# Outline

1. Previous lecture recap: backpropagation, activations, intuition.
2. Optimizers.
3. Data normalization.
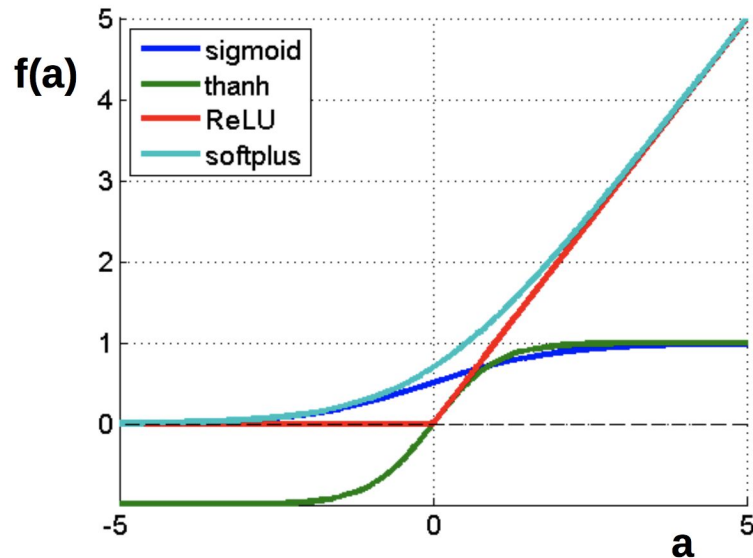4. Regularization.

# Once again: nonlinearities

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

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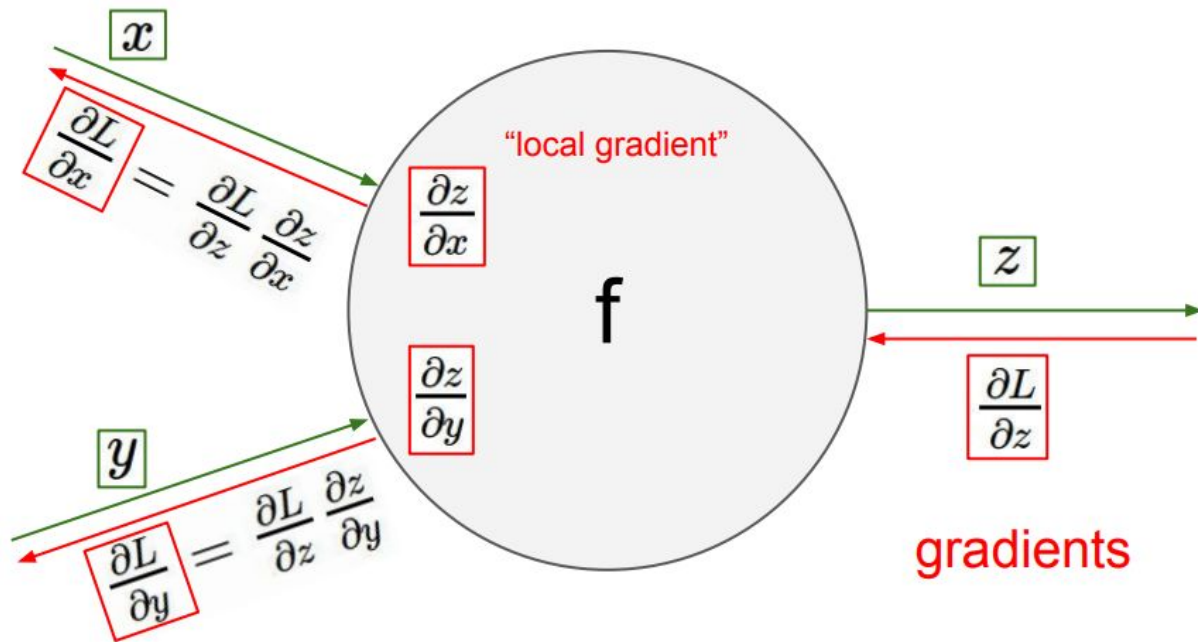
$$f(a) = \max(0, a)$$

