# Lecture 3: Introduction to Deep Learning

**Neychev Radoslav**

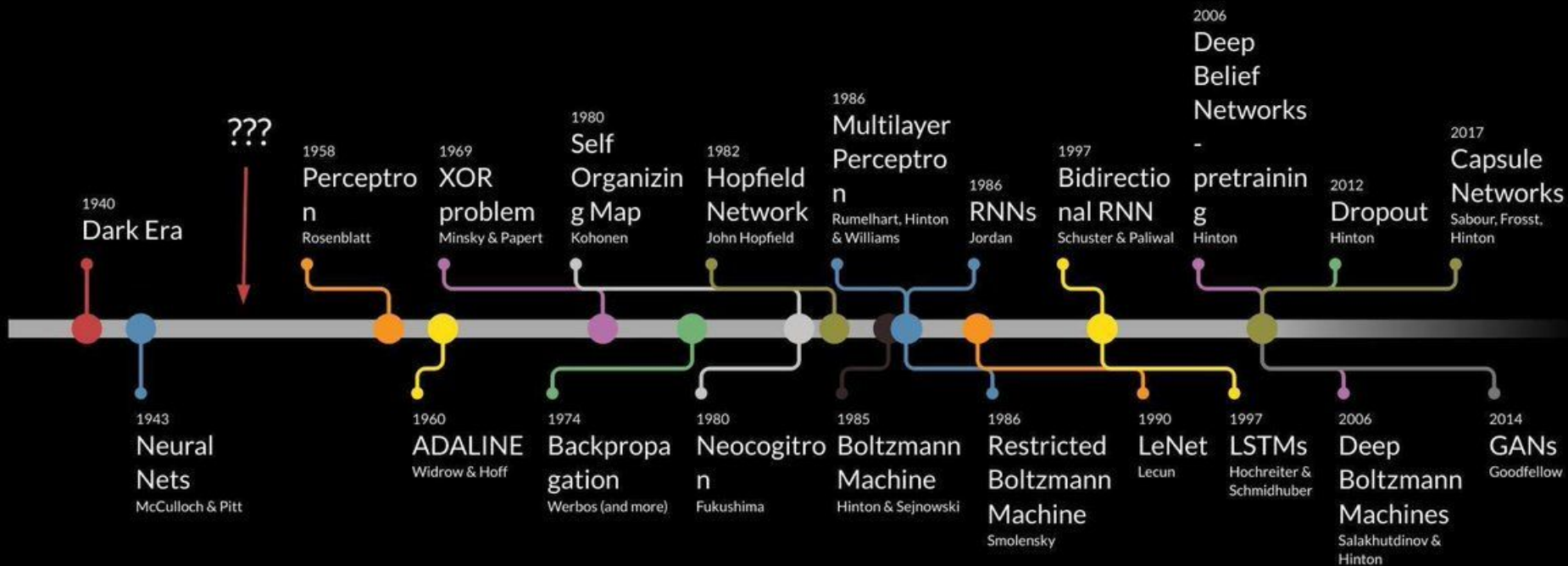ML Instructor (MIPT, HSE, Harbour.Space, BigData Team)
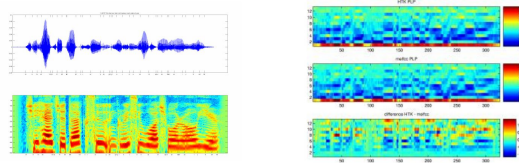
Research Scientist, MIPT

16.10.2019, Moscow, Russia

# Outline

1. Neural Networks in different areas. Historical overview.
2. Backpropagation.
3. Playground.
4. More on backpropagation.
5. Activation functions.
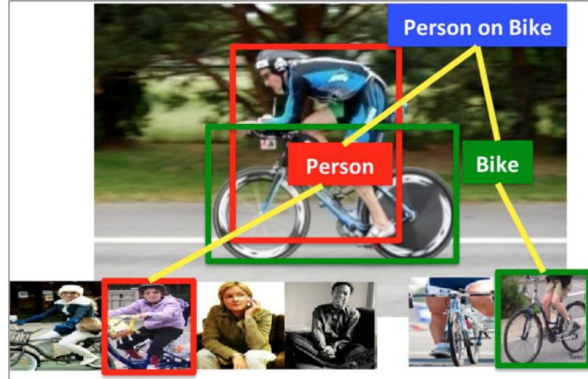6. PyTorch practice

# Deep Learning Timeline



**1940** Dark Era

??? 

**1943** Neural Nets
McCulloch & Pitt

**1958** Perceptron
Rosenblatt

**1960** ADALINE
Widrow & Hoff

**1969** XOR problem
Minsky & Papert

**1974** Backpropagation
Werbos (and more)

**1980** Self Organizing Map
Kohonen

**1980** Neocogitron
Fukushima

**1982** Hopfield Network
John Hopfield

**1985** Boltzmann Machine
Hinton & Sejnowski

**1986** Multilayer Perceptron
Rumelhart, Hinton & Williams

**1986** Restricted Boltzmann Machine
Smolensky

**1986** RNNs
Jordan

**1990** LeNet
Lecun

**1997** Bidirectional RNN
Schuster & Paliwal

**1997** LSTMs
Hochreiter & Schmidhuber

**2006** Deep Belief Networks - pretraining
Hinton

**2006** Deep Boltzmann Machines
Salakhutdinov & Hinton

**2012** Dropout
Hinton

**2014** GANs
Goodfellow

**2017** Capsule Networks
Sabour, Frosst, Hinton

Made by Favio Vázquez

# Real world problems

Audio Features



Spectrogram



MFCC

- Object detection
- Action classification
- Image captioning
- …

person
hammer
flower pot
power drill

Person on Bike

Person

Bike

"man in black shirt is playing guitar."

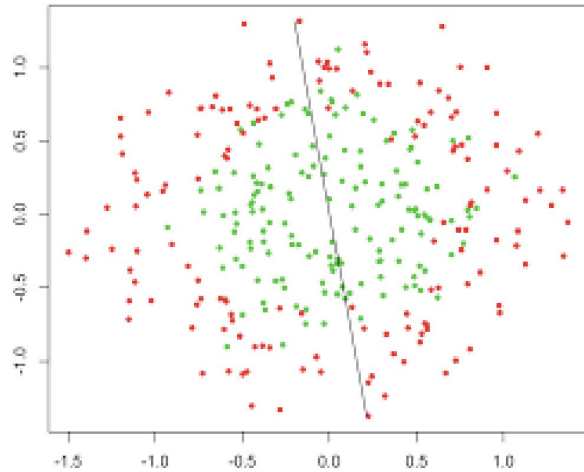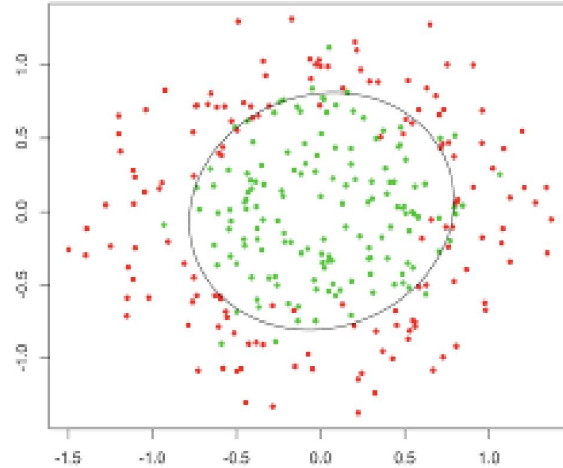$$X \longrightarrow Wx + b \longrightarrow \text{[sigmoid curve]} \longrightarrow P(y)$$

$$P(y|x) = \sigma(w \cdot x + b)$$

$$L = -\sum_i y_i \log P(y|x_i) + (1 - y_i) \log(1 - P(y|x_i))$$
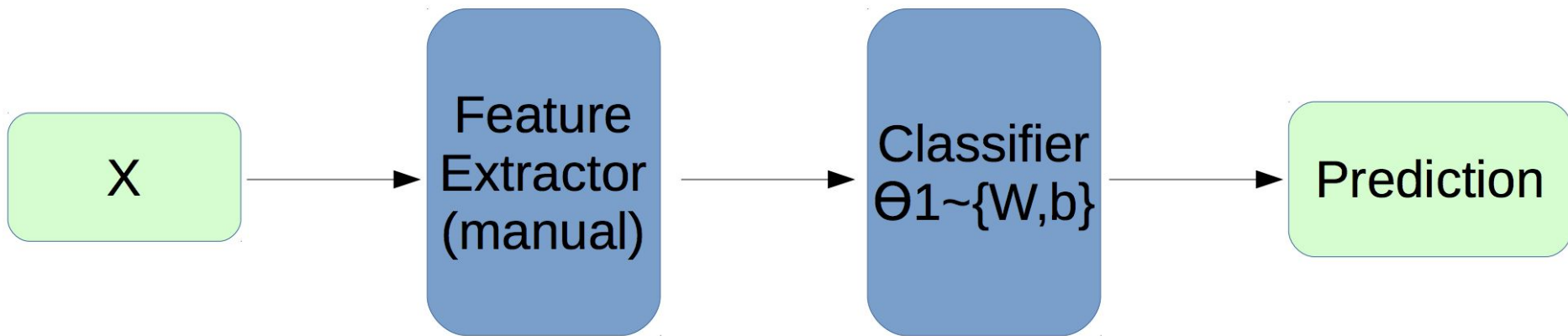
# Problem: nonlinear dependencies



