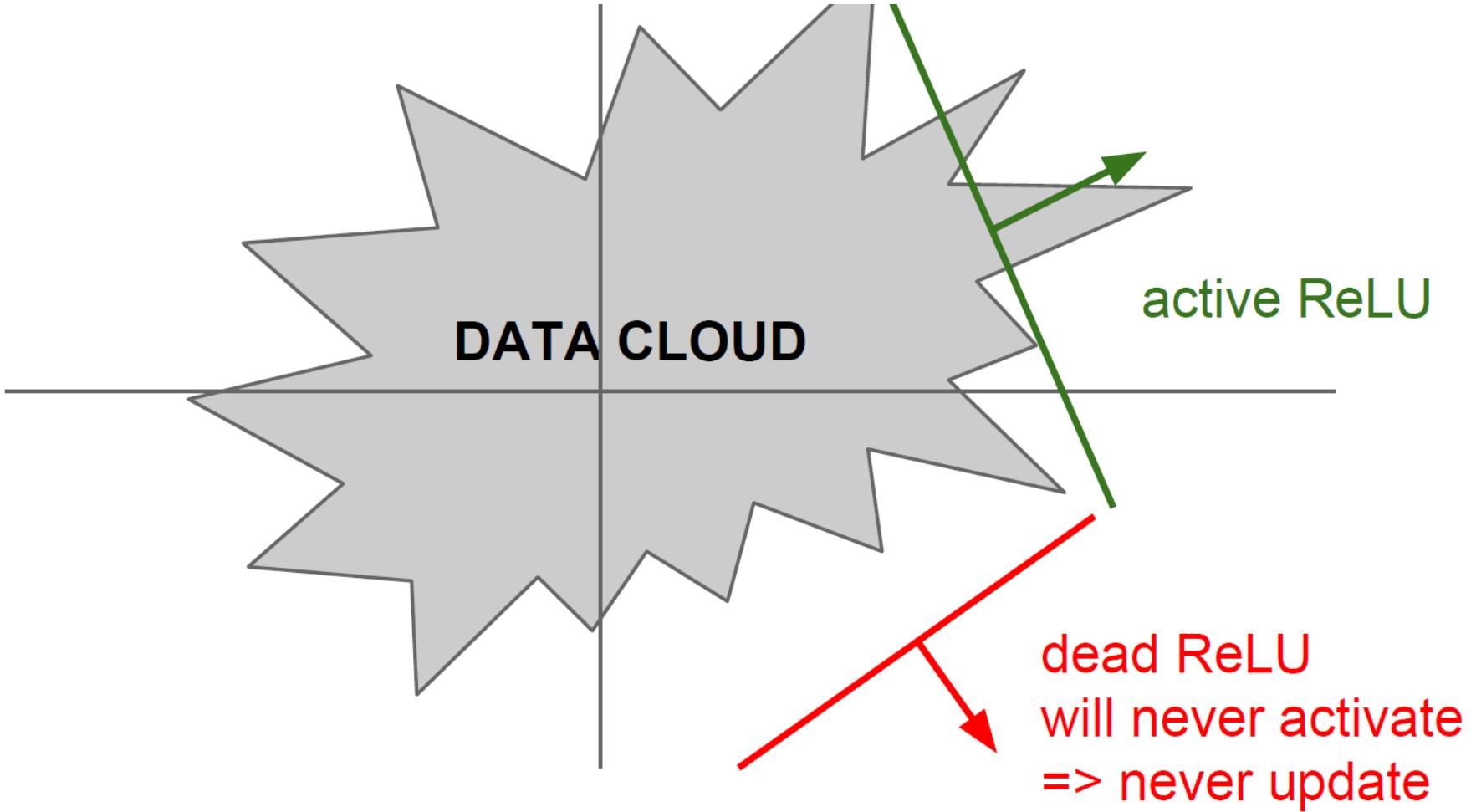
hint: what is the gradient when $x < 0$?

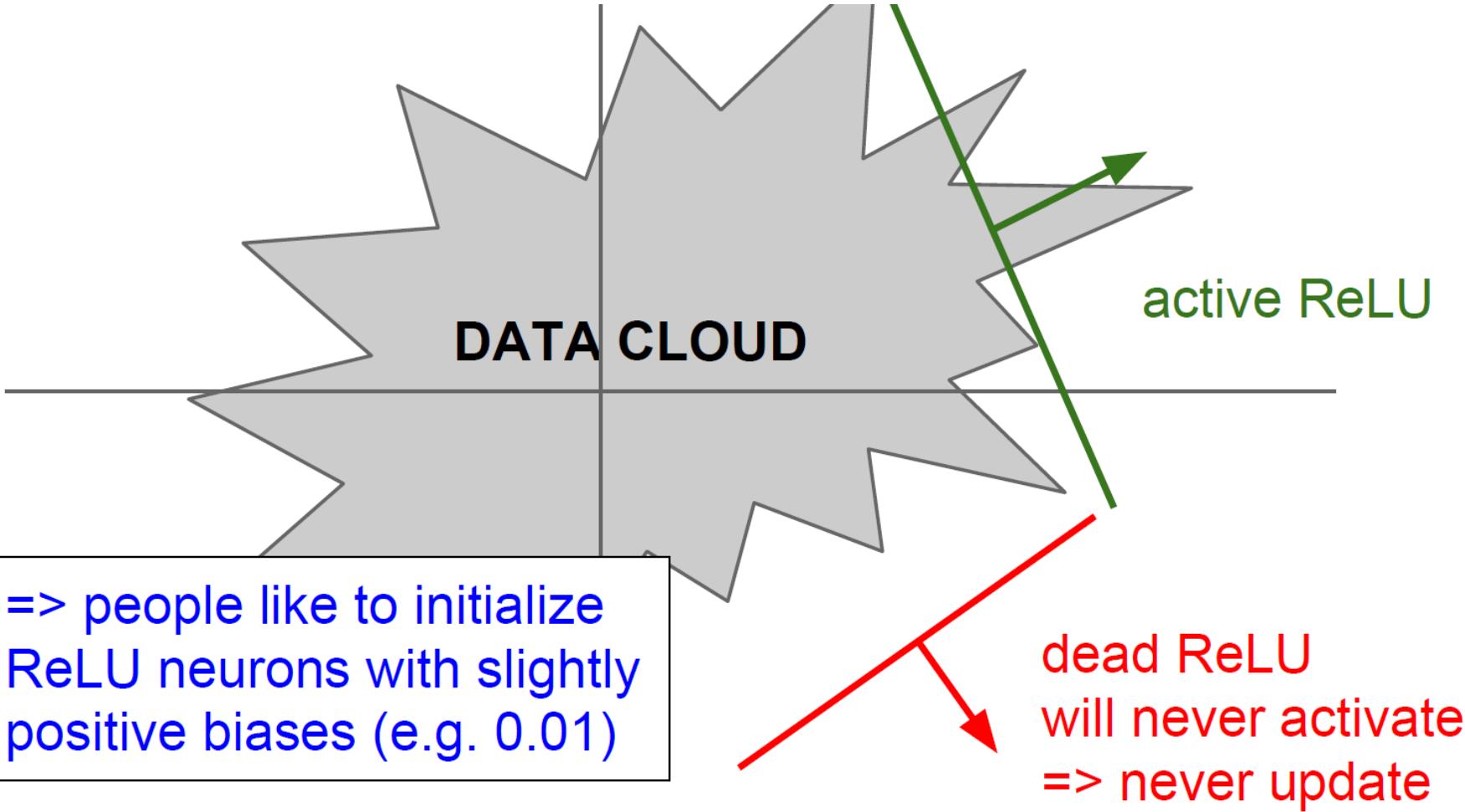


What happens when $x = -10$?

What happens when $x = 0$?

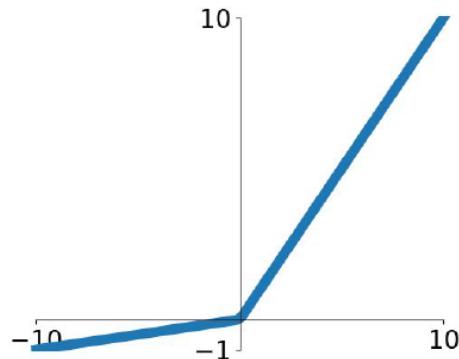
What happens when $x = 10$?





Activation Functions

[Mass et al., 2013]
[He et al., 2015]

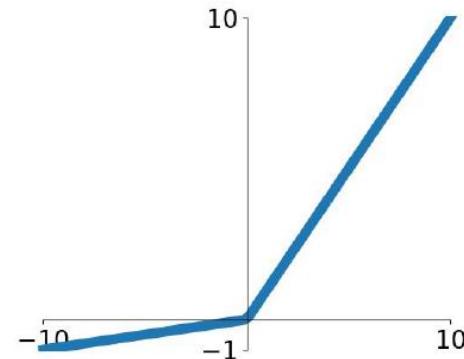


- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Leaky ReLU

$$f(x) = \max(0.01x, x)$$

Activation Functions



Leaky ReLU

$$f(x) = \max(0.01x, x)$$

[Mass et al., 2013]
[He et al., 2015]

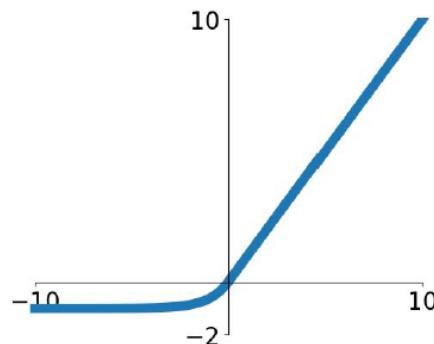
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not “die”.**

Parametric Rectifier (PReLU)

$$f(x) = \max(\alpha x, x)$$

backprop into α
(parameter)

Exponential Linear Units (ELU)

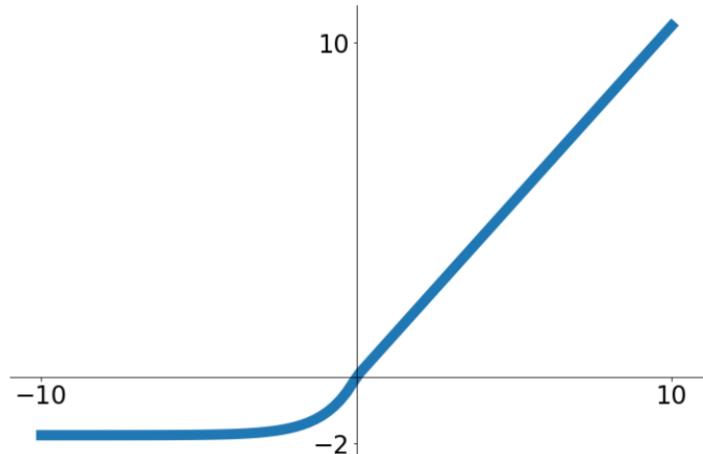


- All benefits of ReLU
- Closer to zero mean outputs
- Negative saturation regime compared with Leaky ReLU adds some robustness to noise

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ \alpha (\exp(x) - 1) & \text{if } x \leq 0 \end{cases}$$

- Computation requires $\exp()$

Activation Functions: Scaled Exponential Linear Unit (SELU)



- Scaled version of ELU that works better for deep networks
- “Self-Normalizing” property; can train deep SELU networks without BatchNorm

$$selu(x) = \begin{cases} \lambda x & \text{if } x > 0 \\ \lambda\alpha(e^x - 1) & \text{if } x \leq 0 \end{cases}$$

$$\alpha = 1.6732632423543772848170429916717$$

$$\lambda = 1.0507009873554804934193349852946$$

Activation Functions: Scaled Exponential Linear Unit (SELU)

$0 \leq \mu < 1$ and $0 \leq \nu \leq 0.1$.
 g is increasing in μ and increasing in ν . We set $\mu = 1$ and $\nu = 0.0180173$.

Therefore the maximal value of g is -0.0180173 .

A.3 Proof of Theorem 2

First we recall Theorem 2. **Theorem 2.** *We consider $\lambda = \lambda_{01} = 0.01$ on and the two domains Ω_1^+ ($\{\mu, \omega, \nu\} | -0.1 \leq \mu \leq 0.1, -0.1 \leq \nu \leq 0.1, 0.05 \leq \omega \leq 0.16, 0.8 \leq \tau \leq 1.25\}$) and $\Omega_2^- = \{\mu, \nu, \tau | -0.1 \leq \mu \leq 0.1, -0.1 \leq \nu \leq 0.1, 0.05 \leq \nu \leq 0.24, 0.9 \leq \tau \leq 1.25\}$. The mapping of the variance $\tilde{\xi}(\mu, \omega, \nu, \tau, \lambda, \alpha_1, \alpha_2)$ given in Eq. (5) increases*

$$\tilde{\xi}(\mu, \omega, \nu, \tau, \lambda, \alpha_1, \alpha_2) > \nu \quad (44)$$

in both Ω_1^+ and Ω_2^- . All fixed points (μ, ν) of mapping Eq. (5) and Eq. (4) ensure for $0.8 \leq \tau$ that $\nu = 0.05$ and for $0.9 \leq \tau$ that $\nu = \nu_{max}$. Consequently, the variance mapping Eq. (5) and Eq. (4) ensures a lower bound on the variance ν .