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

# Backpropagation and chain rule

Chain rule is just simple math:

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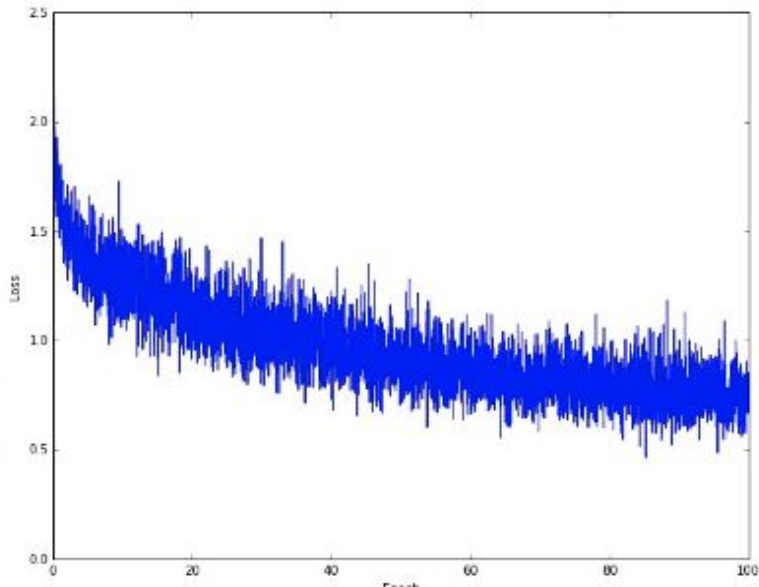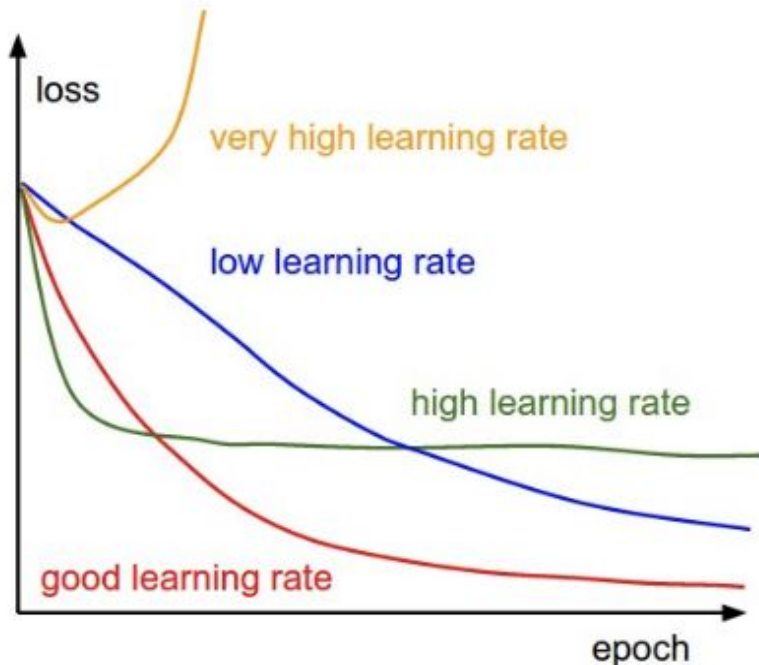
Backprop is just way to use it in NN training.