What we have



What we want

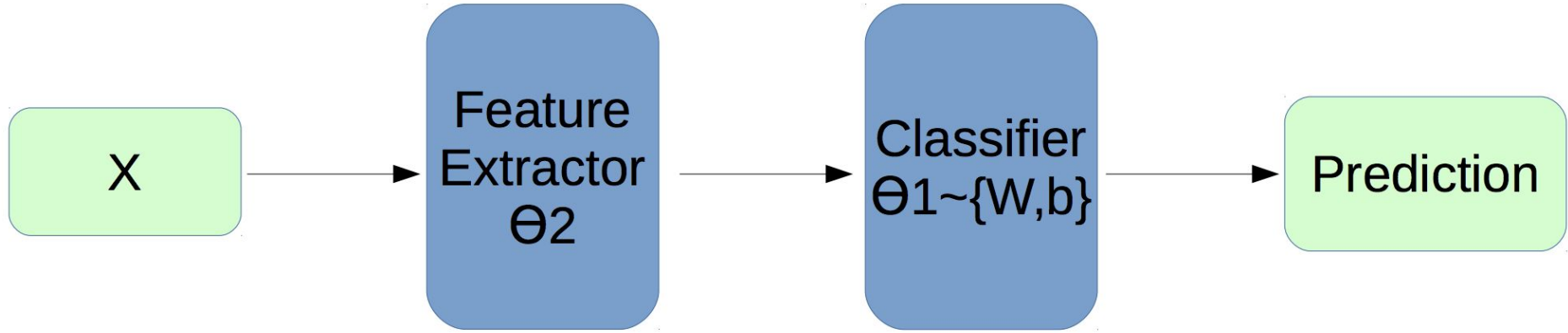Logistic regression (generally, linear model) need feature engineering to show good results.
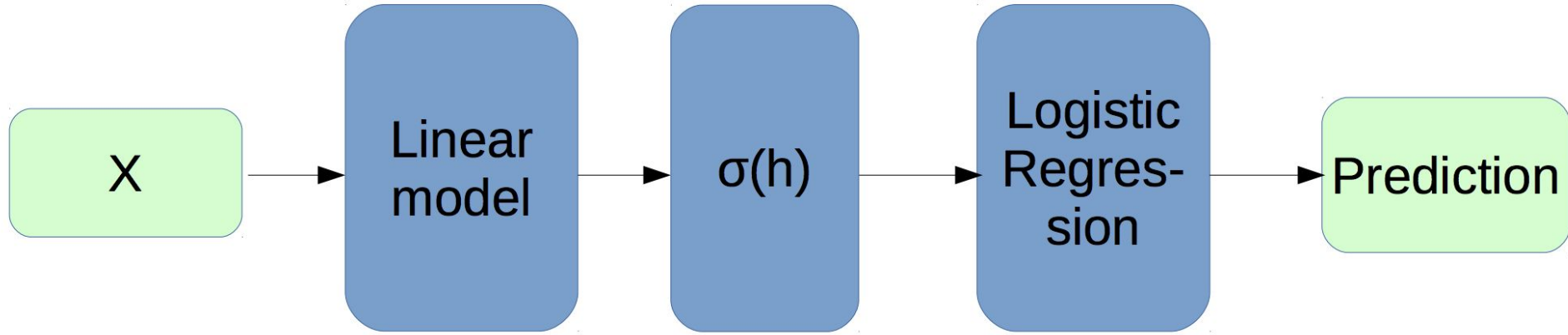
And feature engineering is an *art*.

# Classic pipeline

X → Feature Extractor (manual) → Classifier Θ1~{W,b} → Prediction

Handcrafted features,  generated by experts.

# NN pipeline

$$X \rightarrow \boxed{\text{Feature Extractor } \Theta2} \rightarrow \boxed{\text{Classifier } \Theta1 \sim \{W,b\}} \rightarrow \text{Prediction}$$

Automatically extracted features.

X → Linear model → $\sigma(h)$ → Logistic Regression → Prediction

E.g. two logistic regressions one after another.

```
┌─────────┐      ┌───────────┐      ┌─────────┐      ┌───────────┐      ┌────────────┐
│    X    │ ───► │  Linear   │ ───► │  σ(h)   │ ───► │  Logistic │ ───► │ Prediction │
│         │      │  model    │      │         │      │  Regres-  │      │            │
└─────────┘      └───────────┘      └─────────┘      │   sion    │      └────────────┘
                                                     └───────────┘
```

Actually, it's a neural network.

# XOR problem



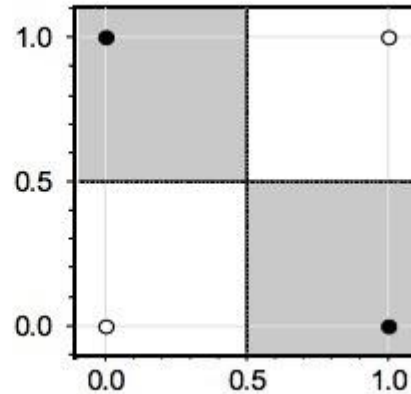This 2-layer NN (on the left) implements XOR with only x1 and x2 features.

1-layer NN also can succeed, but only with extra feature x1*x2.
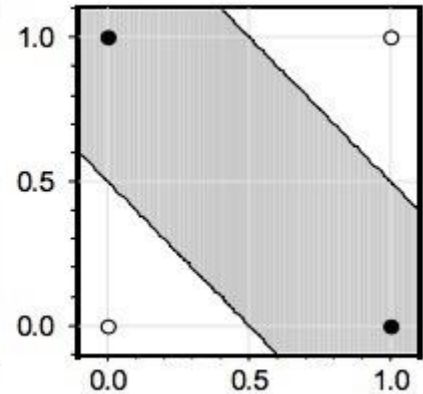
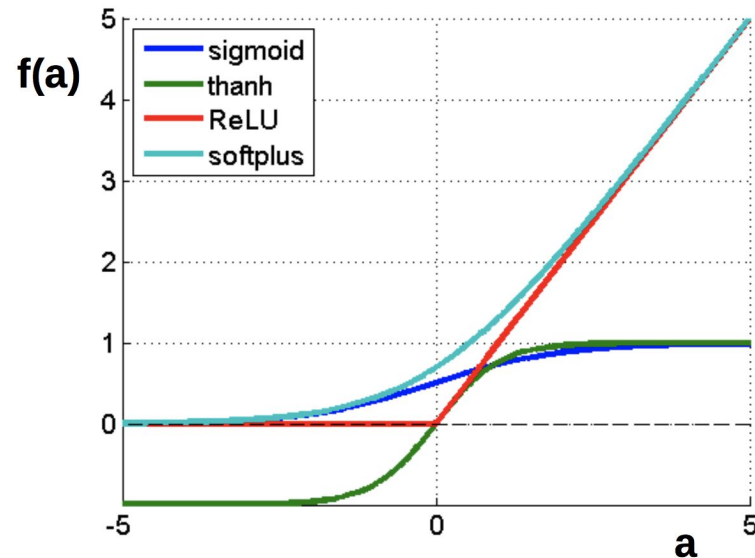AND       OR       XOR(with x1*x2)       XOR

# Activation functions: nonlinearities

$$f(a) = \frac{1}{1 + e^a}$$
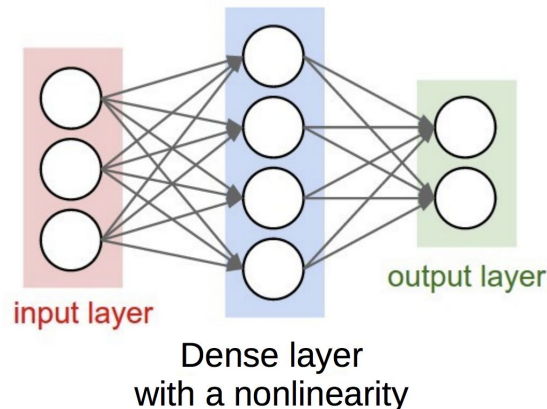
$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$
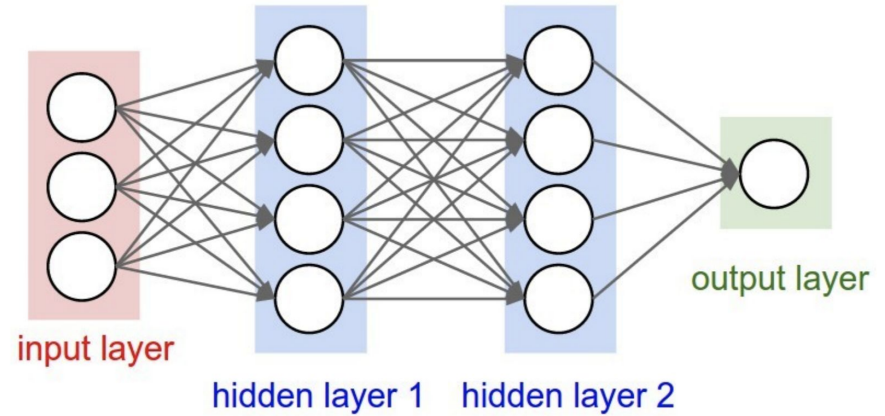
$$f(a) = \log(1 + e^a)$$

# Some generally accepted terms

- Layer – a building block for NNs :
  - Dense layer: $f(x) = Wx+b$
  - Nonlinearity layer: $f(x) = \sigma(x)$
  - Input layer, output layer
  - A few more we will cover later
- Activation function – function applied to layer output
  - Sigmoid
  - tanh
  - ReLU
  - Any other function to get nonlinear intermediate signal in NN
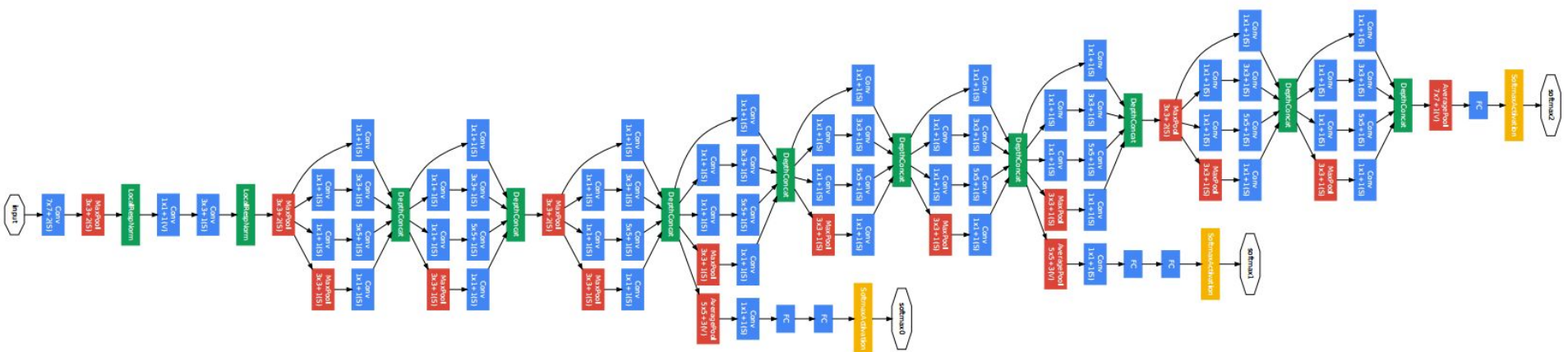- Backpropagation – a fancy word for "chain rule"



Dense layer
with a nonlinearity

"Train it via backprop!"