Proof. The mean value theorem states that there exists $\epsilon \in [0, 1]$ for which

$$\begin{aligned} \tilde{\xi}(\mu, \omega, \nu, \tau, \lambda, \alpha_1, \alpha_2) - \tilde{\xi}(\mu, \omega, \nu_{max}, \tau, \lambda, \alpha_1, \alpha_2) &= \\ \frac{\partial}{\partial \nu} \tilde{\xi}(\mu, \omega, \nu + (\nu_{max} - \nu), \tau, \lambda, \alpha_1, \alpha_2) (\nu - \nu_{max}). \end{aligned} \quad (45)$$

Therefore

$$\begin{aligned} \tilde{\xi}(\mu, \omega, \nu, \tau, \lambda, \alpha_1, \alpha_2) &= \tilde{\xi}(\mu, \omega, \nu_{max}, \tau, \lambda, \alpha_1, \alpha_2) + \\ \frac{\partial}{\partial \nu} \tilde{\xi}(\mu, \omega, \nu + (\nu_{max} - \nu), \tau, \lambda, \alpha_1, \alpha_2) (\nu - \nu_{max}). \end{aligned} \quad (46)$$

Therefore we are interested to bound the derivative of the ξ -mapping Eq. (13) with respect to ν :

$$\begin{aligned} \frac{\partial}{\partial \nu} \tilde{\xi}(\mu, \omega, \nu, \tau, \lambda, \alpha_1, \alpha_2) &= \\ \frac{3}{2} \lambda^2 \nu^{-\frac{1}{2}} e^{\frac{\mu \nu}{2}} \left(\sigma^2 \left(-\left(e^{\frac{\mu \nu}{2} + 0.05}{\sqrt{2\nu}} \right)^2 \operatorname{erfc}\left(\frac{\mu \nu + \nu \tau}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) \right) - \right. \\ \left. 2e^{\left(\frac{\mu \nu}{2} + 0.05 \right)^2} \operatorname{erfc}\left(\frac{\mu \nu + 2\nu \tau}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) \right) - \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) + 2. \end{aligned} \quad (47)$$

The sub-term Eq. (309) enters the derivative Eq. (47) with a negative sign! According to Lemma 18, the minimal value of sub-term Eq. (308) is obtained by the largest largest ν , by the smallest τ , and the largest $\nu = \mu \omega = 0.05$. Also the positive term $\operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu} \sqrt{\nu \tau}} \right)$ is multiplied by τ , which is minimized by using the smallest τ . Therefore we can use the smaller τ in formula Eq. (47) to obtain the lower bound:

$$\frac{1}{2} \lambda^2 \nu^{-\frac{1}{2}} e^{\frac{\mu \nu}{2}} \left(\sigma^2 \left(-\left(e^{\frac{\mu \nu}{2} + 0.05}{\sqrt{2\nu}} \right)^2 \operatorname{erfc}\left(\frac{\mu \nu + \nu \tau}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) \right) - \right. \\ \left. 2e^{\left(\frac{\mu \nu}{2} + 0.05 \right)^2} \operatorname{erfc}\left(\frac{\mu \nu + 2\nu \tau}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) \right) - \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) + 2. \quad (48)$$

□

$$\begin{aligned} 0.8 - 0.01 &- \\ 2\sqrt{0.16 - 0.8} &\times \left(\frac{0.01}{\sqrt{2\sqrt{0.05 - 0.8}}} + 2 \right) > 0.960231. \end{aligned}$$

Now we follow the proof of Lemma 8, al and $x = \nu \tau$ must be minimal. Thus, the $\nu_{max}, \alpha_1, \alpha_2 = 0.0962727$ for $0.5 \leq \nu$ and

(Lemma 43) provide:

$$\begin{aligned} (\mu_1)^2 &> \\ 0.001281115 &+ 0.960231 \nu > \\ > \nu. \end{aligned} \quad (49)$$

$\nu \leq 0.125$. The factor consisting of the ν is $\frac{0.01}{2\sqrt{0.05 - 0.8}}$. Since erfc is monotonically increasing in ν in order to obtain the maximal

$$\begin{aligned} \text{the derivative:} \\ \nu^{\left(\frac{\mu \nu}{2} + 0.05 \right)^2} \operatorname{erfc}\left(\frac{\mu \nu + 2\nu \tau}{\sqrt{2\nu} \sqrt{\nu \tau}} \right) - \\ (50) \end{aligned}$$

$$\begin{aligned} 4 - 0.9 + 0.01 &- \\ 2\sqrt{0.24 - 0.9} &\times \left(\frac{0.01}{\sqrt{2\sqrt{0.05 - 0.9}}} + 2 \right) > 0.976932. \end{aligned} \quad (51)$$

Now we follow the proof of Lemma 8, al and $x = \nu \tau$ must be minimal. Thus, the $\nu_{max}, \alpha_1, \alpha_2 = 0.07378409$ for $0.05 \leq \nu$ and 1 on $[\mu^2]$. (Lemma 43) gives

$$\begin{aligned} (\mu_0)^2 &> \\ 0.01995928 &+ 0.976932 > \\ > \nu. \end{aligned} \quad (52)$$

□

softs

cobian norm smaller than one
The Jacobian of the mapping g is smaller than one in a larger domain than the original extend to $\tau \in [0.8, 1.25]$. The range of the following domain throughout this section: $[0.8, 1.25]$.

19

In the following, we denote two Jacobians: (1) the Jacobian J of the mapping $g: (\mu, \nu) \rightarrow (\tilde{\mu}, \tilde{\nu})$ and (2) the Jacobian H of the mapping $g: (\mu, \nu, \tau) \rightarrow (\tilde{\mu}, \tilde{\nu}, \tilde{\tau})$. The many properties of the system can already be seen on J .

$$\begin{aligned} J_{11} &= \frac{\partial \tilde{\mu}}{\partial \mu} \left(\frac{\partial \tilde{\mu}}{\partial \mu} \frac{\partial \tilde{\mu}}{\partial \nu} \right) \\ J_{12} &= \frac{\partial \tilde{\mu}}{\partial \nu} \left(\frac{\partial \tilde{\mu}}{\partial \mu} \frac{\partial \tilde{\mu}}{\partial \nu} \right) \end{aligned} \quad (52)$$

$$\begin{aligned} H_{11} &= \frac{\partial \tilde{\mu}}{\partial \mu} \left(\frac{\partial \tilde{\mu}}{\partial \mu} \frac{\partial \tilde{\mu}}{\partial \nu} \right) \\ H_{12} &= \frac{\partial \tilde{\mu}}{\partial \nu} \left(\frac{\partial \tilde{\mu}}{\partial \mu} \frac{\partial \tilde{\mu}}{\partial \nu} \right) \end{aligned} \quad (53)$$

of the Jacobian J is:

$$\begin{aligned} \frac{\partial}{\partial \mu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \mu \nu + \nu \tau &- \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + 2 \end{aligned} \quad (54)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \frac{\mu \nu + \nu \tau}{\sqrt{2\nu \tau}} &- (\alpha - 1) \sqrt{\frac{2}{\pi \nu \tau}} e^{-\frac{\mu^2 \nu^2}{2\nu \tau}} \end{aligned} \quad (55)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + & \\ \frac{2\nu \tau}{2\sqrt{\nu \tau}} &+ \mu \left(2 - \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) \right) + \sqrt{\frac{2}{\pi} \nu \tau - \frac{\mu^2 \nu^2}{2\nu \tau}} \end{aligned} \quad (56)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + & \\ \frac{\mu \nu + \nu \tau}{\sqrt{2\nu \tau}} &- \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + 2 \end{aligned} \quad (57)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + & \\ \frac{\mu \nu + \nu \tau}{\sqrt{2\nu \tau}} &- \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + 2 \end{aligned} \quad (58)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + & \\ \frac{\mu \nu + \nu \tau}{\sqrt{2\nu \tau}} &- \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + 2 \end{aligned} \quad (59)$$

$$\begin{aligned} \frac{\partial}{\partial \nu} \mu(\mu, \nu, \tau, \lambda, \alpha) &= \\ \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + & \\ \frac{\mu \nu + \nu \tau}{\sqrt{2\nu \tau}} &- \operatorname{erfc}\left(\frac{\mu \nu}{\sqrt{2\nu \tau}} \right) + 2 \end{aligned} \quad (60)$$

largest singular value of the Jacobian. If the largest singular value 1 , then the spectral norm of the Jacobian is smaller than 1 . Then the mean and variance to the mean and variance in the next layer is

The value is smaller than 1 by evaluating the function $S(\mu, \nu, \tau, \lambda, \alpha)$ can Value Theorem to bound the deviation of the function S between we have to bound the gradient of S with respect to (μ, ω, ν, τ) . If all terms in the difference between points and evaluated points) have proved that the function is below in

2 matrix

$$A = \begin{pmatrix} a_{11} & a_{12} \\ a_{21} & a_{22} \end{pmatrix} \quad (58)$$

$$\begin{aligned} a_{11}^2 + (a_{11} - a_{12})^2 &+ \sqrt{(a_{11} - a_{12})^2 + (a_{12} + a_{21})^2} \\ a_{22}^2 + (a_{21} - a_{12})^2 &- \sqrt{(a_{11} - a_{12})^2 + (a_{12} + a_{21})^2}. \end{aligned} \quad (59)$$

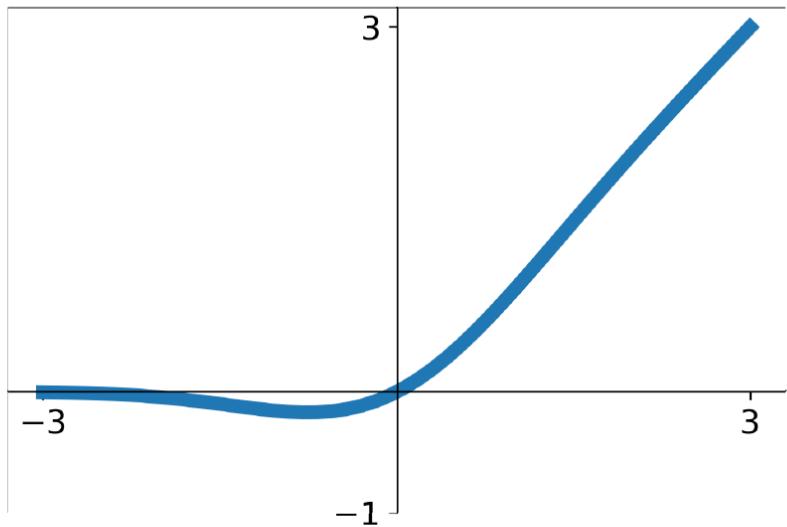
$$\begin{aligned} a_{11}^2 + (a_{11} - a_{12})^2 &- \sqrt{(a_{11} - a_{12})^2 + (a_{12} + a_{21})^2} \\ a_{22}^2 + (a_{21} - a_{12})^2 &- \sqrt{(a_{11} - a_{12})^2 + (a_{12} + a_{21})^2}. \end{aligned} \quad (60)$$

Scaled version of ELU that works better for deep networks “Self-Normalizing” property; can train deep SELU networks without BatchNorm

Derivation takes
91 pages of math
in appendix...

$$\begin{aligned} \alpha &= 1.673263242354377284170429916717 \\ \lambda &= 1.0507009873554804934193349852946 \end{aligned}$$

Activation Functions: Gaussian Error Linear Unit (GELU)



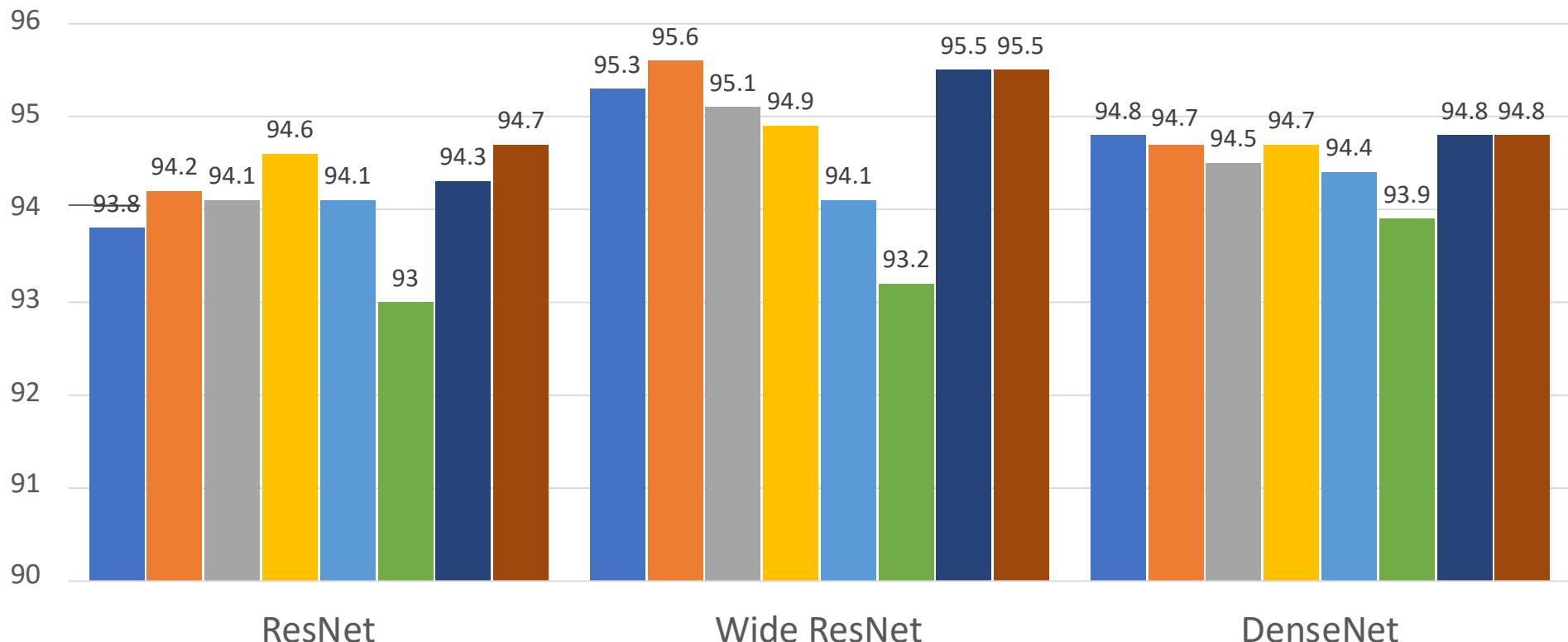
$X \sim N(0, 1)$

$$\begin{aligned} \text{gelu}(x) &= xP(X \leq x) = \frac{x}{2} \left(1 + \operatorname{erf}(x/\sqrt{2}) \right) \\ &\approx x\sigma(1.702x) \end{aligned}$$