# Optimizers

Stochastic gradient descent is used to optimize NN parameters.
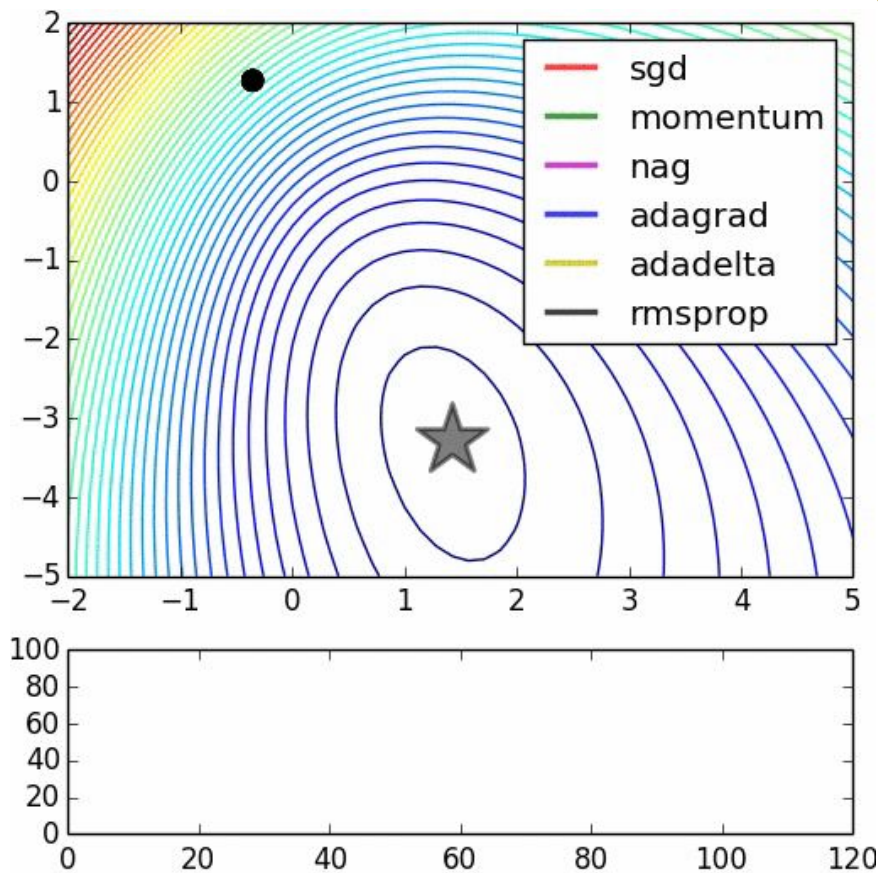
$$x_{t+1} = x_t - \text{learning rate} \cdot dx$$

# Optimizers

There are much more optimizers:

- Momentum
- Adagrad
- Adadelta
- RMSprop
- Adam
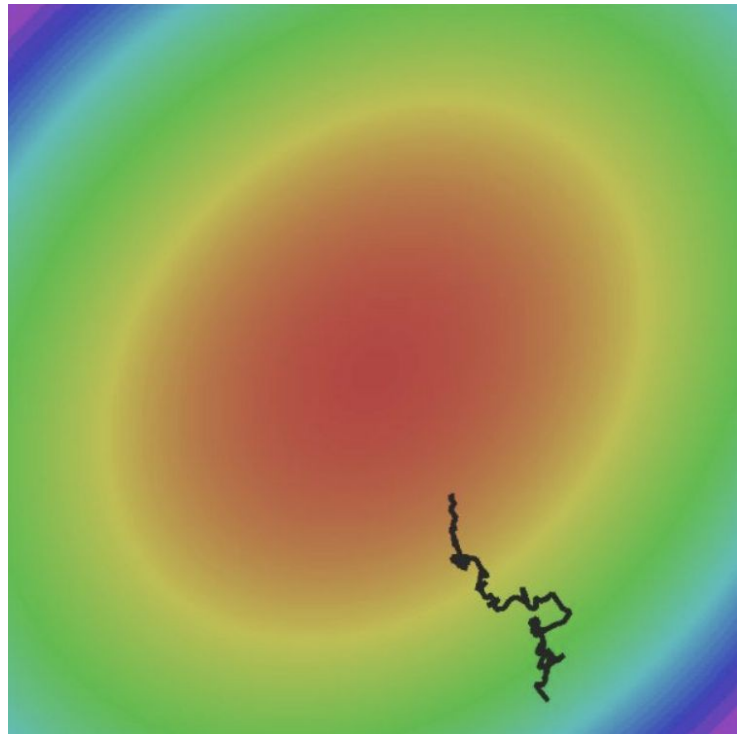- …
- even other NNs