# Actually, it can be deeper



input layer

hidden layer 1    hidden layer 2

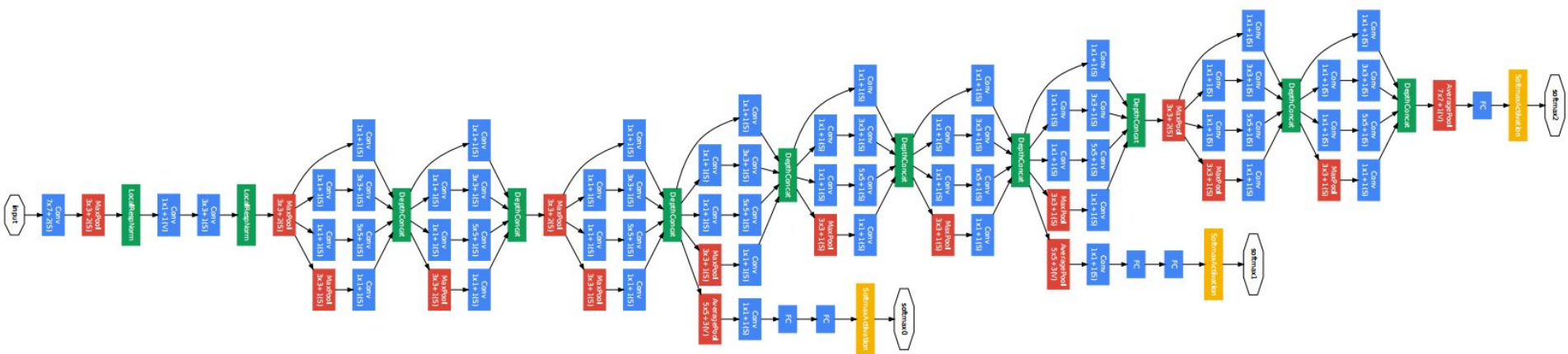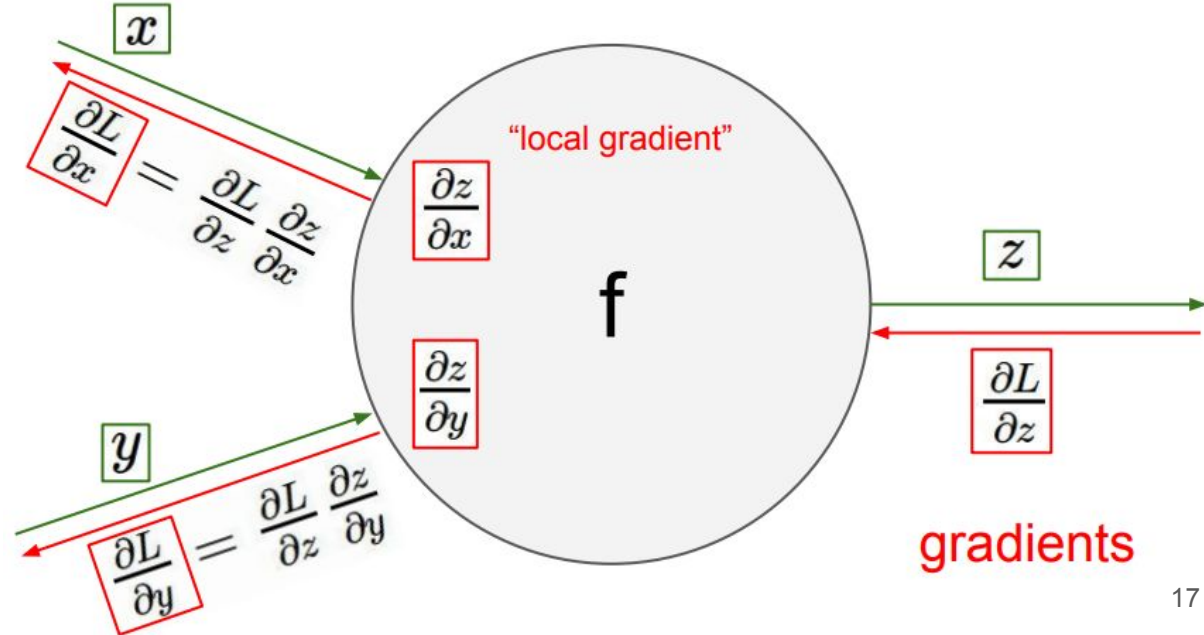output layer

Much deeper...

# Much deeper...



How to train it?

# Backpropagation and chain rule

Chain rule is just simple math:

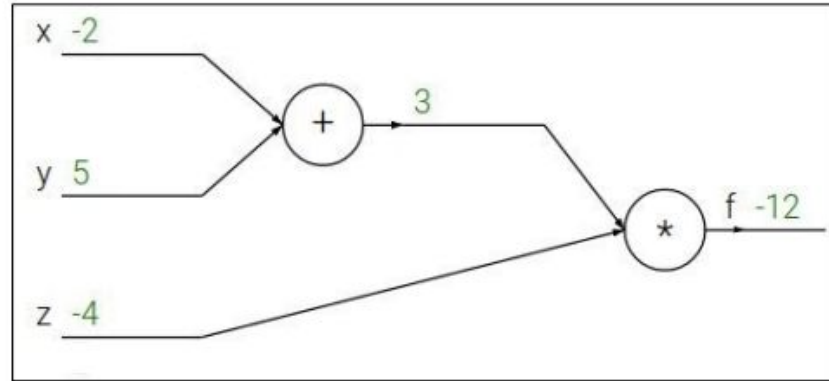$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$

Backprop is just way to use it in NN training.



source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

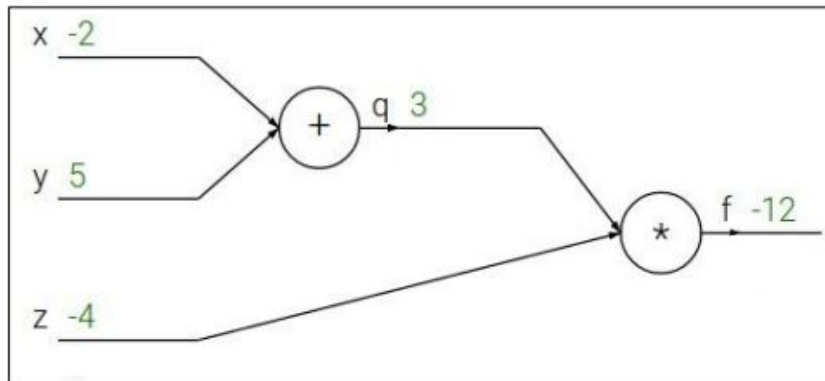source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$

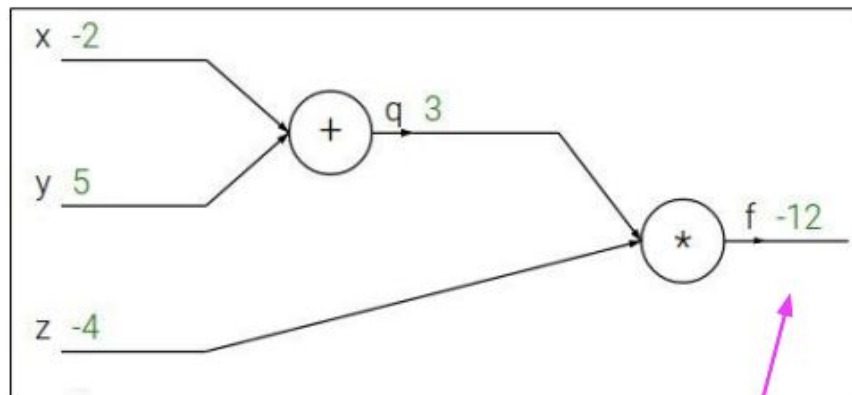source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial f}$$

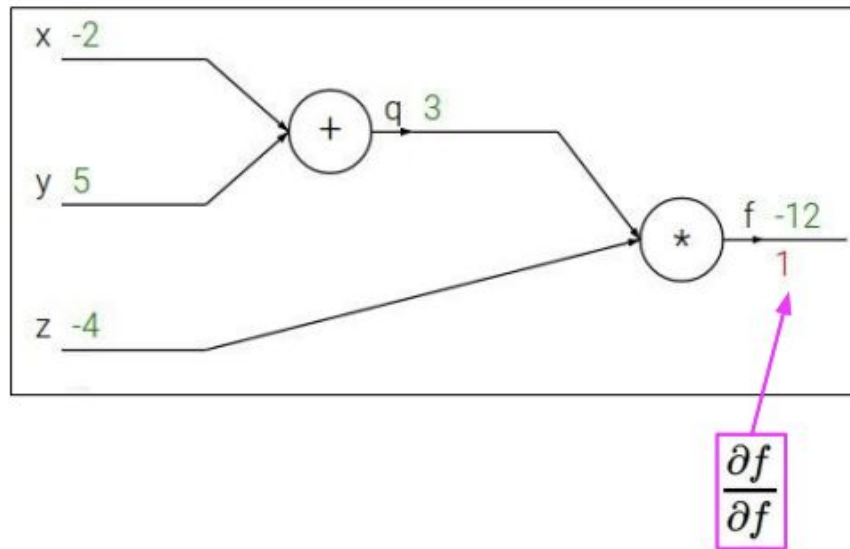source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



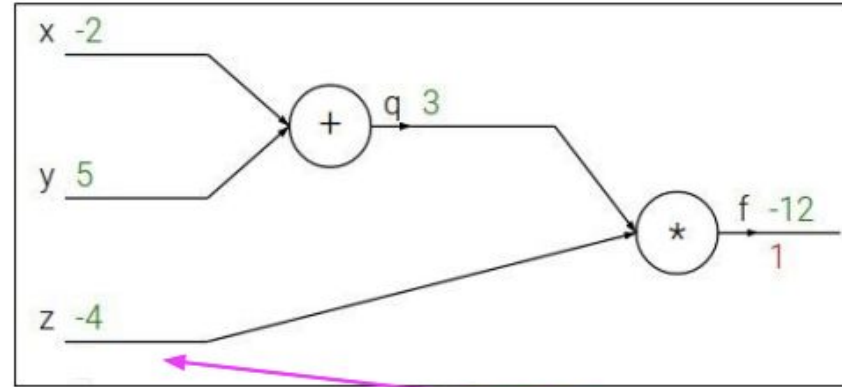$$\frac{\partial f}{\partial f}$$

21

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

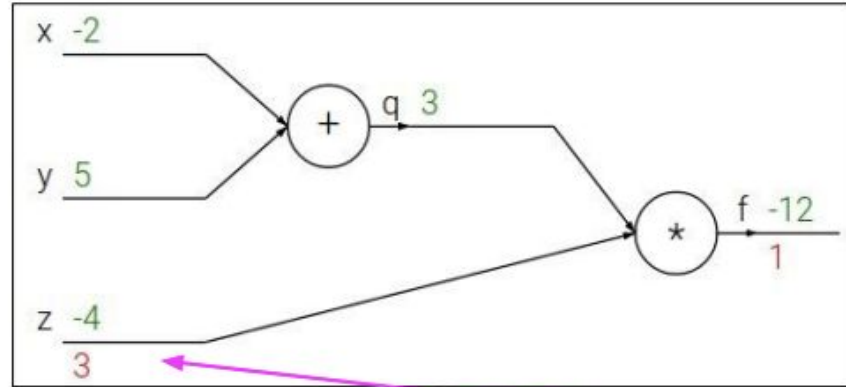source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial z}$$