- **Idea:** Multiply input by 0 or 1 at random; large values more likely to be multiplied by 1, small values more likely to be multiplied by 0 (data-dependent dropout)
- Take expectation over randomness
- Very common in Transformers (BERT, GPT, GPT-2, GPT-3)

Accuracy on CIFAR10

ReLU Leaky ReLU Parametric ReLU Softplus ELU SELU GELU Swish



Activation Functions: Summary

- Don't think too hard. Just use ReLU
- Try out Leaky ReLU / ELU / SELU / GELU if you need to squeeze that last 0.1%
- Don't use sigmoid or tanh

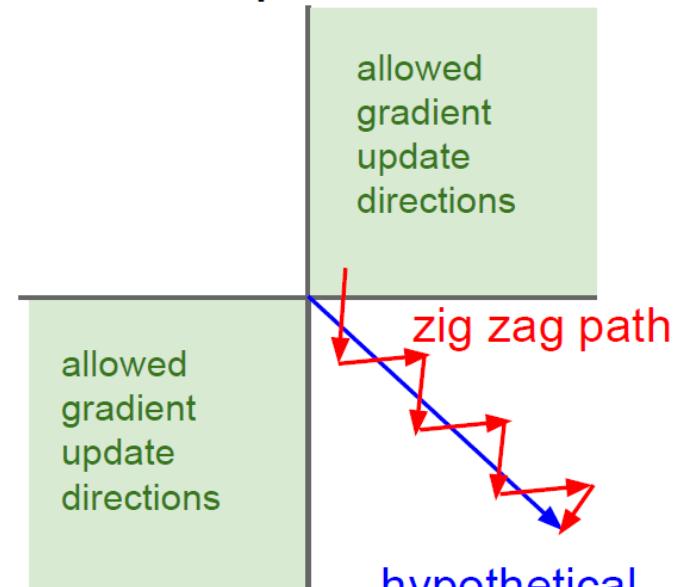
Some (very) recent architectures use GeLU instead of ReLU, but the gains are minimal

Dosovitskiy et al, "An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale", ICLR 2021
Liu et al, "A ConvNet for the 2020s", arXiv 2022

Data Preprocessing

Remember: Consider what happens when the input to a neuron is always positive...

$$f \left(\sum_i w_i x_i + b \right)$$



What can we say about the gradients on w ?

Always all positive or all negative :(

(this is also why you want zero-mean data!)

hypothetical
optimal w
vector

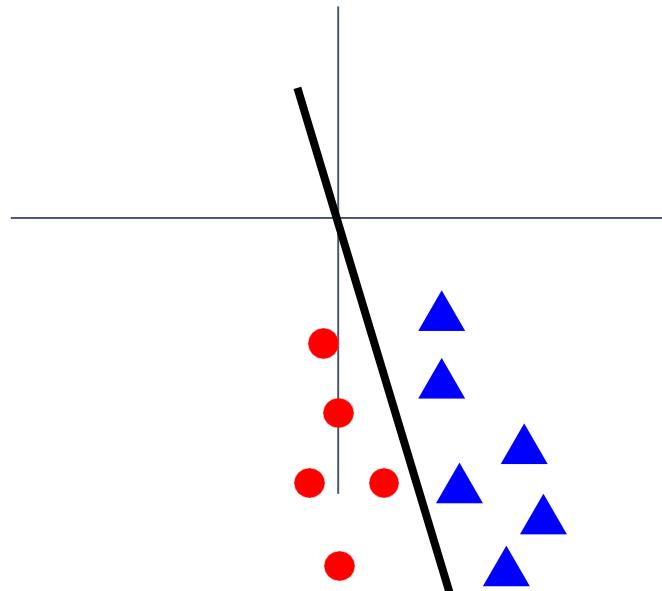
allowed
gradient
update
directions

allowed
gradient
update
directions

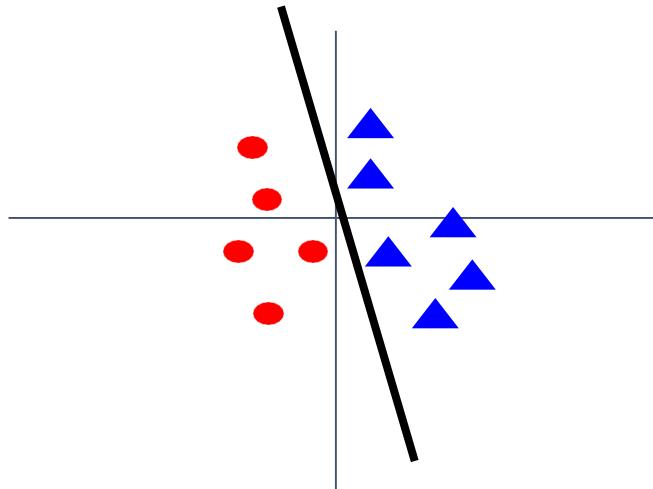
zig zag path

Data Preprocessing

Before normalization: classification loss very sensitive to changes in weight matrix; hard to optimize

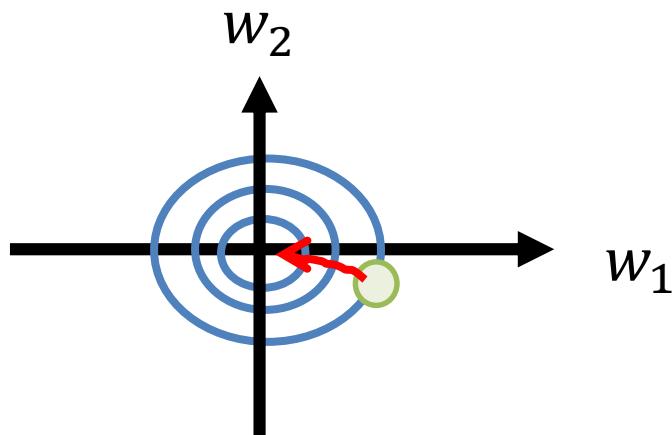
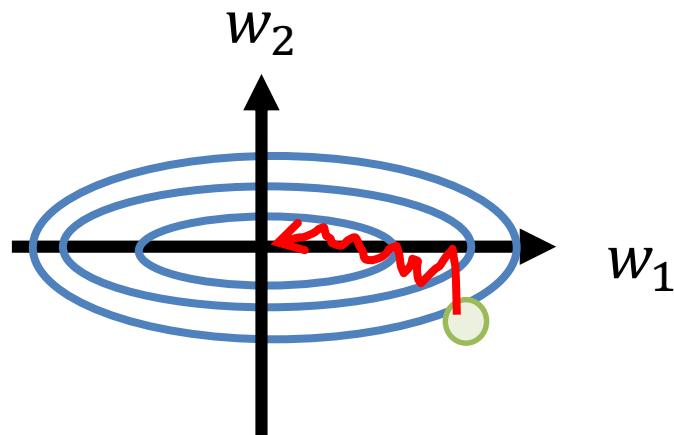


After normalization: less sensitive to small changes in weights; easier to optimize

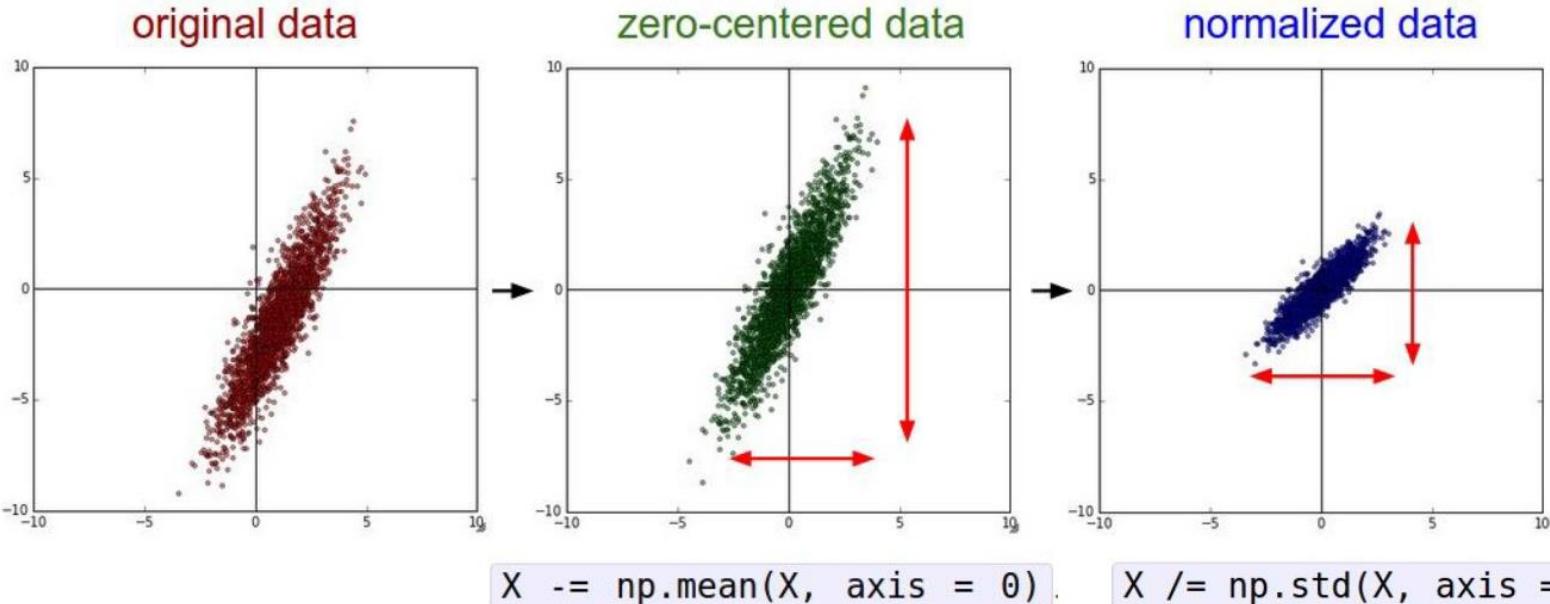


Loss Function

$$\begin{aligned} L &= (w_1 x_1 + w_2 x_2 - y)^2 \\ &= x_1^2 w_1^2 + x_2^2 w_2^2 + \dots \end{aligned}$$



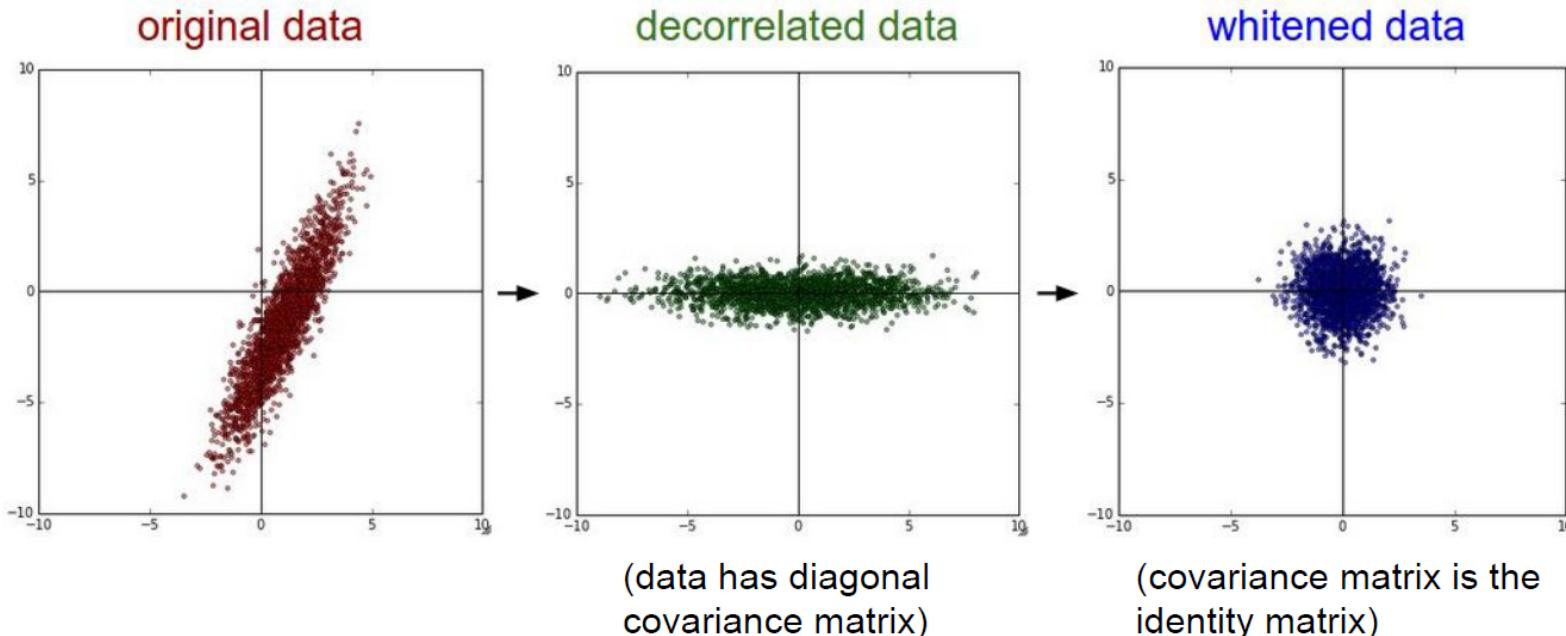
Step 1: Preprocess the data



(Assume $X [NxD]$ is data matrix,
each example in a row)