source: link

# Optimization: SGD

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

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

Averaging over mini batches => noisy gradient

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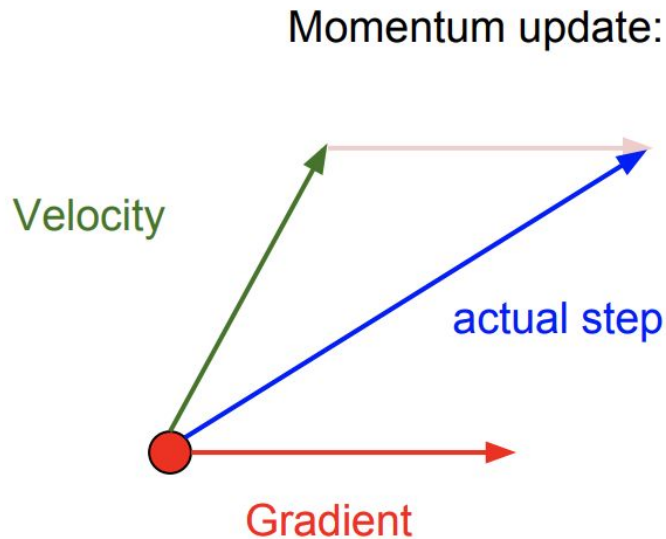
# First idea: momentum

Simple SGD

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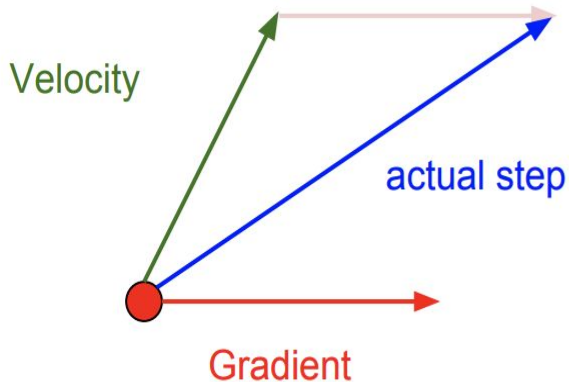
SGD with momentum

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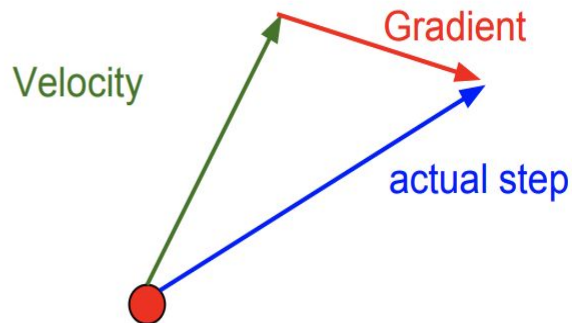$$x_{t+1} = x_t - \alpha v_{t+1}$$



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

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

9

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

# Comparing momentums

source:
https://ruder.io/content/images/2016/09/saddle_point_evaluation_optimizers.gif

# 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

13

# Second idea: different dimensions are different

Adagrad: SGD with cache

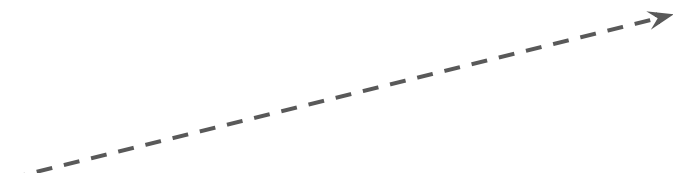$$\text{cache}_{t+1} = \text{cache}_t + (\nabla f(x_t))^2$$