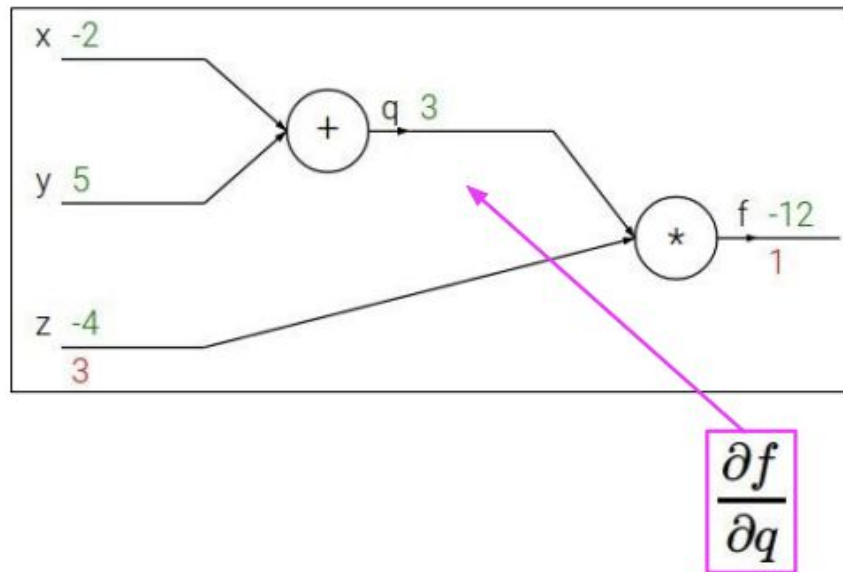source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$

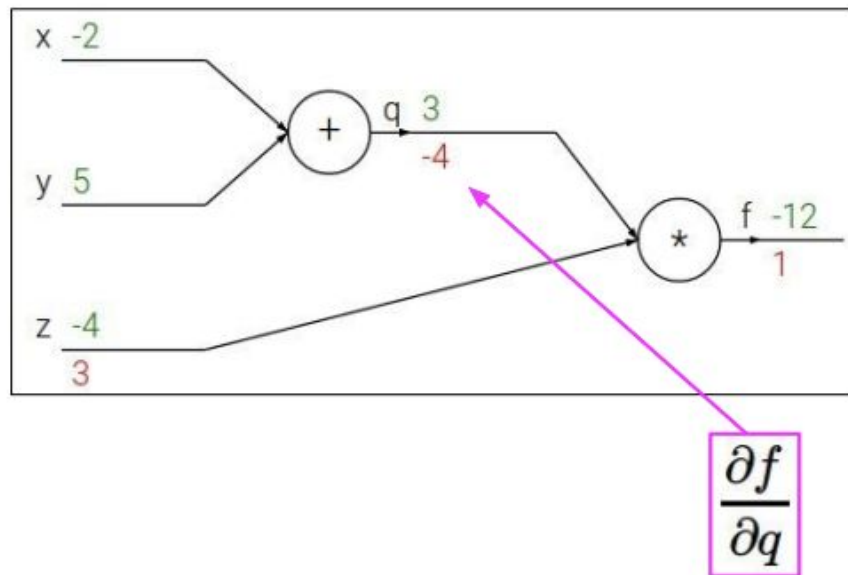source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial q}$$
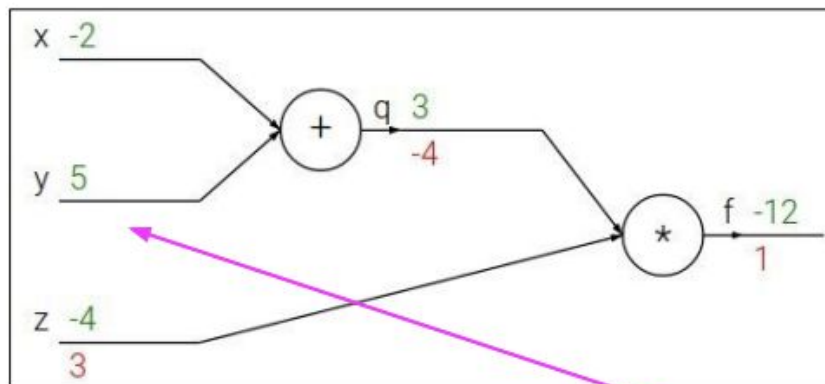
source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial y}$$

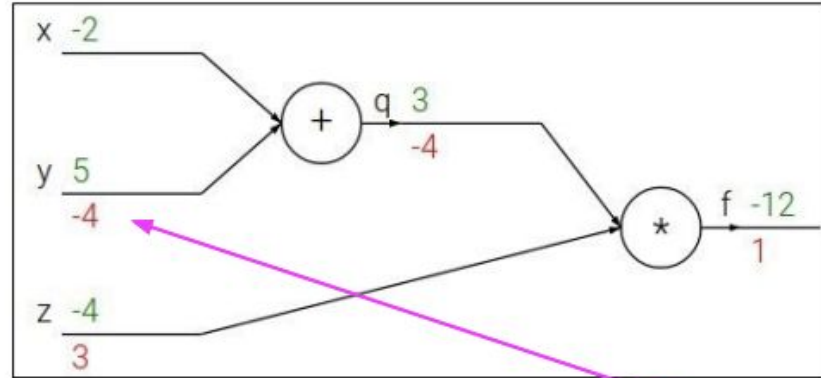source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$q = x + y \qquad \dfrac{\partial q}{\partial x} = 1, \dfrac{\partial q}{\partial y} = 1$

$f = qz \qquad \dfrac{\partial f}{\partial q} = z, \dfrac{\partial f}{\partial z} = q$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\dfrac{\partial f}{\partial y}$$

Chain rule:

$$\dfrac{\partial f}{\partial y} = \dfrac{\partial f}{\partial q} \dfrac{\partial q}{\partial y}$$

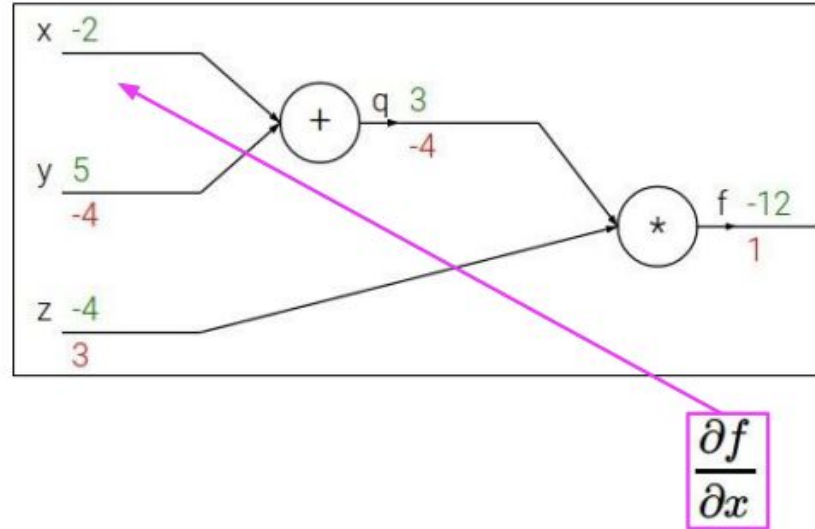source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\dfrac{\partial f}{\partial x}, \dfrac{\partial f}{\partial y}, \dfrac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

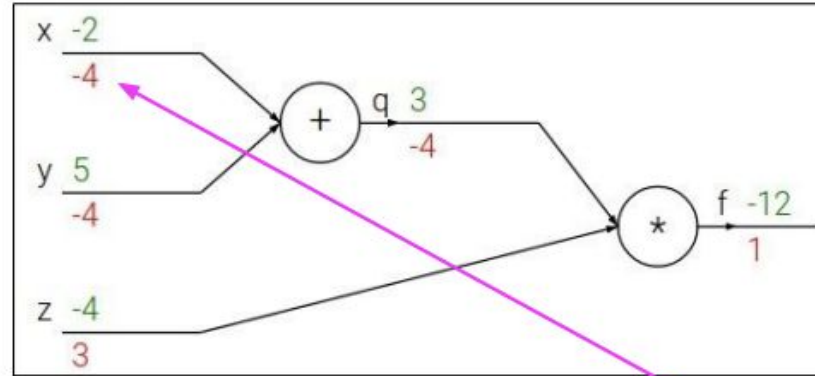source: http://cs231n.github.io

# Backpropagation example

$$f(x, y, z) = (x + y)z$$

e.g. x = -2, y = 5, z = -4

$$q = x + y \qquad \frac{\partial q}{\partial x} = 1, \frac{\partial q}{\partial y} = 1$$

$$f = qz \qquad \frac{\partial f}{\partial q} = z, \frac{\partial f}{\partial z} = q$$

Want: $\frac{\partial f}{\partial x}, \frac{\partial f}{\partial y}, \frac{\partial f}{\partial z}$



$$\frac{\partial f}{\partial x}$$

Chain rule:

$$\frac{\partial f}{\partial x} = \frac{\partial f}{\partial q} \frac{\partial q}{\partial x}$$