Step 1: Preprocess the data

In practice, you may also see **PCA** and **Whitening** of the data



Data Preprocessing for Images

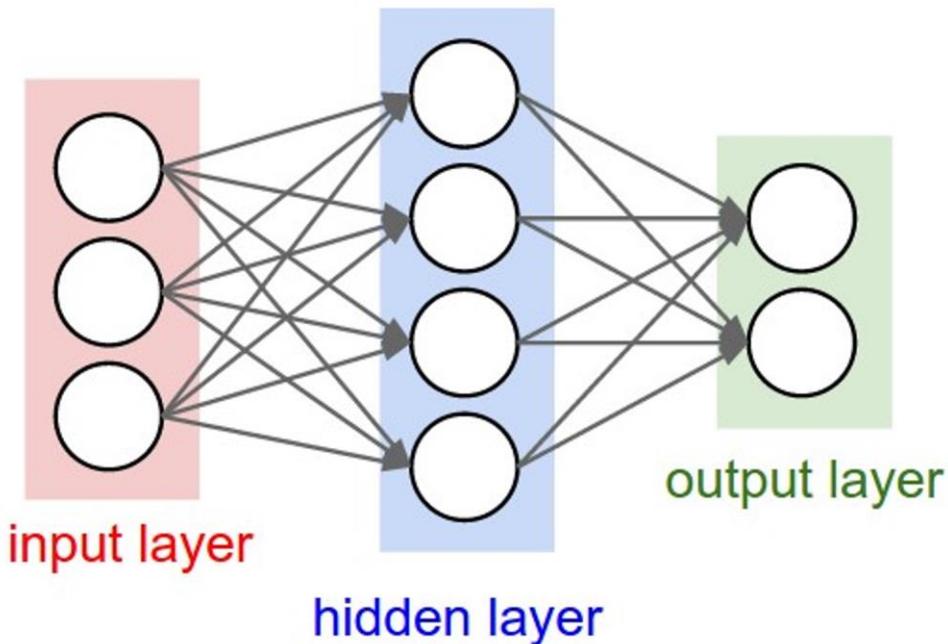
e.g. consider CIFAR-10 example with [32,32,3] images

- Subtract the mean image (e.g. AlexNet)
(mean image = [32,32,3] array)
- Subtract per-channel mean (e.g. VGGNet)
(mean along each channel = 3 numbers)
- Subtract per-channel mean and
Divide by per-channel std (e.g. ResNet)
(mean along each channel = 3 numbers)

Not common to
do PCA or
whitening

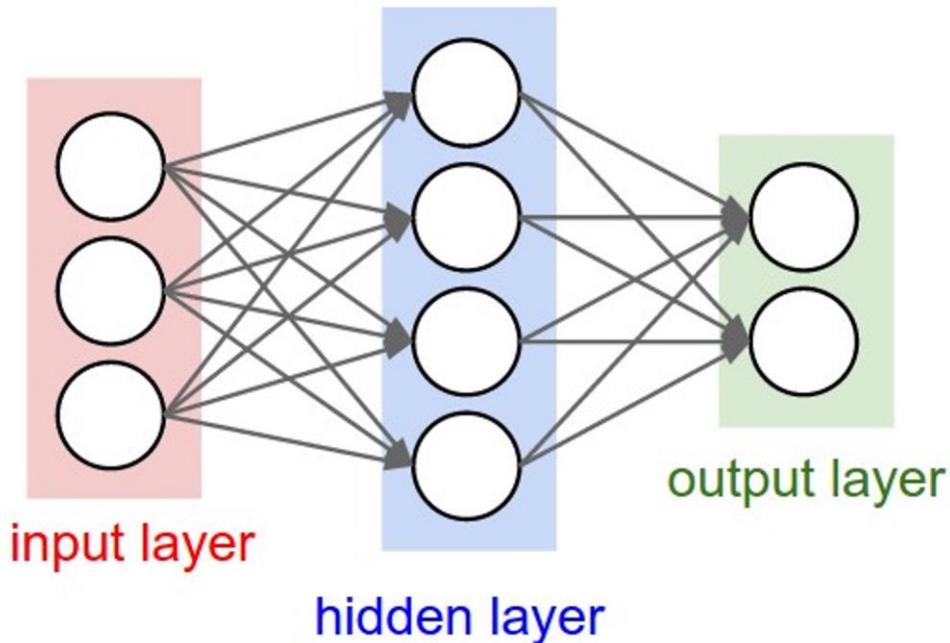
Weight Initialization

Weight Initialization



Q: What happens if we initialize all $W=0$, $b=0$?

Weight Initialization



Q: What happens if we initialize all $W=0$, $b=0$?

A: All outputs are 0, all gradients are the same!
No “symmetry breaking”

Weight Initialization

Next idea: **small random numbers**
(Gaussian with zero mean, std=0.01)

```
W = 0.01 * np.random.randn(Din, Dout)
```

Weight Initialization

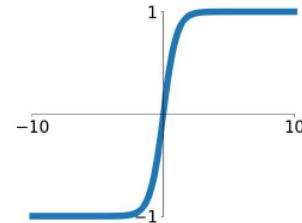
Next idea: **small random numbers**
(Gaussian with zero mean, std=0.01)

```
W = 0.01 * np.random.randn(Din, Dout)
```

Works ~okay for small networks, but
problems with deeper networks.

Weight Initialization: Activation Statistics

```
dims = [4096] * 7      Forward pass for a 6-layer  
hs = []                  net with hidden size 4096  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = 0.01 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```



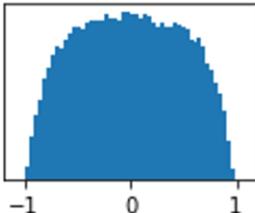
Weight Initialization: Activation Statistics

```
dims = [4096] * 7      Forward pass for a 6-layer  
hs = []                  net with hidden size 4096  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = 0.01 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

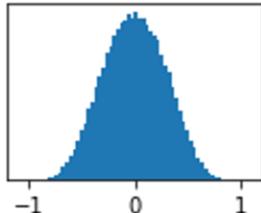
All activations tend to zero for deeper network layers

Q: What do the gradients dL/dW look like?

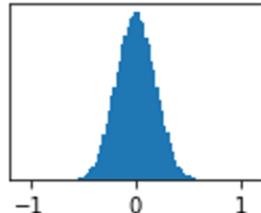
Layer 1
mean=-0.00
std=0.49



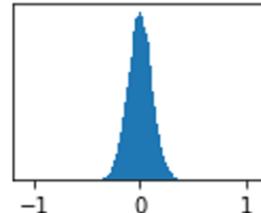
Layer 2
mean=0.00
std=0.29



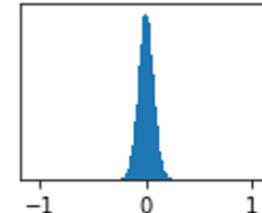
Layer 3
mean=0.00
std=0.18



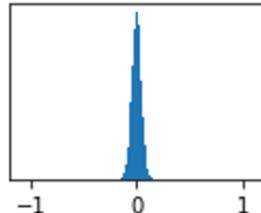
Layer 4
mean=-0.00
std=0.11



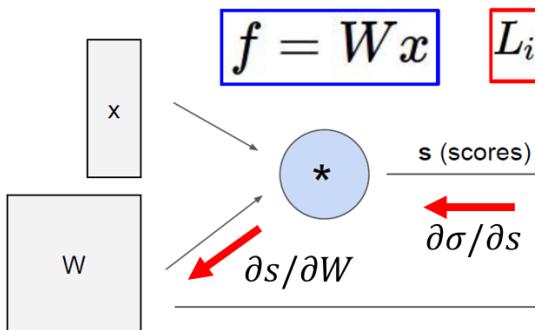
Layer 5
mean=-0.00
std=0.07



Layer 6
mean=0.00
std=0.05



Weight Initialization: Activation Statistics

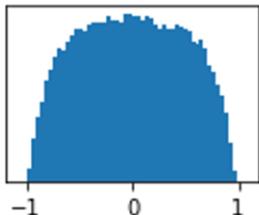


All activations tend to zero for deeper network layers

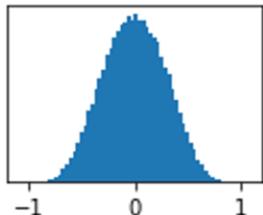
Q: What do the gradients dL/dW look like?

A: All zero, no learning =(

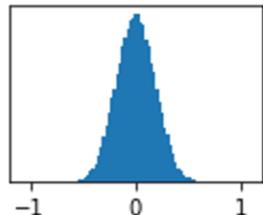
Layer 1
mean=-0.00
std=0.49



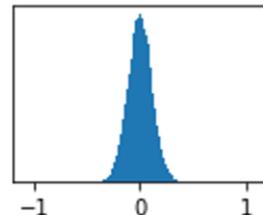
Layer 2
mean=0.00
std=0.29



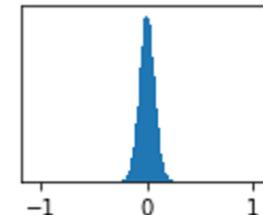
Layer 3
mean=0.00
std=0.18



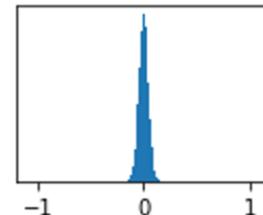
Layer 4
mean=-0.00
std=0.11



Layer 5
mean=-0.00
std=0.07



Layer 6
mean=0.00
std=0.05



Weight Initialization: Activation Statistics

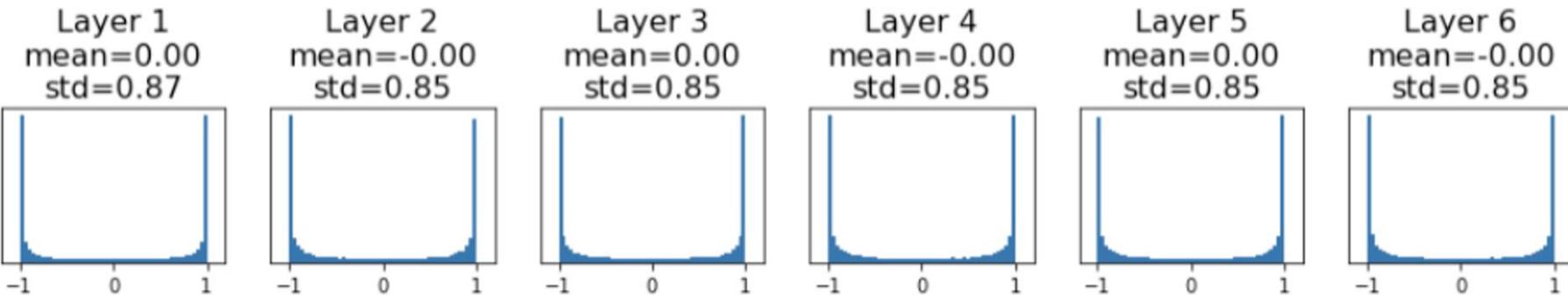
```
dims = [4096] * 7    Increase std of initial weights
hs = []                  from 0.01 to 0.05
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = 0.05 * np.random.randn(Din, Dout)
    x = np.tanh(x.dot(W))
    hs.append(x)
```