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

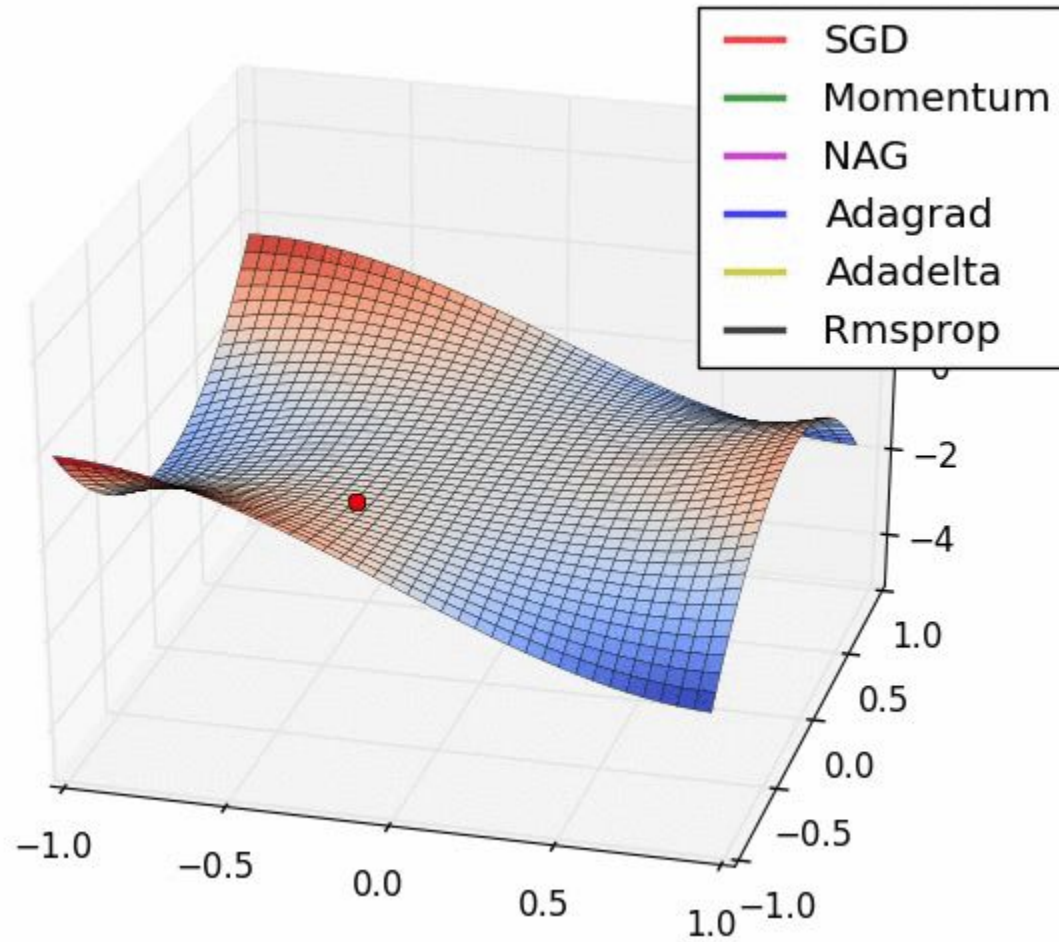RMSProp: SGD with cache with exp. Smoothing