source: http://cs231n.github.io

# Practice time: interactive playground
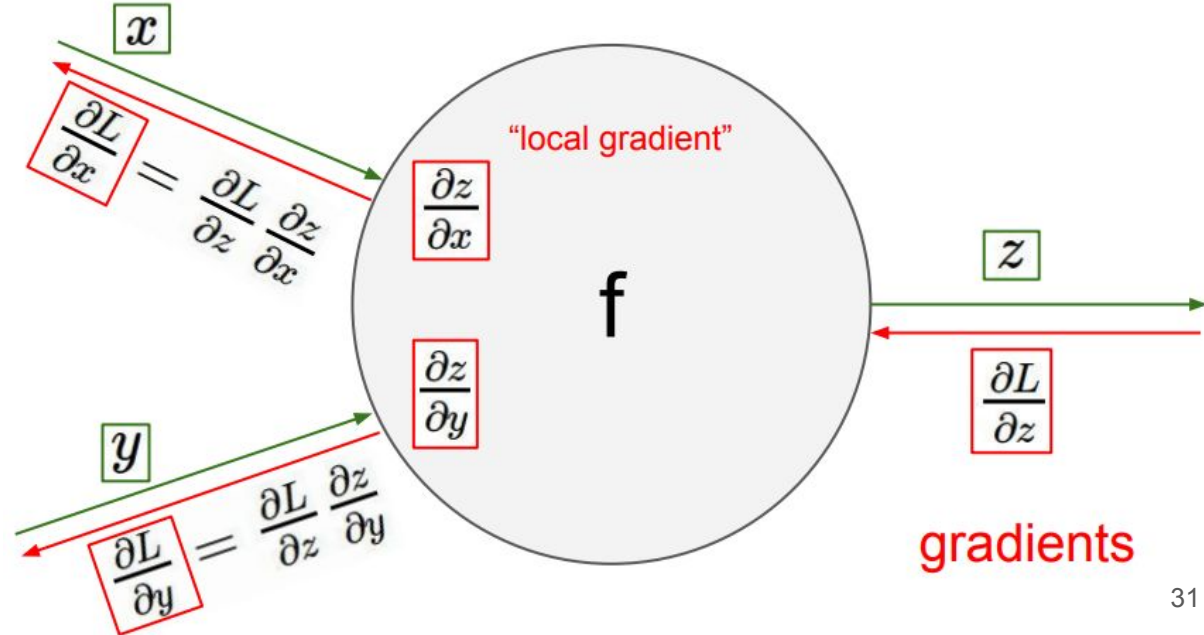
https://playground.tensorflow.org/

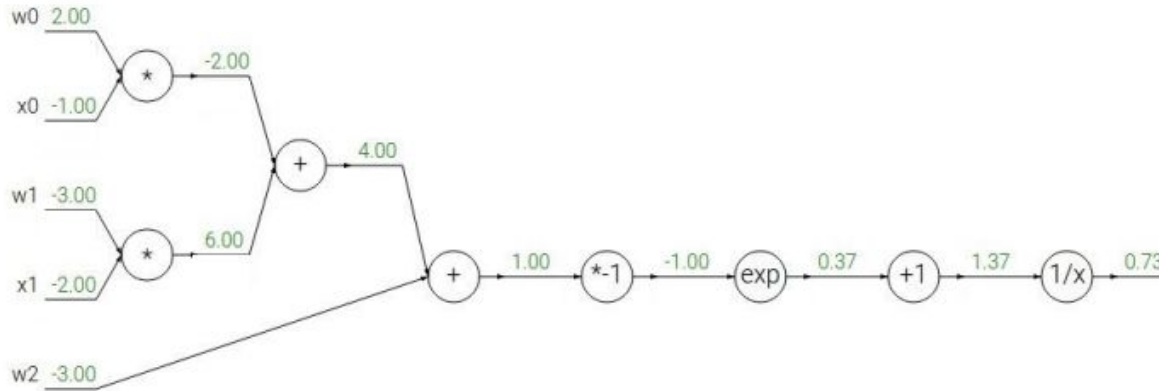# Backpropagation and chain rule

Chain rule is just simple math:

$$\frac{\partial L}{\partial x} = \frac{\partial L}{\partial z}\frac{\partial z}{\partial x}$$
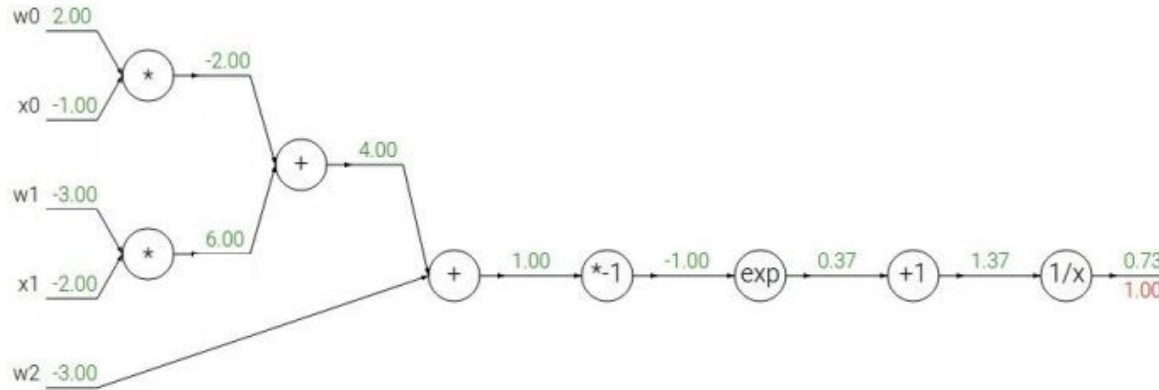
Backprop is just way to use it in NN training.



source: http://cs231n.github.io

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

source: http://cs231n.github.io

# Backpropagation example

Another example:  $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$

w0  2.00

x0  -1.00

$*$  -2.00

w1  -3.00

x1  -2.00

$*$  6.00

$+$  4.00

$+$  1.00  →  $*$-1  -1.00  →  exp  0.37  →  +1  1.37  →  1/x  0.73 / 1.00

w2  -3.00

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

source: http://cs231n.github.io

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

34

source: http://cs231n.github.io

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$\left(\dfrac{-1}{1.37^2}\right)(1.00) = -0.53$

| | | | |
|---|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | | $\dfrac{df}{dx} = a$ |

| | | | |
|---|---|---|---|
| $f(x) = \dfrac{1}{x}$ | $\rightarrow$ | | $\dfrac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | | $\dfrac{df}{dx} = 1$ |

source: http://cs231n.github.io

Another example: $f(w, x) = \dfrac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

source: http://cs231n.github.io

# Backpropagation example

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$(1)(-0.53) = -0.53$$

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Big| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



| | | | |
|---|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ | |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ | |

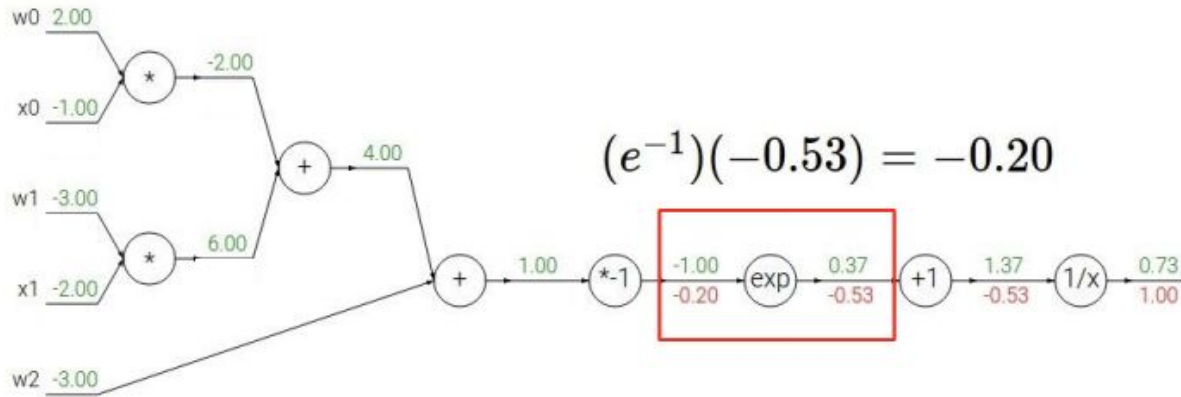| | | |
|---|---|---|
| $f(x) = \dfrac{1}{x}$ | $\rightarrow$ | $\dfrac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\dfrac{df}{dx} = 1$ |

source: http://cs231n.github.io

Another example:

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$(e^{-1})(-0.53) = -0.20$$

| | | |
|---|---|---|
| $f(x) = e^x$ | $\rightarrow$ | $\dfrac{df}{dx} = e^x$ |
| $f_a(x) = ax$ | $\rightarrow$ | $\dfrac{df}{dx} = a$ |

| | | |
|---|---|---|
| $f(x) = \dfrac{1}{x}$ | $\rightarrow$ | $\dfrac{df}{dx} = -1/x^2$ |
| $f_c(x) = c + x$ | $\rightarrow$ | $\dfrac{df}{dx} = 1$ |

source: http://cs231n.github.io

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \bigg| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \bigg| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

40

source: http://cs231n.github.io

# Backpropagation example

Another example:
$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



(-1) * (-0.20) = 0.20

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

source: http://cs231n.github.io

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



w0  2.00
x0  -1.00
-2.00

w1  -3.00
x1  -2.00
6.00

4.00

w2  -3.00

+ 1.00 / 0.20

*-1  -1.00 / -0.20

exp  0.37 / -0.53

+1  1.37 / -0.53

1/x  0.73 / 1.00

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

42

source: http://cs231n.github.io

Another example:
$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$

[local gradient] x [its gradient]
[1] x [0.2] = 0.2
[1] x [0.2] = 0.2  (both inputs!)

w0 2.00
x0 -1.00
-2.00
w1 -3.00
x1 -2.00
6.00
4.00
0.20
w2 -3.00
0.20
1.00
0.20
*-1
-1.00
-0.20
exp
0.37
-0.53
+1
1.37
-0.53
1/x
0.73
1.00

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

43

source: http://cs231n.github.io

# Backpropagation example

Another example:
$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x \qquad \Big| \qquad f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a \qquad \Big| \qquad f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

44

Another example:

$$f(w,x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$



[local gradient] x [its gradient]
x0: [2] x [0.2] = 0.4
w0: [-1] x [0.2] = -0.2

$$f(x) = e^x \qquad \rightarrow \qquad \frac{df}{dx} = e^x$$

$$f_a(x) = ax \qquad \rightarrow \qquad \frac{df}{dx} = a$$

$$f(x) = \frac{1}{x} \qquad \rightarrow \qquad \frac{df}{dx} = -1/x^2$$

$$f_c(x) = c + x \qquad \rightarrow \qquad \frac{df}{dx} = 1$$

source: http://cs231n.github.io

# Backpropagation example

$$f(w, x) = \frac{1}{1 + e^{-(w_0 x_0 + w_1 x_1 + w_2)}}$$
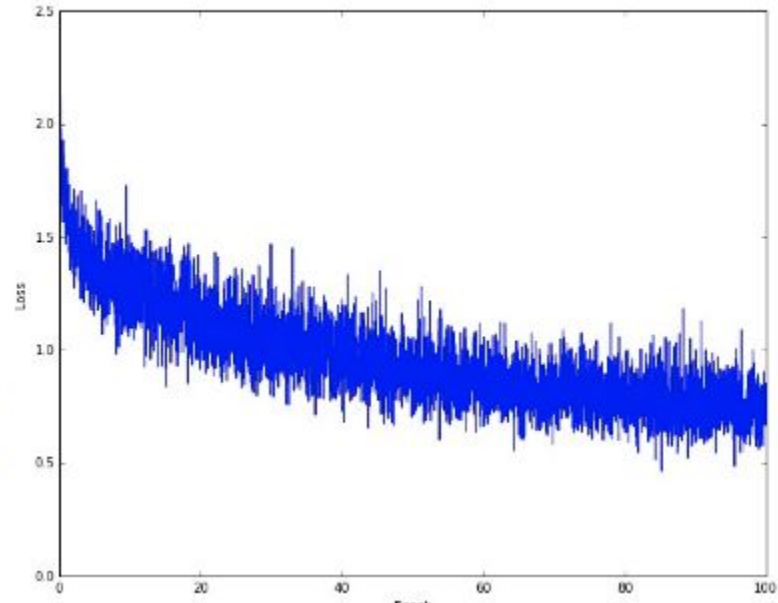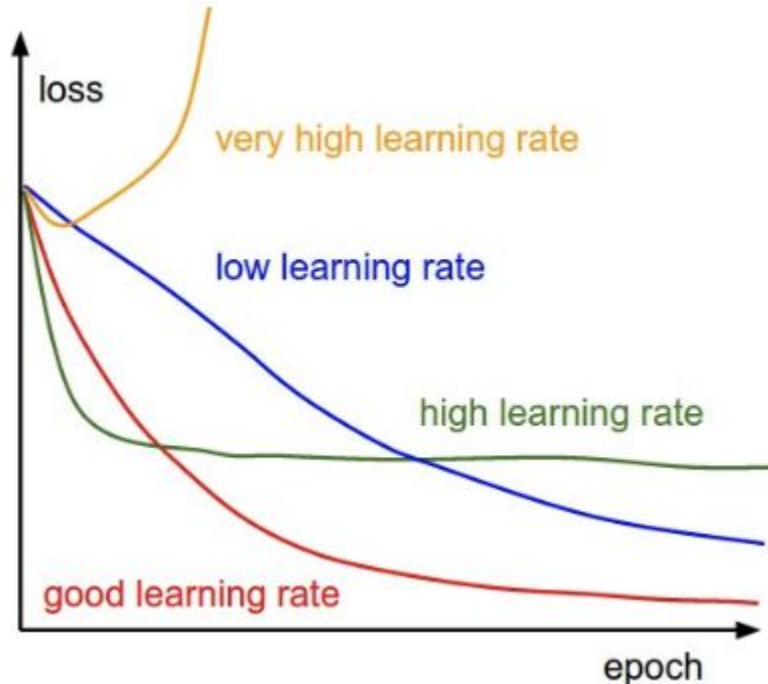
$$\sigma(x) = \frac{1}{1 + e^{-x}}$$ sigmoid function

$$\frac{d\sigma(x)}{dx} = \frac{e^{-x}}{(1 + e^{-x})^2} = \left(\frac{1 + e^{-x} - 1}{1 + e^{-x}}\right)\left(\frac{1}{1 + e^{-x}}\right) = (1 - \sigma(x))\,\sigma(x)$$



sigmoid gate

(0.73) * (1 - 0.73) = 0.2

source: http://cs231n.github.io

Stochastic gradient descent (and variations)
is used to optimize NN parameters.

$$x_{t+1} = x_t - \text{learning rate} \cdot dx$$



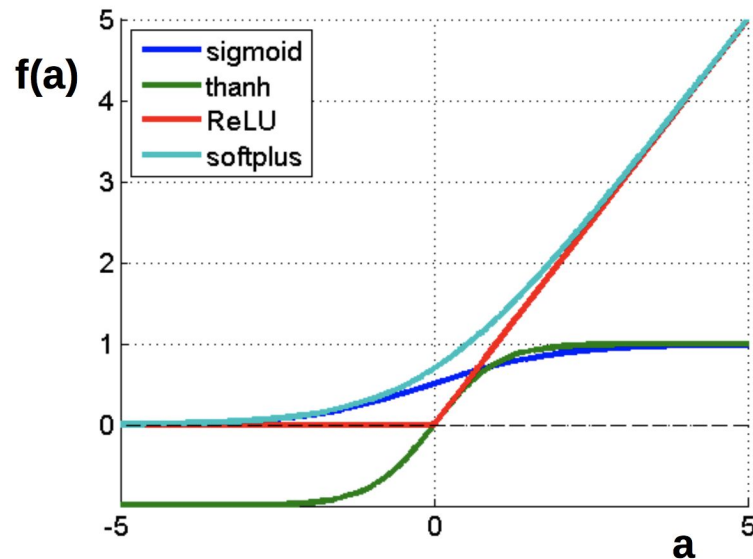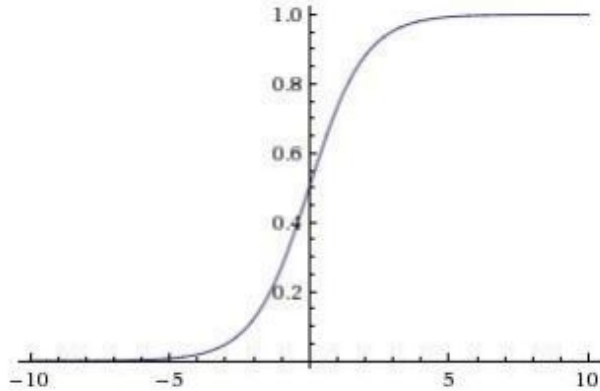source: http://cs231n.github.io/neural-networks-3/

$$f(a) = \frac{1}{1 + e^a}$$

$$f(a) = \tanh(a)$$

$$f(a) = \max(0, a)$$

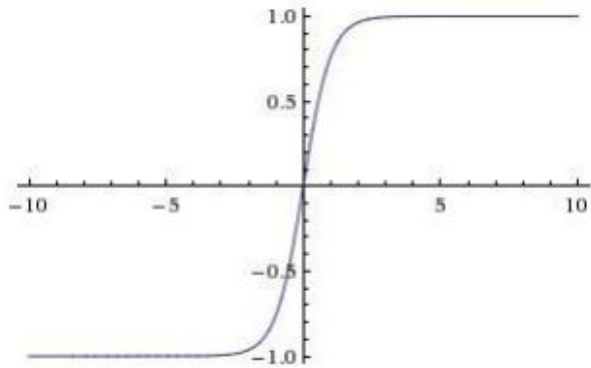$$f(a) = \log(1 + e^a)$$

**Sigmoid**

$$f(a) = \frac{1}{1 + e^a}$$

- 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

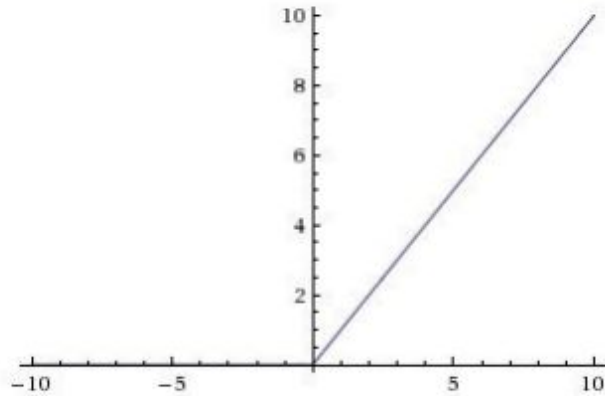**tanh(x)**

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

$$f(a) = \tanh(a)$$

- Does not saturate (in +region)
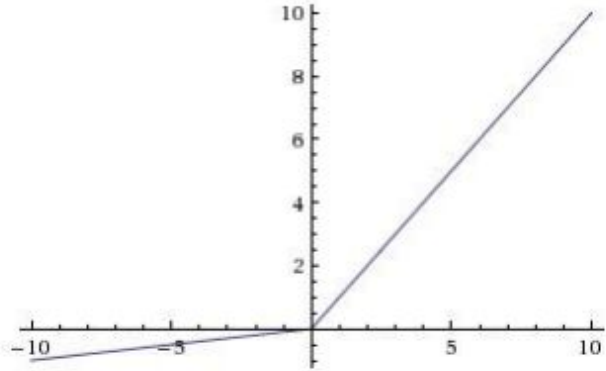- Very computationally efficient
- Converges much faster than sigmoid/tanh in practice (e.g. 6x)

- Not zero-centered output
- An annoyance:

hint: what is the gradient when x < 0?

**ReLU**
(Rectified Linear Unit)

$$f(a) = \max(0, a)$$

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

**Leaky ReLU**

$$f(x) = \max(0.01x, x)$$
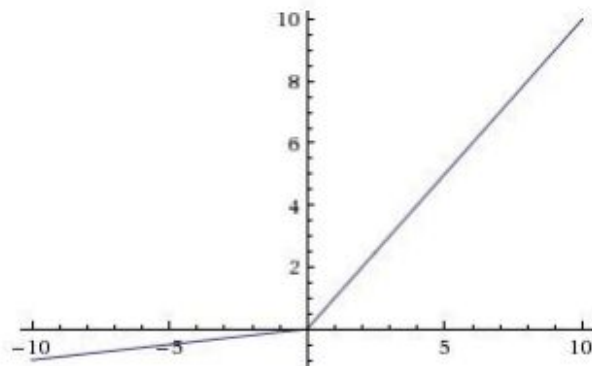
- Does not saturate
- Computationally efficient
- Converges much faster than sigmoid/tanh in practice! (e.g. 6x)
- **will not "die".**

**Leaky ReLU**
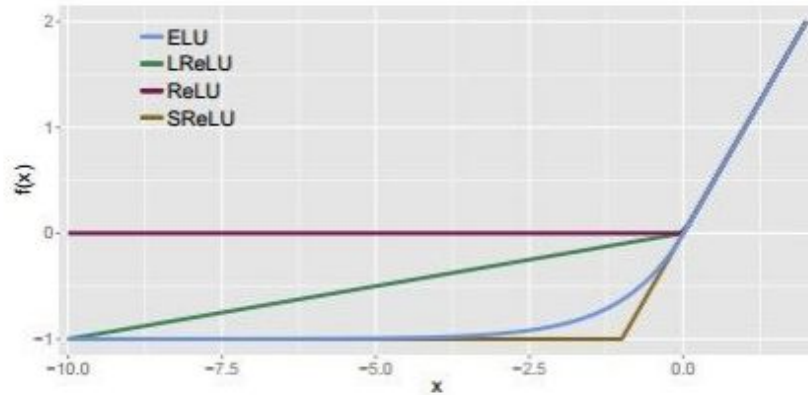
$$f(x) = \max(0.01x, x)$$

**Parametric Rectifier (PReLU)**

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

backprop into \alpha
(parameter)

## Exponential Linear Units (ELU)



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

- All benefits of ReLU
- Does not die
- Closer to zero mean outputs

- Computation requires exp()