Weight Initialization: Activation Statistics

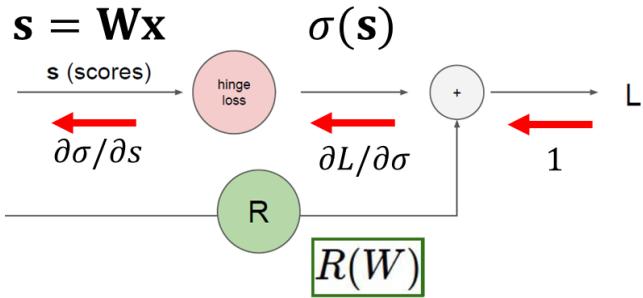
```
dims = [4096] * 7      Increase std of initial weights  
hs = []                  from 0.01 to 0.05  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = 0.05 * np.random.randn(Din, Dout)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

All activations saturate

Q: What do the gradients look like?



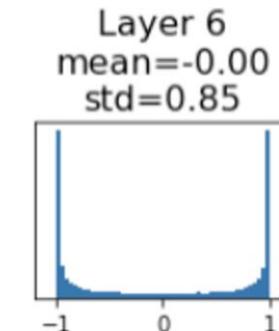
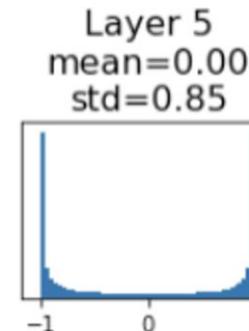
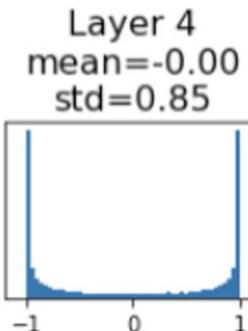
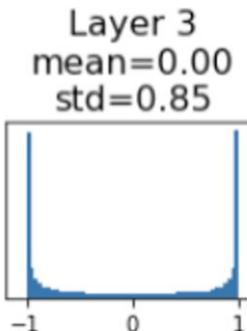
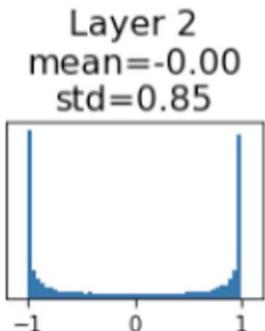
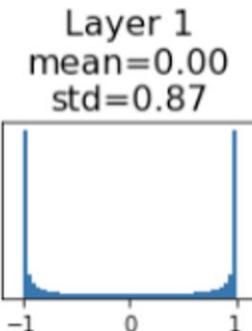
Weight Initialization: Activation Statistics



All activations saturate

Q: What do the gradients look like?

A: Local gradients all zero, no learning =(



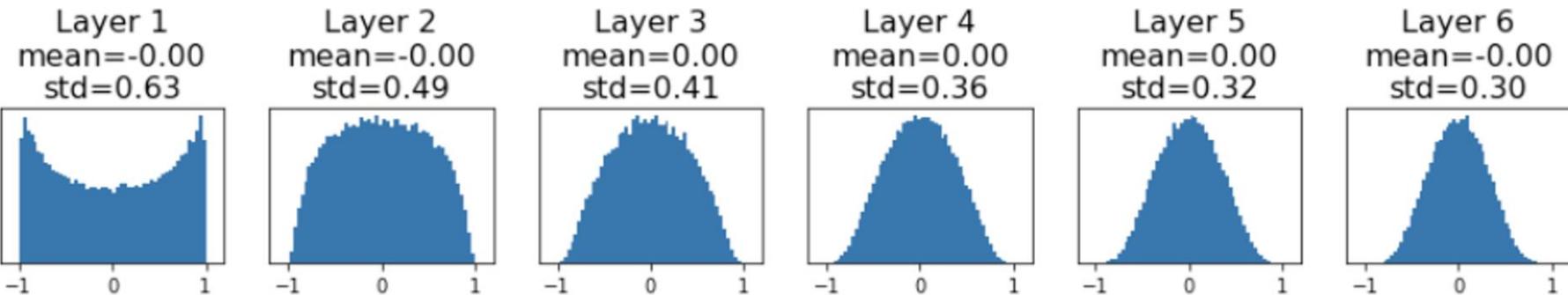
Weight Initialization: Xavier Initialization

```
dims = [4096] * 7          "Xavier" initialization:  
hs = []                      std = 1/sqrt(Din)  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = np.random.randn(Din, Dout) / np.sqrt(Din)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

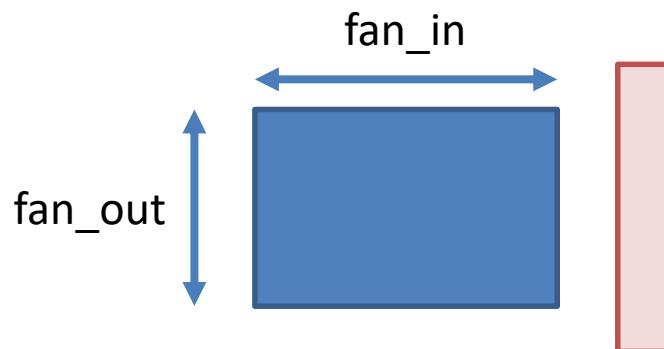
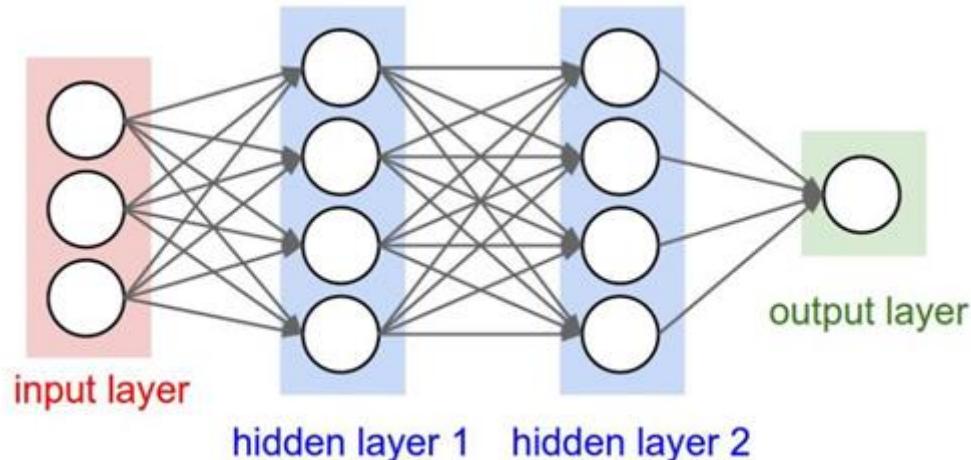
Weight Initialization: Xavier Initialization

```
dims = [4096] * 7          "Xavier" initialization:  
hs = []                      std = 1/sqrt(Din)  
x = np.random.randn(16, dims[0])  
for Din, Dout in zip(dims[:-1], dims[1:]):  
    W = np.random.randn(Din, Dout) / np.sqrt(Din)  
    x = np.tanh(x.dot(W))  
    hs.append(x)
```

"Just right": Activations are nicely scaled for all layers!



Glorot and Bengio, "Understanding the difficulty of training deep feedforward neural networks", AISTAT 2010



$$w_1x_1 + w_2x_2 + \cdots + w_Nx_N$$

$$w_1 \sim N(0, 1)$$

$$w_2 \sim N(0, 1)$$

$$w_1 + w_2 \sim N(0, 2)$$

$$\frac{w_1}{\sqrt{N}} + \frac{w_2}{\sqrt{N}} \sim N\left(0, \frac{2}{N}\right)$$

Weight Initialization: Xavier Initialization

“Xavier” initialization:
std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

Weight Initialization: Xavier Initialization

“Xavier” initialization:
std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\text{Var}(y_i) = \text{Din} * \text{Var}(x_i w_i)$$

[Assume x, w are iid]

Weight Initialization: Xavier Initialization

“Xavier” initialization:
std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\text{Var}(y_i) = \text{Din} * \text{Var}(x_i w_i) \quad [\text{Assume } x, w \text{ are iid}]$$

$$= \text{Din} * (\mathbb{E}[x_i^2] \mathbb{E}[w_i^2] - \mathbb{E}[x_i]^2 \mathbb{E}[w_i]^2) \quad [\text{Assume } x, w \text{ independent}]$$

Weight Initialization: Xavier Initialization

“Xavier” initialization:
std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\begin{aligned} \text{Var}(y_i) &= \text{Din} * \text{Var}(x_i w_i) && [\text{Assume } x, w \text{ are iid}] \\ &= \text{Din} * (\mathbb{E}[x_i^2] \mathbb{E}[w_i^2] - \mathbb{E}[x_i]^2 \mathbb{E}[w_i]^2) && [\text{Assume } x, w \text{ independent}] \\ &= \text{Din} * \text{Var}(x_i) * \text{Var}(w_i) && [\text{Assume } x, w \text{ are zero-mean}] \end{aligned}$$

Weight Initialization: Xavier Initialization

“Xavier” initialization:
std = 1/sqrt(Din)

Derivation: Variance of output = Variance of input

$$y = Wx$$

$$y_i = \sum_{j=1}^{Din} x_j w_j$$

$$\begin{aligned} \text{Var}(y_i) &= \text{Din} * \text{Var}(x_i w_i) && [\text{Assume } x, w \text{ are iid}] \\ &= \text{Din} * (\mathbb{E}[x_i^2] \mathbb{E}[w_i^2] - \mathbb{E}[x_i]^2 \mathbb{E}[w_i]^2) && [\text{Assume } x, w \text{ independent}] \\ &= \text{Din} * \text{Var}(x_i) * \text{Var}(w_i) && [\text{Assume } x, w \text{ are zero-mean}] \end{aligned}$$

If $\text{Var}(w_i) = 1/\text{Din}$ then $\text{Var}(y_i) = \text{Var}(x_i)$

Weight Initialization: What about ReLU?

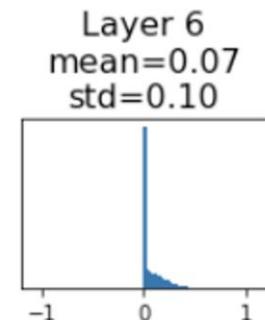
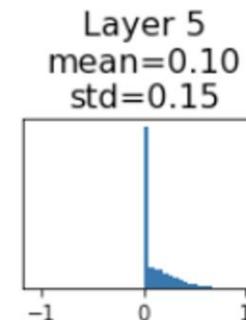
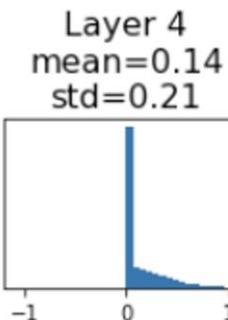
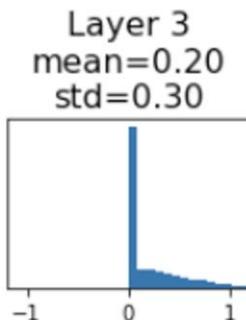
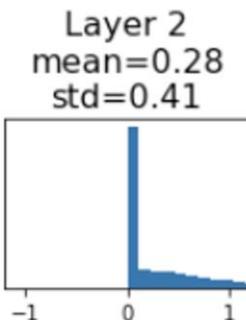
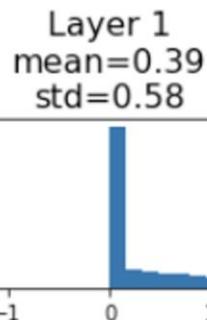
```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

Weight Initialization: What about ReLU?

```
dims = [4096] * 7      Change from tanh to ReLU
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

Xavier assumes zero centered activation function

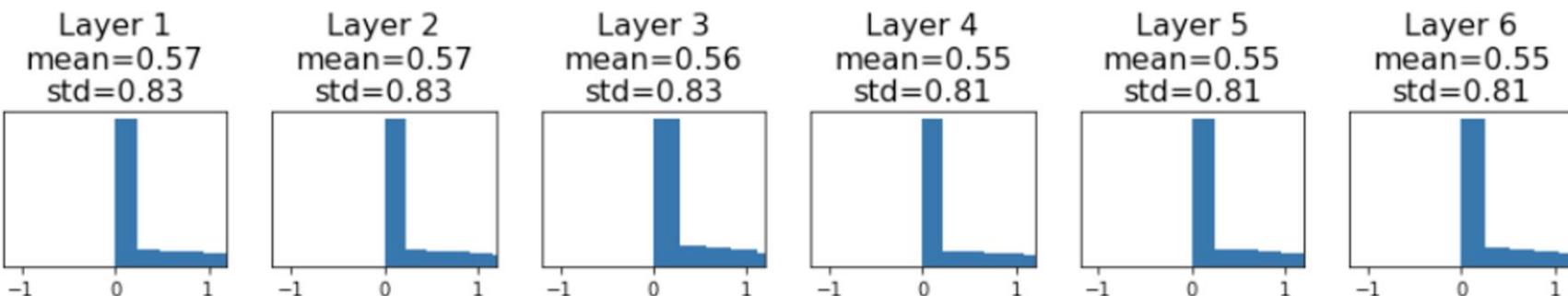
Activations collapse to zero again, no learning =(



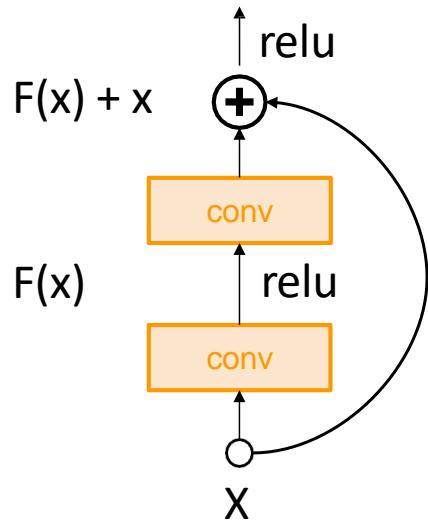
Weight Initialization: Kaiming / MSRA Initialization

```
dims = [4096] * 7 # ReLU correction: std = sqrt(2 / Din)
hs = []
x = np.random.randn(16, dims[0])
for Din, Dout in zip(dims[:-1], dims[1:]):
    W = np.random.randn(Din, Dout) / np.sqrt(Din)
    x = np.maximum(0, x.dot(W))
    hs.append(x)
```

"Just right" – activations nicely scaled for all layers



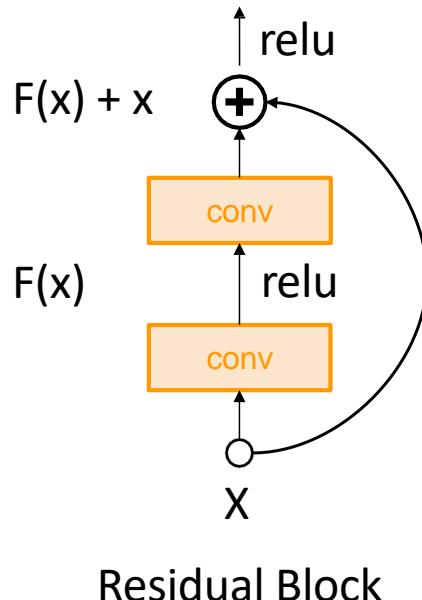
Weight Initialization: Residual Networks



Residual Block

If we initialize with MSRA:
then $\text{Var}(F(x)) = \text{Var}(x)$
But then $\text{Var}(F(x) + x) > \text{Var}(x)$
variance grows with each block!