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

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

14

source: https://imgur.com/a/Hqolp#NKsFHJb
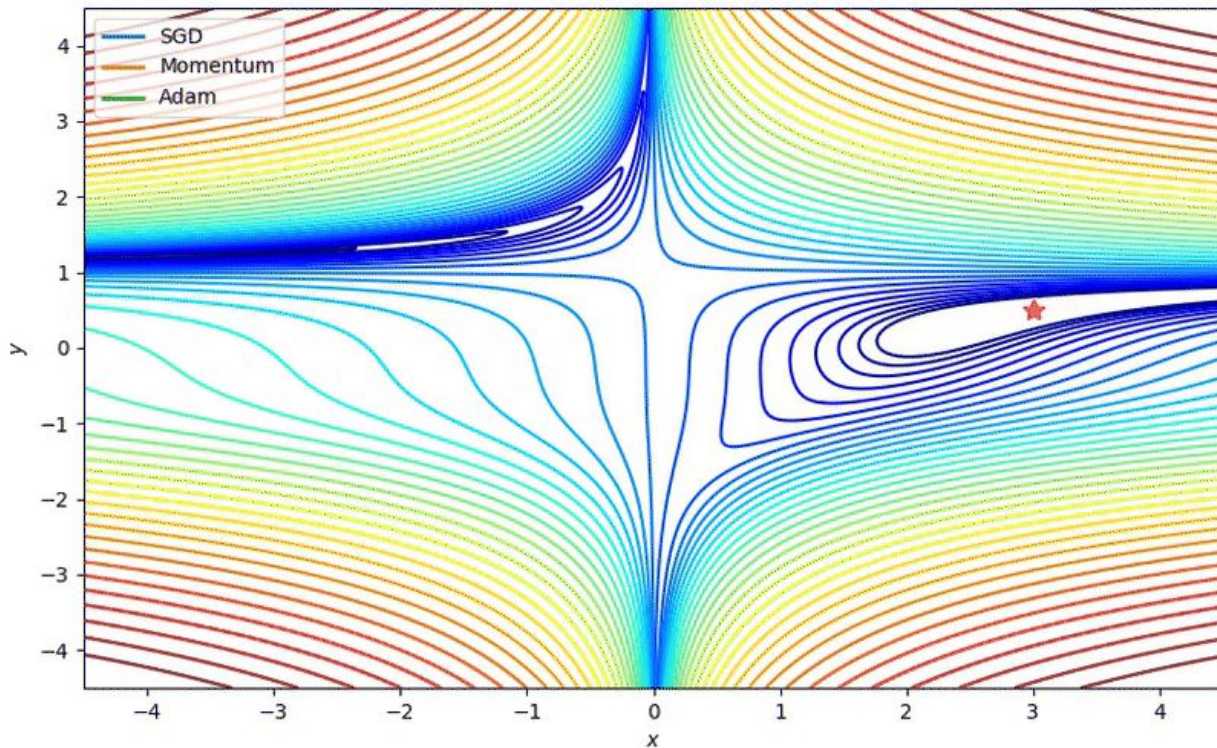
# Adam

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



source:
https://joshvarty.com/2018/02/27/ltfn-7-a-quick-look-at-tensorflow-optimizers/

Andrej Karpathy ✔
@karpathy

3e-4 is the best learning rate for Adam, hands down.

6:01 AM · Nov 24, 2016 · Twitter Web Client

108 Retweets    461 Likes

Andrej Karpathy ✔  @karpathy · Nov 24, 2016
Replying to @karpathy
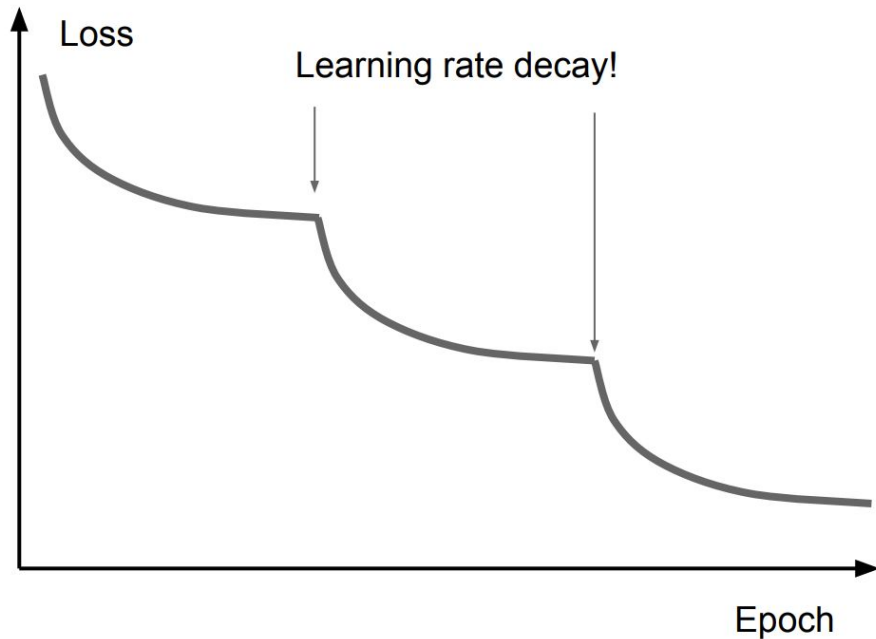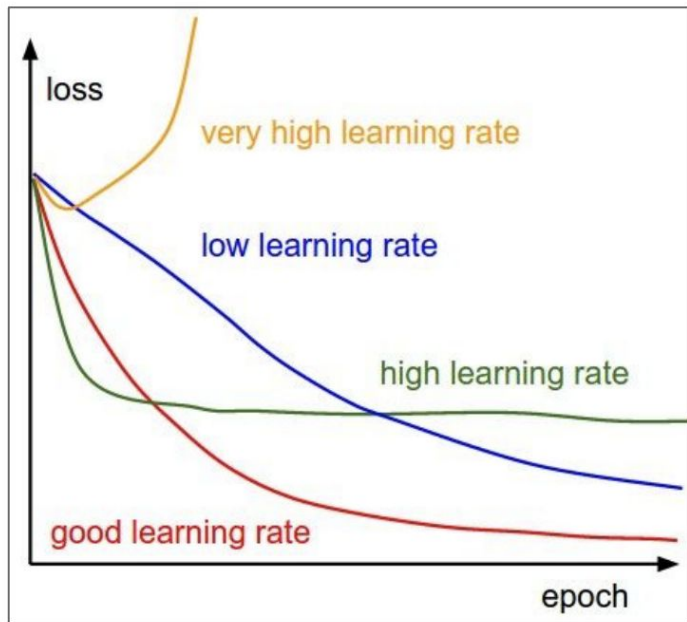(i just wanted to make sure that people understand that this is a joke...)