- Use ReLU as baseline approach
- Be careful with the learning rates
- Try out Leaky ReLU or ELU
- Try out tanh but do not expect much from it
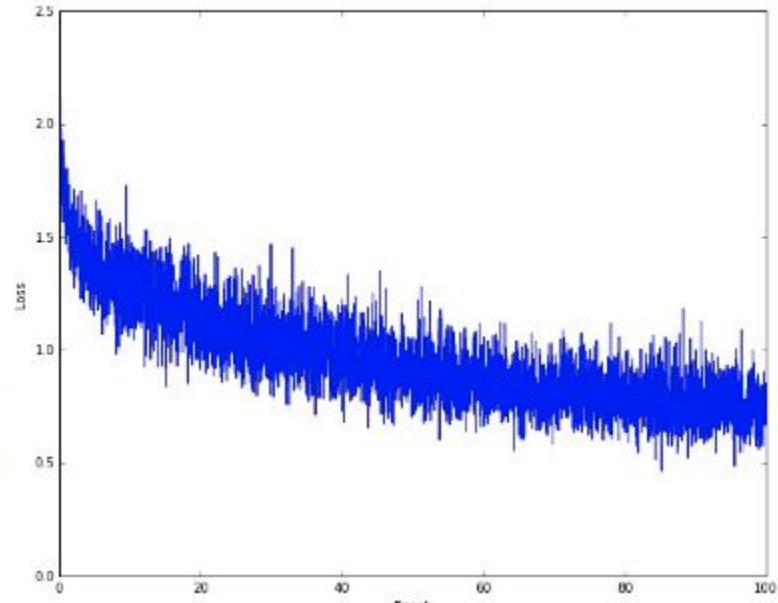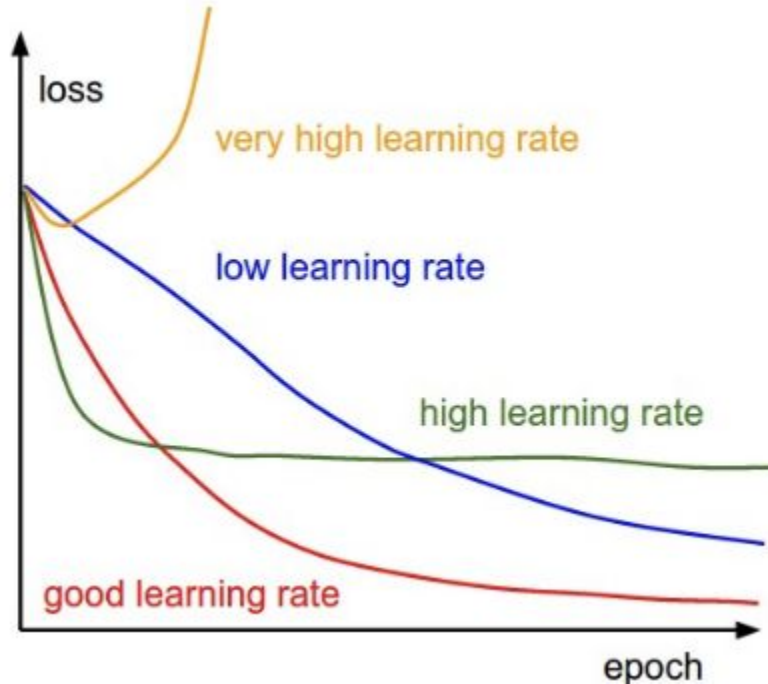- Do not use Sigmoid

# That's all. Time to build some NN.

# Backup

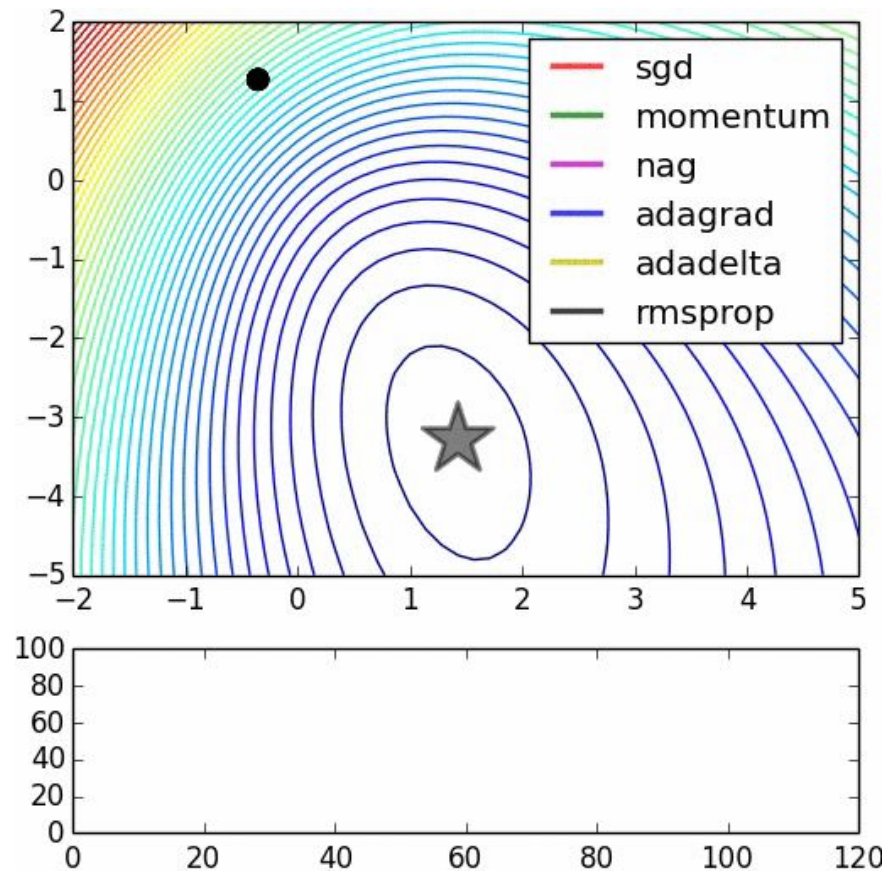Stochastic gradient descent is used to optimize NN parameters.

$$x_{t+1} = x_t - \text{learning rate} \cdot dx$$





source: http://cs231n.github.io/neural-networks-3/

There are much more optimizers:
- Momentum
- Adagrad
- Adadelta
- RMSprop
- Adam
- …
- even other NNs
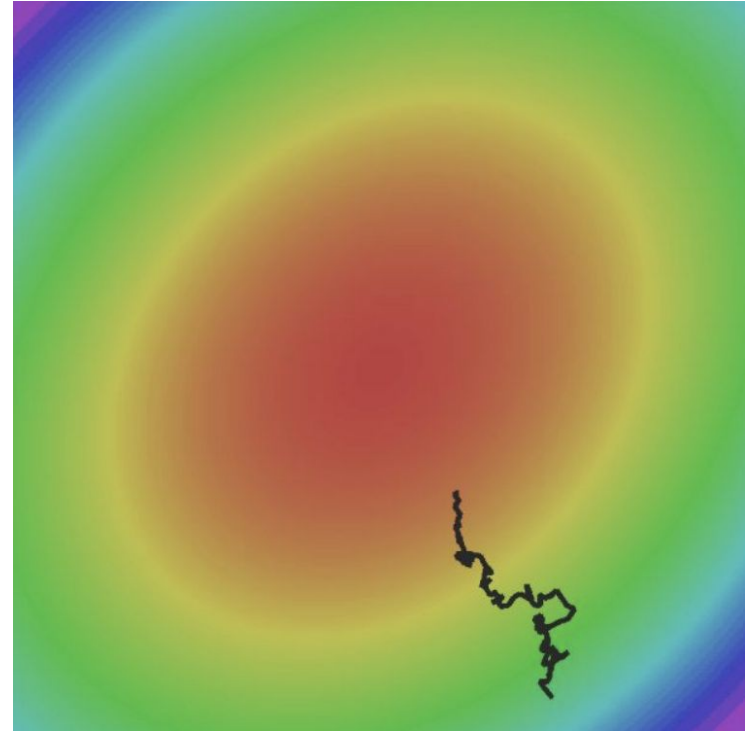
$$L(W) = \frac{1}{N} \sum_{i=1}^{N} L_i(x_i, y_i, W)$$

$$\nabla_W L(W) = \frac{1}{N} \sum_{i=1}^{N} \nabla_W L_i(x_i, y_i, W)$$

Averaging over minibatches ---> noisy gradient

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf
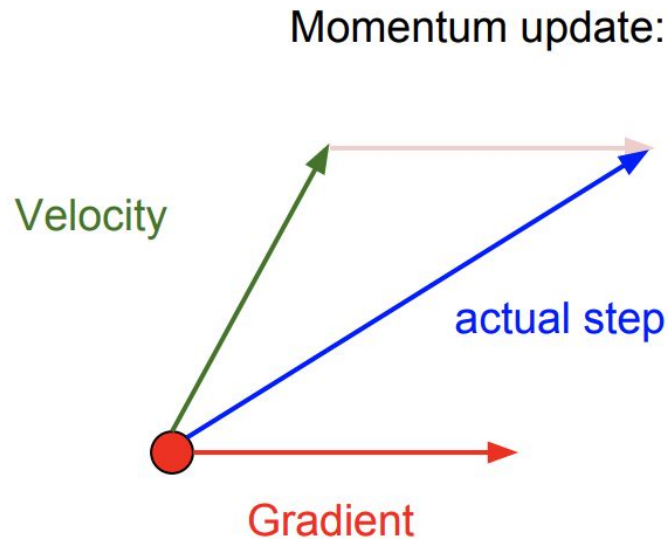
# First idea: momentum

Simple SGD

$$x_{t+1} = x_t - \alpha \nabla f(x_t)$$

SGD with momentum

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

Momentum update:



Velocity

actual step

Gradient

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

# Nesterov momentum

Momentum update:

Velocity

actual step

Gradient

$$v_{t+1} = \rho v_t + \nabla f(x_t)$$
$$x_{t+1} = x_t - \alpha v_{t+1}$$

Nesterov Momentum

Gradient

Velocity

actual step

$$v_{t+1} = \rho v_t - \alpha \nabla f(\boxed{x_t + \rho v_t})$$
$$x_{t+1} = x_t + v_{t+1}$$

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

# Comparing momentums



SGD

SGD+Momentum

Nesterov

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

# Second idea: different dimensions are different

Adagrad: SGD with cache

$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$
$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

# Second idea: different dimensions are different

Adagrad: SGD with cache

$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$
$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

*Problem: gradient fades with time*

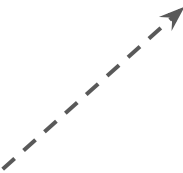# Second idea: different dimensions are different