Weight Initialization: Residual Networks



If we initialize with MSRA:
then $\text{Var}(F(x)) = \text{Var}(x)$
But then $\text{Var}(F(x) + x) > \text{Var}(x)$
variance grows with each block!

Solution: Initialize first conv with MSRA, initialize second conv to zero. Then $\text{Var}(x + F(x)) = \text{Var}(x)$

Proper initialization is an active area of research

Understanding the difficulty of training deep feedforward neural networks by Glorot and Bengio, 2010

Exact solutions to the nonlinear dynamics of learning in deep linear neural networks by Saxe et al, 2013

Random walk initialization for training very deep feedforward networks by Sussillo and Abbott, 2014

Delving deep into rectifiers: Surpassing human-level performance on ImageNet classification by He et al., 2015

Data-dependent Initializations of Convolutional Neural Networks by Krähenbühl et al., 2015

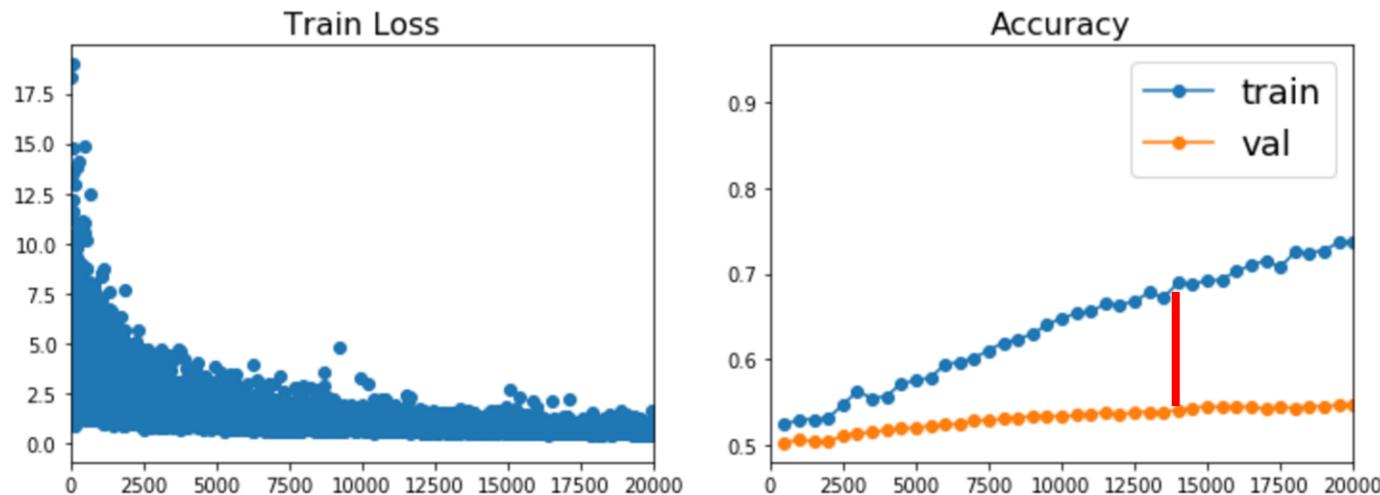
All you need is a good init, Mishkin and Matas, 2015

Fixup Initialization: Residual Learning Without Normalization, Zhang et al, 2019

The Lottery Ticket Hypothesis: Finding Sparse, Trainable Neural Networks, Frankle and Carbin, 2019

Regularization

Now your model is training ... but it overfits!



Regularization

Regularization: Add term to the loss

$$L = \frac{1}{N} \sum_{i=1}^N \sum_{j \neq y_i} \max(0, f(x_i; W)_j - f(x_i; W)_{y_i} + 1) + \boxed{\lambda R(W)}$$

In common use:

L2 regularization

$$R(W) = \sum_k \sum_l W_{k,l}^2 \quad (\text{Weight decay})$$

L1 regularization

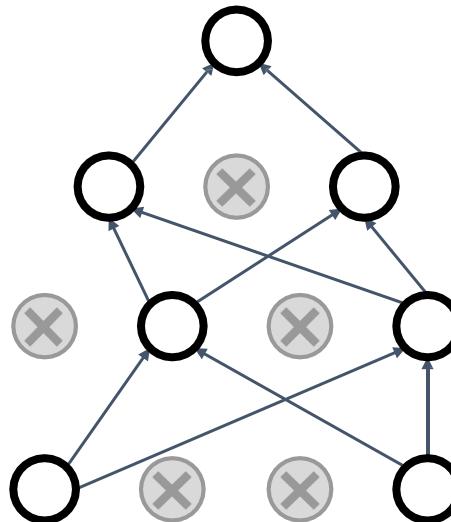
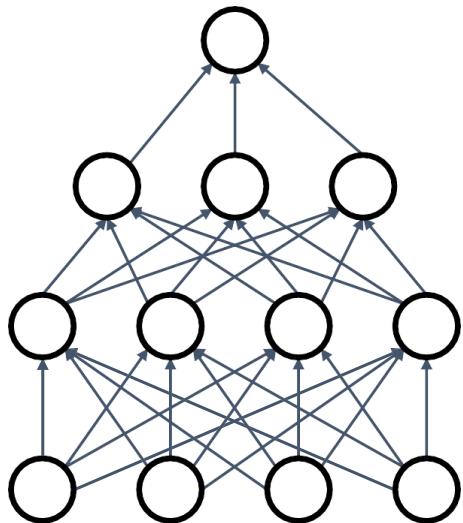
$$R(W) = \sum_k \sum_l |W_{k,l}|$$

Elastic net (L1 + L2)

$$R(W) = \sum_k \sum_l \beta W_{k,l}^2 + |W_{k,l}|$$

Regularization: Dropout

In each forward pass, randomly set some neurons to zero
Probability of dropping is a hyperparameter; 0.5 is common



Regularization: Dropout

```
def train_step(X):
    """ X contains the data """

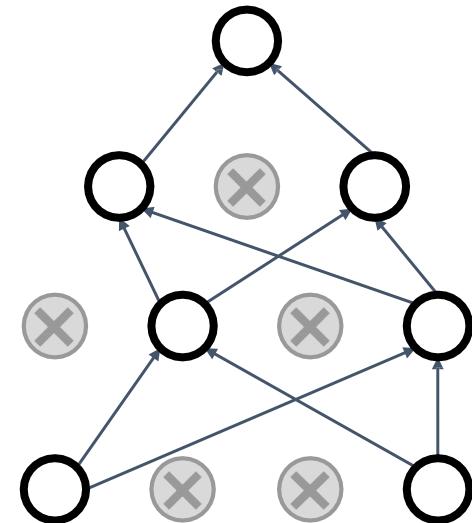
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)

    H2 = np.maximum(0, np.dot(W2, H1) + b2)

    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout



Regularization: Dropout

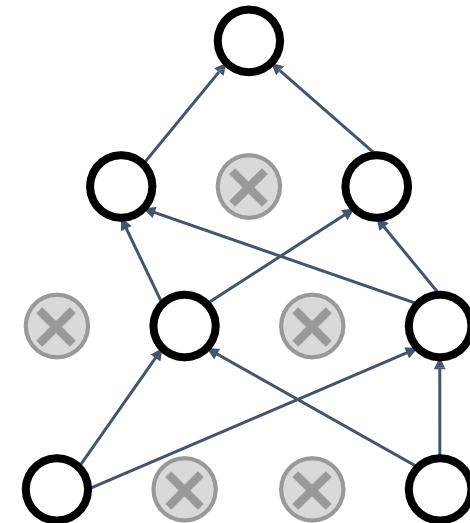
```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    """ X contains the data """

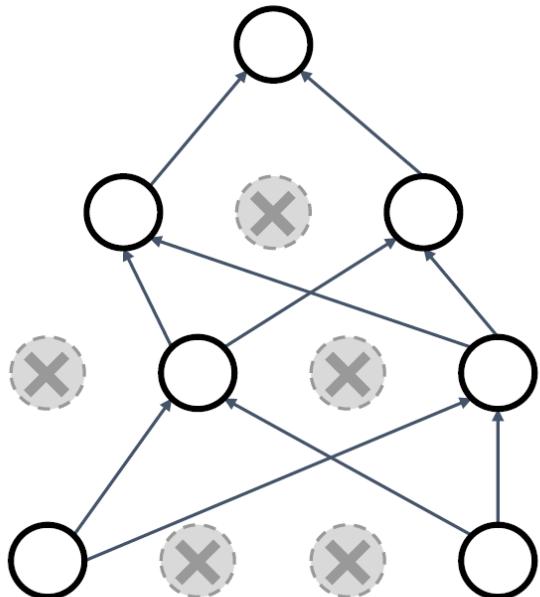
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = np.random.rand(*H1.shape) < p # first dropout mask
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = np.random.rand(*H2.shape) < p # second dropout mask
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)
```

Example forward pass with a 3-layer network using dropout



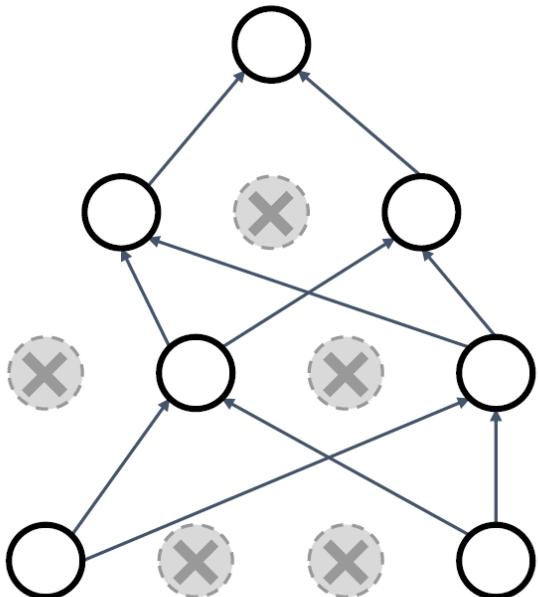
Regularization: Dropout



Forces the network to have a redundant representation; Prevents **co-adaptation** of features



Regularization: Dropout



Another interpretation:

Dropout is training a large **ensemble** of models (that share parameters).

Each binary mask is one model

An FC layer with 4096 units has
 $2^{4096} \sim 10^{1233}$ possible masks!
Only $\sim 10^{82}$ atoms in the universe...

Dropout: Test Time

Output
(label) Input
(image)

Dropout makes our output random!

$$\textcolor{red}{y} = f_W(\textcolor{blue}{x}, \textcolor{green}{z}) \quad \text{Random mask}$$

Want to “average out” the randomness at test-time

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

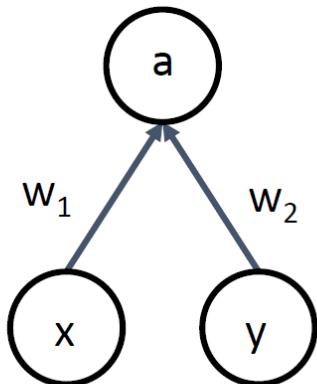
But this integral seems hard ...

Dropout: Test Time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Consider a single neuron:

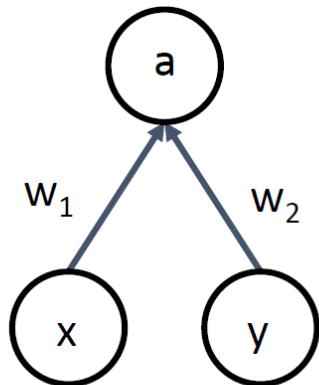


At test time we have: $E[a] = w_1 x + w_2 y$

Dropout: Test Time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$



Consider a single neuron:

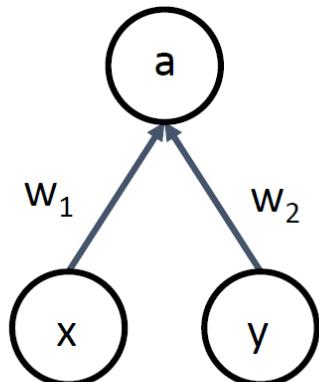
At test time we have: $E[a] = w_1x + w_2y$

$$\begin{aligned} \text{During training we have: } E[a] &= \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y) \\ &\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y) \\ &= \frac{1}{2}(w_1x + w_2y) \end{aligned}$$

Dropout: Test Time

Want to approximate
the integral

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$



Consider a single neuron:

At test time we have: $E[a] = w_1x + w_2y$

During training we have: $E[a] = \frac{1}{4}(w_1x + w_2y) + \frac{1}{4}(w_1x + 0y)$
 $\quad\quad\quad + \frac{1}{4}(0x + 0y) + \frac{1}{4}(0x + w_2y)$
 $\quad\quad\quad = \frac{1}{2}(w_1x + w_2y)$

At test time, drop
nothing and multiply
by dropout probability

Dropout: Test Time

```
def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations
    out = np.dot(W3, H2) + b3
```

At test time all neurons are active always

=> We must scale the activations so that for each neuron:

output at test time = expected output at training time

Dropout

Summary

```
""" Vanilla Dropout: Not recommended implementation (see notes below) """  
  
p = 0.5 # probability of keeping a unit active. higher = less dropout  
  
def train_step(X):  
    """ X contains the data """  
  