Adagrad: SGD with cache

$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$

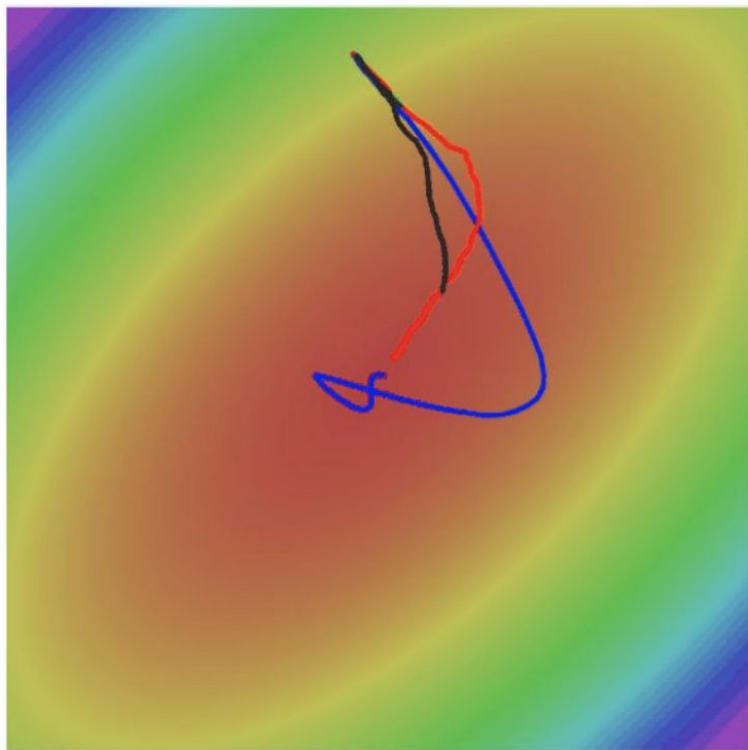$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

RMSProp: SGD with cache with exp. Smoothing

$$\text{cache}_{t+1} = \beta \text{cache}_t + (1 - \beta)(\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha \frac{\nabla f(x_t)}{\text{cache}_{t+1} + \varepsilon}$$

Slide 29 Lecture 6 of Geoff Hinton's Coursera class
http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf

**SGD**

**SGD+Momentum**

**RMSProp**

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

Let's combine the momentum idea and RMSProp normalization:

$$v_{t+1} = \gamma v_t + (1 - \gamma)\nabla f(x_t)$$

$$\text{cache}_{t+1} = \beta\text{cache}_t + (1 - \beta)(\nabla f(x_t))^2$$

$$x_{t+1} = x_t - \alpha\frac{v_{t+1}}{\text{cache}_{t+1} + \varepsilon}$$

Adam full form involves bias correction term. See http://cs231n.github.io/neural-networks-3/ for more info.

Let's combine the momentum idea and RMSProp normalization:

$$v_{t+1} = \gamma v_t + (1 - \gamma)\nabla f(x_t)$$

$$\text{cache}_{t+1} = \beta \text{cache}_t + (1 - \beta)(\nabla f(x_t))^2$$

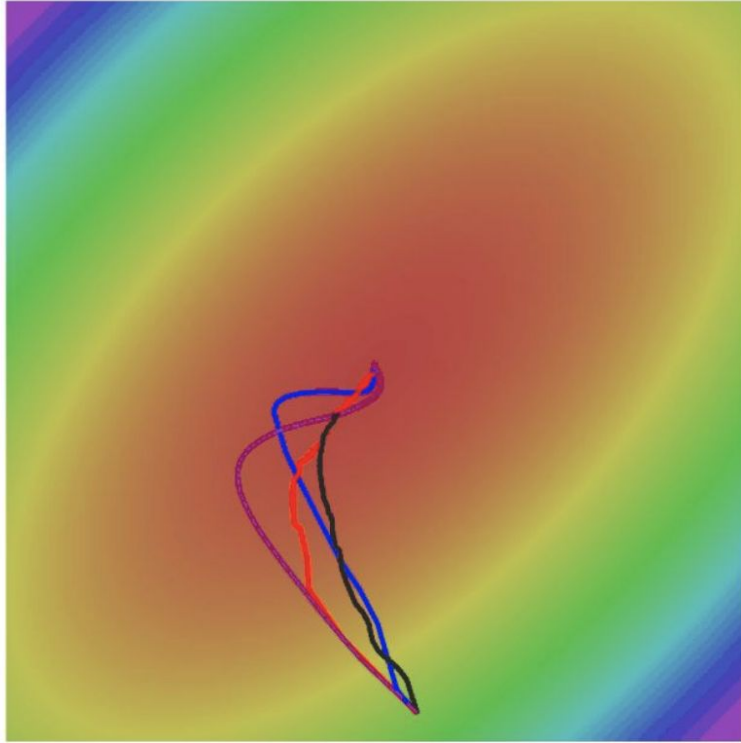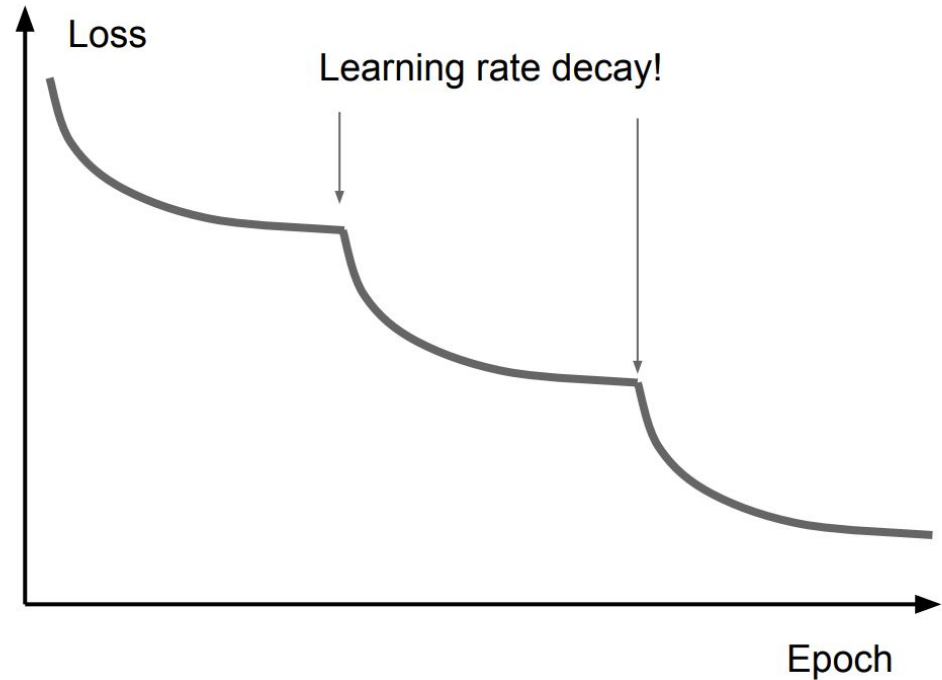$$x_{t+1} = x_t - \alpha \frac{v_{t+1}}{\text{cache}_{t+1} + \varepsilon}$$

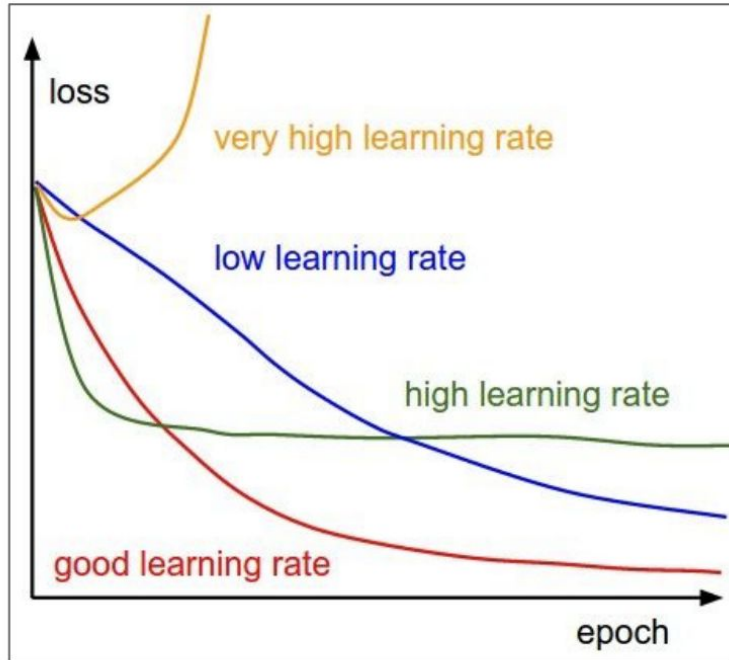*Actually, that's not quite Adam.*

Adam full form involves bias correction term. See http://cs231n.github.io/neural-networks-3/ for more info.

# Comparing optimizers



SGD

SGD+Momentum

RMSProp

Adam

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

# Once more: learning rate



loss

very high learning rate

low learning rate

high learning rate

good learning rate

epoch

Loss

Learning rate decay!

Epoch

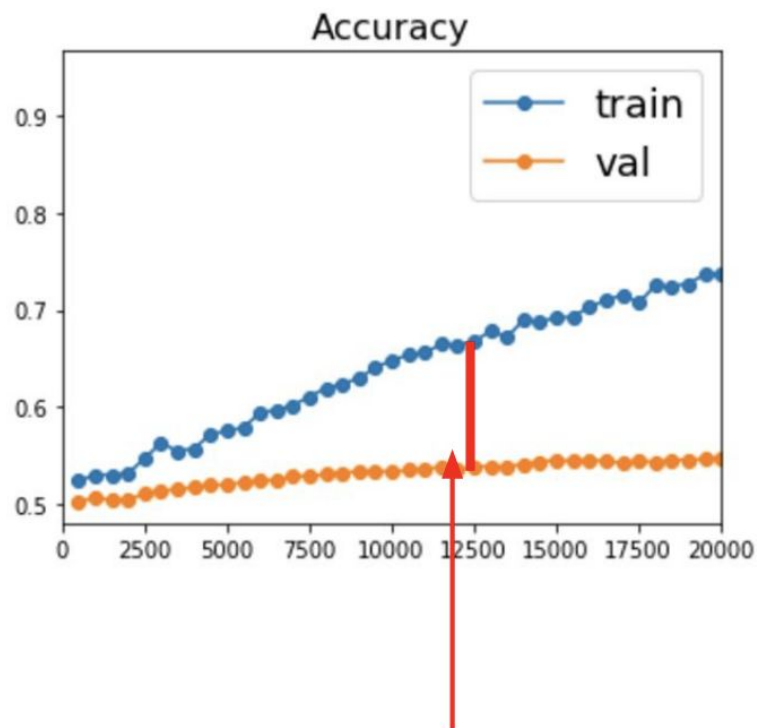source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf

# Sum up: optimization

- Adam is great basic choice
- Even for Adam/RMSProp learning rate matters
- Use learning rate decay
- Monitor your model quality

## Train Loss

## Accuracy

Better optimization algorithms
help reduce training loss

But we really care about error on new
data - how to reduce the gap?

source: http://cs231n.stanford.edu/slides/2017/cs231n_2017_lecture7.pdf