    # forward pass for example 3-layer neural network  
    H1 = np.maximum(0, np.dot(W1, X) + b1)  
    U1 = np.random.rand(*H1.shape) < p # first dropout mask  
    H1 *= U1 # drop!  
    H2 = np.maximum(0, np.dot(W2, H1) + b2)  
    U2 = np.random.rand(*H2.shape) < p # second dropout mask  
    H2 *= U2 # drop!  
    out = np.dot(W3, H2) + b3  
  
    # backward pass: compute gradients... (not shown)  
    # perform parameter update... (not shown)  
  
def predict(X):  
    # ensembled forward pass  
    H1 = np.maximum(0, np.dot(W1, X) + b1) * p # NOTE: scale the activations  
    H2 = np.maximum(0, np.dot(W2, H1) + b2) * p # NOTE: scale the activations  
    out = np.dot(W3, H2) + b3
```

drop in forward pass

scale at test time

More common: “Inverted dropout”

```
p = 0.5 # probability of keeping a unit active. higher = less dropout

def train_step(X):
    # forward pass for example 3-layer neural network
    H1 = np.maximum(0, np.dot(W1, X) + b1)
    U1 = (np.random.rand(*H1.shape) < p) / p # first dropout mask. Notice /p!
    H1 *= U1 # drop!
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    U2 = (np.random.rand(*H2.shape) < p) / p # second dropout mask. Notice /p!
    H2 *= U2 # drop!
    out = np.dot(W3, H2) + b3

    # backward pass: compute gradients... (not shown)
    # perform parameter update... (not shown)

def predict(X):
    # ensembled forward pass
    H1 = np.maximum(0, np.dot(W1, X) + b1) # no scaling necessary
    H2 = np.maximum(0, np.dot(W2, H1) + b2)
    out = np.dot(W3, H2) + b3
```

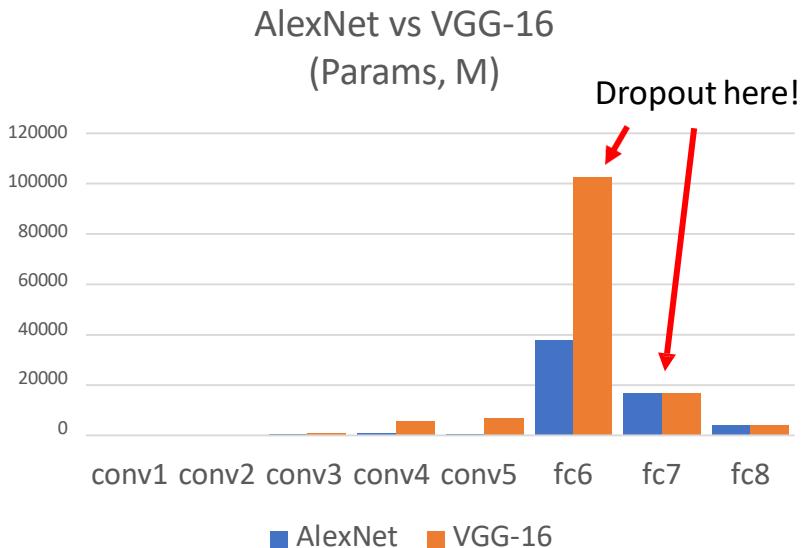
Drop and scale
during training

test time is unchanged!



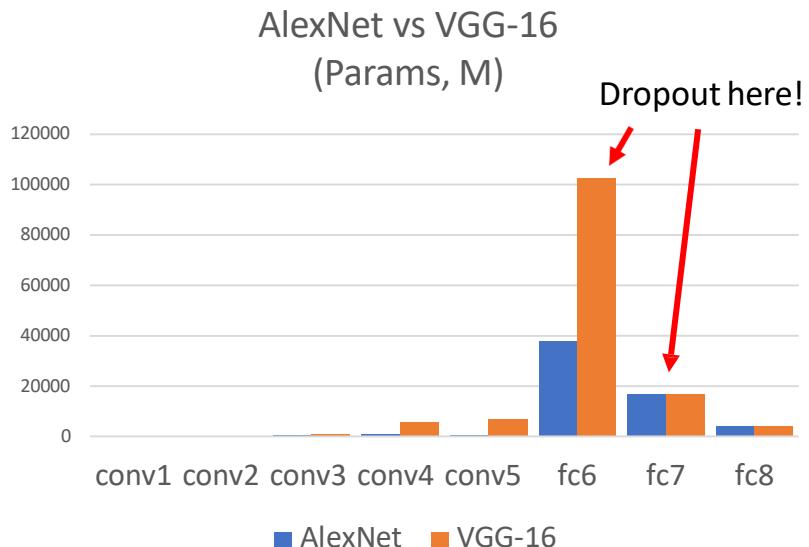
Dropout architectures

Recall AlexNet, VGG have most of their parameters in **fully-connected layers**; usually Dropout is applied there



Dropout architectures

Recall AlexNet, VGG have most of their parameters in **fully-connected layers**; usually Dropout is applied there



Later architectures (GoogLeNet, ResNet, etc) use global average pooling instead of fully-connected layers: they don't use dropout at all!

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness
(sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

Testing: Average out randomness
(sometimes approximate)

$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Example: Batch Normalization

Training: Normalize using stats from random minibatches

Testing: Use fixed stats to normalize

Regularization: A common pattern

Training: Add some kind of randomness

$$y = f_W(x, z)$$

For ResNet and later,
often L2 and Batch
Normalization are
the only regularizers!

Testing: Average out randomness
(sometimes approximate)

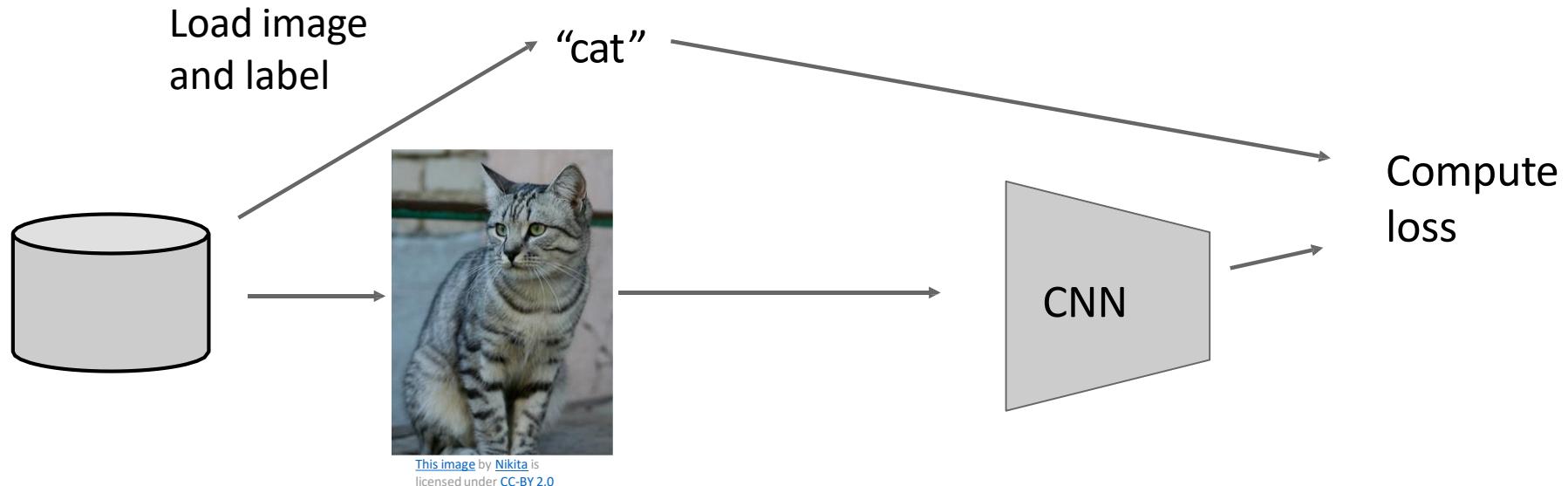
$$y = f(x) = E_z[f(x, z)] = \int p(z)f(x, z)dz$$

Example: Batch Normalization

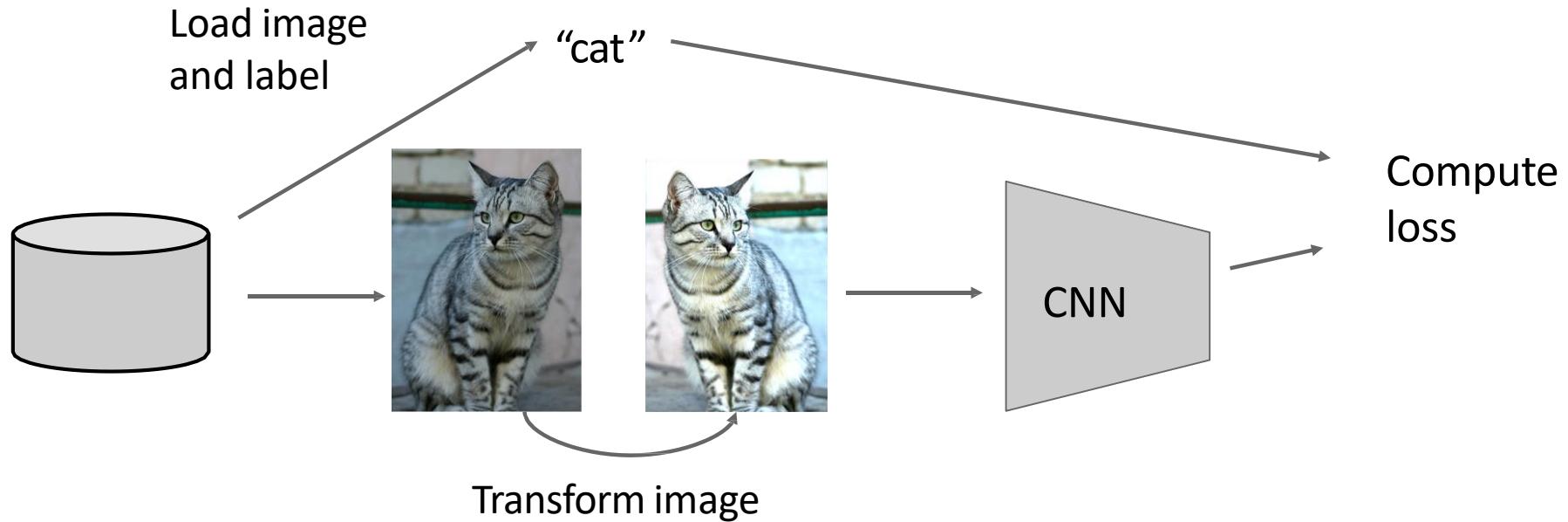
Training: Normalize using stats from random minibatches

Testing: Use fixed stats to normalize

Data Augmentation



Data Augmentation



Data Augmentation: Horizontal Flips

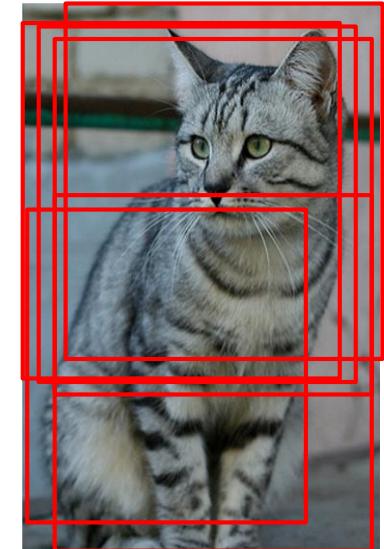


Data Augmentation: Random Crops and Scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range $[256, 480]$
2. Resize training image, short side = L
3. Sample random 224×224 patch

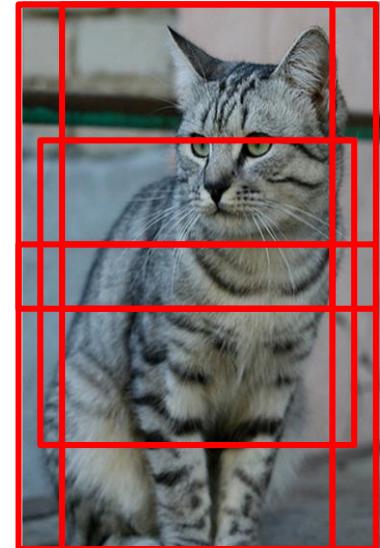


Data Augmentation: Random Crops and Scales

Training: sample random crops / scales

ResNet:

1. Pick random L in range [256, 480]
2. Resize training image, short side = L
3. Sample random 224×224 patch



Testing: average a fixed set of crops

ResNet:

1. Resize image at 5 scales: {224, 256, 384, 480, 640}
2. For each size, use 10 224×224 crops: 4 corners + center, + flips

Data Augmentation: Color Jitter

Simple: Randomize
contrast and brightness



More Complex:

1. Apply PCA to all [R, G, B] pixels in training set
2. Sample a “color offset” along principal component directions
3. Add offset to all pixels of a training image

(Used in AlexNet, ResNet, etc)

Data Augmentation: RandAugment

Apply random combinations
of transforms:

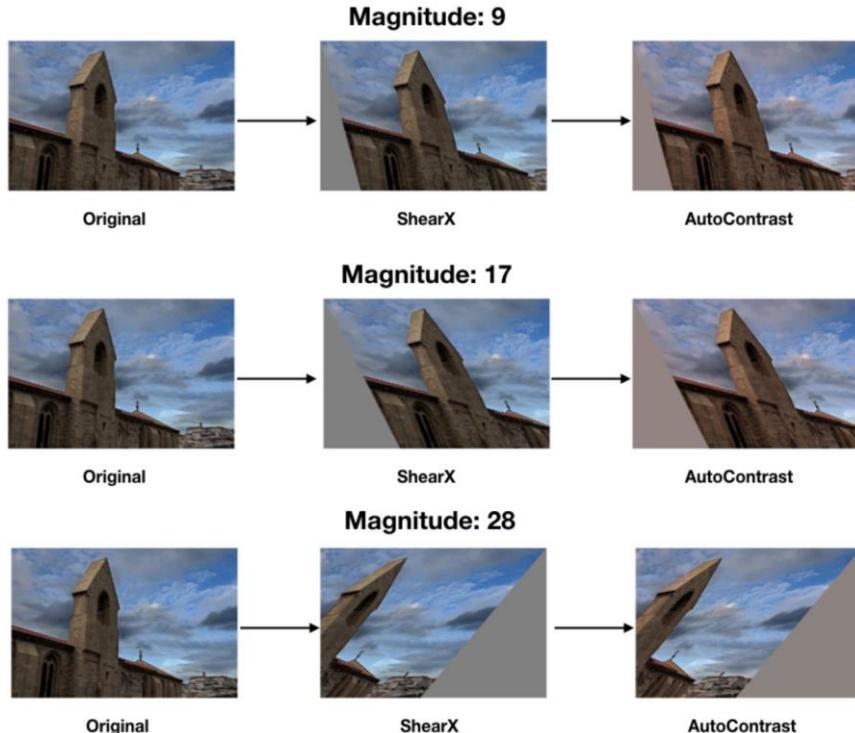
- **Geometric:** Rotate, translate, shear
- **Color:** Sharpen, contrast, brightness, solarize, posterize, color

```
transforms = [  
    'Identity', 'AutoContrast', 'Equalize',  
    'Rotate', 'Solarize', 'Color', 'Posterize',  
    'Contrast', 'Brightness', 'Sharpness',  
    'ShearX', 'ShearY', 'TranslateX', 'TranslateY']  
  
def randaugment(N, M):  
    """Generate a set of distortions.  
  
    Args:  
        N: Number of augmentation transformations to  
            apply sequentially.  
        M: Magnitude for all the transformations.  
    """  
  
    sampled_ops = np.random.choice(transforms, N)  
    return [(op, M) for op in sampled_ops]
```

Data Augmentation: RandAugment

Apply random combinations
of transforms:

- **Geometric:** Rotate, translate, shear
- **Color:** Sharpen, contrast, brightness, solarize, posterize, color



Cubuk et al, “RandAugment: Practical augmented data augmentation with a reduced search space”, NeurIPS 2020

Data Augmentation: Get creative for your problem!

Data augmentation encodes **invariances** in your model

Think for your problem: what changes to the image
should **not** change the network output?

May be different for different tasks!

Regularization: A common pattern

Training: Add some randomness

Testing: Marginalize over randomness

Examples:

Dropout

Batch Normalization

Data Augmentation

Regularization: DropConnect

Training: Drop random connections between neurons (set weight=0)

Testing: Use all the connections

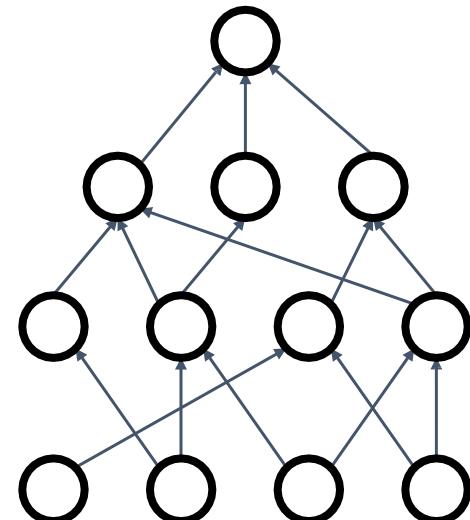
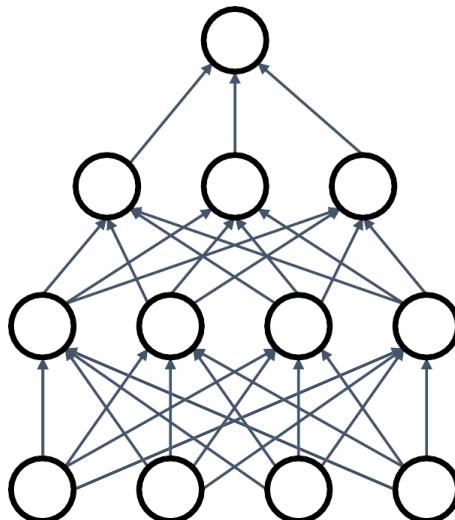
Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect



Regularization: Fractional Pooling

Training: Use randomized pooling regions

Testing: Average predictions over different samples

Examples:

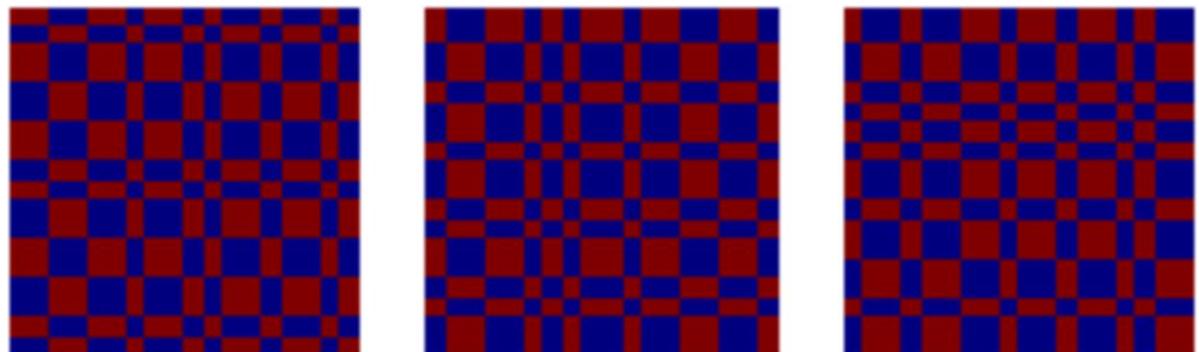
Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling



Regularization: Stochastic Depth

Training: Skip some residual blocks in ResNet

Testing: Use the whole network

Examples:

Dropout

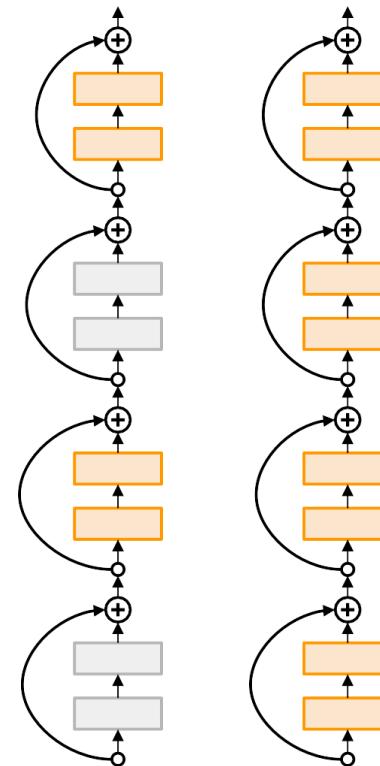
Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth



Regularization: Stochastic Depth

Training: Skip some residual blocks in ResNet

Testing: Use the whole network

Examples:

Dropout

Batch Normalization

Data Augmentation

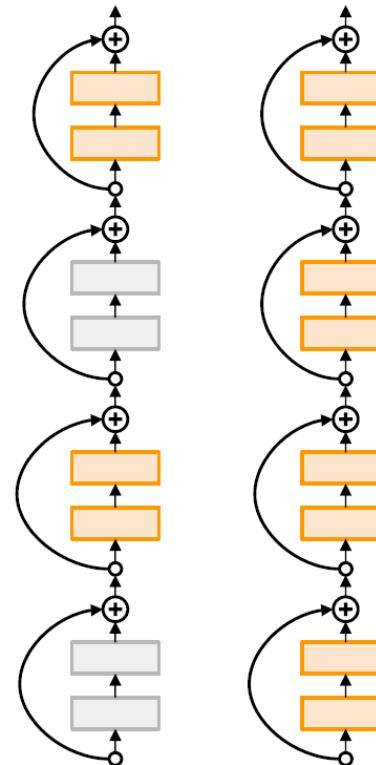
DropConnect

Fractional Max Pooling

Stochastic Depth

Starting to become common in recent architectures!

- Pham et al, "Very Deep Self-Attention Networks for End-to-End Speech Recognition", INTERSPEECH 2019
- Tan and Le, "EfficientNetV2: Smaller Models and Faster Training", ICML 2021
- Fan et al, "Multiscale Vision Transformers", ICCV 2021
- Bello et al, "Revisiting ResNets: Improved Training and Scaling Strategies", NeurIPS 2021
- Steiner et al, "How to train your ViT? Data, Augmentation, and Regularization in Vision Transformers", arXiv 2021



Regularization: CutOut

Training: Set random images regions to 0

Testing: Use the whole image

Examples:

Dropout

Batch Normalization

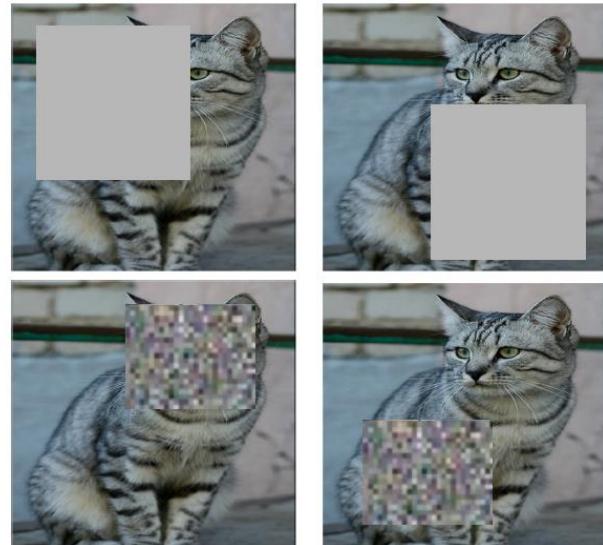
Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing



Replace random regions with
mean value or random values

DeVries and Taylor, "Improved Regularization of Convolutional Neural Networks with Cutout", arXiv 2017

Zhong et al, "Random Erasing Data Augmentation", AAAI 2020

Regularization: Mixup

- **Training:** Train on random blends of images
- **Testing:** Use original images

- **Examples:**

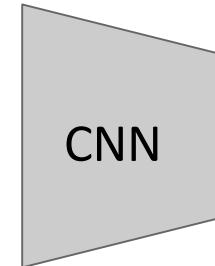
- Dropout
- Batch Normalization
- Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

- Cutout
- Mixup



Target label:
cat: 0.4
dog: 0.6

Randomly blend the pixels of pairs of training images, e.g.
40% cat, 60% dog

Regularization: Mixup

- **Training:** Train on random blends of images
- **Testing:** Use original images

- **Examples:**

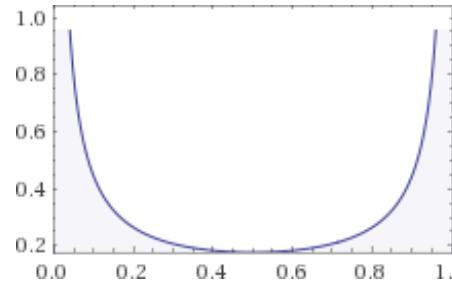
- Dropout
- Batch Normalization
- Data Augmentation

DropConnect

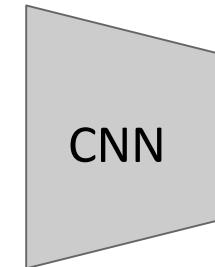
Fractional Max Pooling

Stochastic Depth

- Cutout
- Mixup



Sample blend probability from a beta distribution $\text{Beta}(a, b)$ with $a=b\approx 0$ so blend weights are close to 0/1



Target label:
cat: 0.4
dog: 0.6

Randomly blend the pixels of pairs of training images, e.g. 40% cat, 60% dog

Regularization: CutMix

Training: Train on random blends of images

Testing: Use original images

Examples:

Dropout

Batch Normalization

Data Augmentation

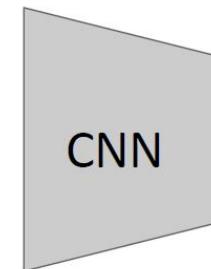
DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix



Target label:
cat: 0.6
dog: 0.4

Replace random crops of one image with another:
e.g. 60% of pixels from cat, 40% from dog

Regularization: Label Smoothing

Training: Train on random blends of images

Testing: Use original images

Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix

Label Smoothing



Target Distribution

Standard Training

Cat: 100%

Dog: 0%

Fish: 0%

Label Smoothing

Cat: 90%

Dog: 5%

Fish: 5%

Set target distribution to be $1 - \frac{K-1}{K}\epsilon$ on the correct category and ϵ/K on all other categories, with K categories and $\epsilon \in (0,1)$. Loss is cross-entropy between predicted and target distribution.

Regularization: Summary

Training: Train on random blends of images

Testing: Use original images

Examples:

Dropout

Batch Normalization

Data Augmentation

DropConnect

Fractional Max Pooling

Stochastic Depth

Cutout / Random Erasing

Mixup / CutMix

Label Smoothing

- Use DropOut for large fully-connected layers
- Data augmentation always a good idea
- Use BatchNorm for CNNs (but not ViTs)
- Try Cutout, MixUp, CutMix, Stochastic Depth, Label Smoothing to squeeze out a bit of extra performance

Summary

1. One time setup

Activation functions, data preprocessing, weight initialization, regularization

Today

2. Training dynamics

Learning rate schedules; large-batch training;
hyperparameter optimization

Next time

3. After training

Model ensembles, transfer learning