9    3    119

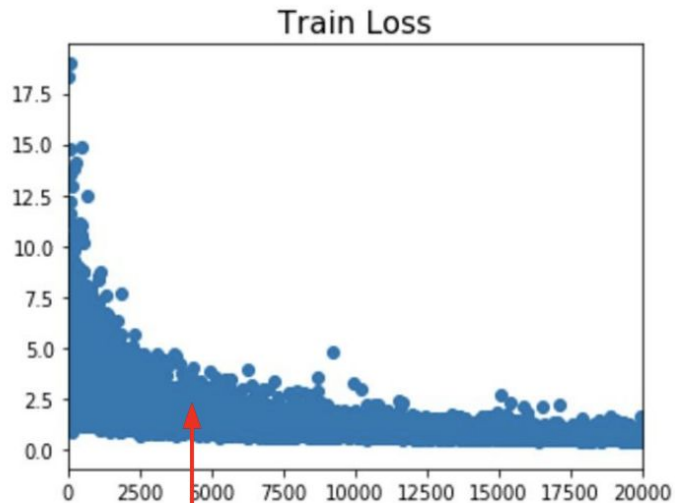source: https://twitter.com/karpathy/status/801621764144971776
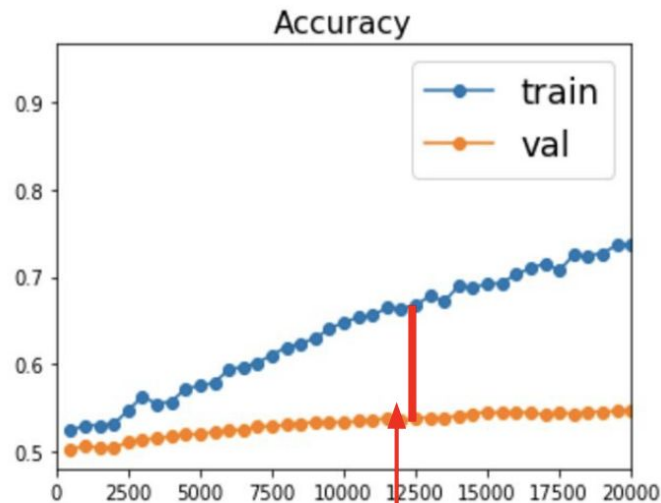
18

# Once more: learning rate

# Sum up: optimization

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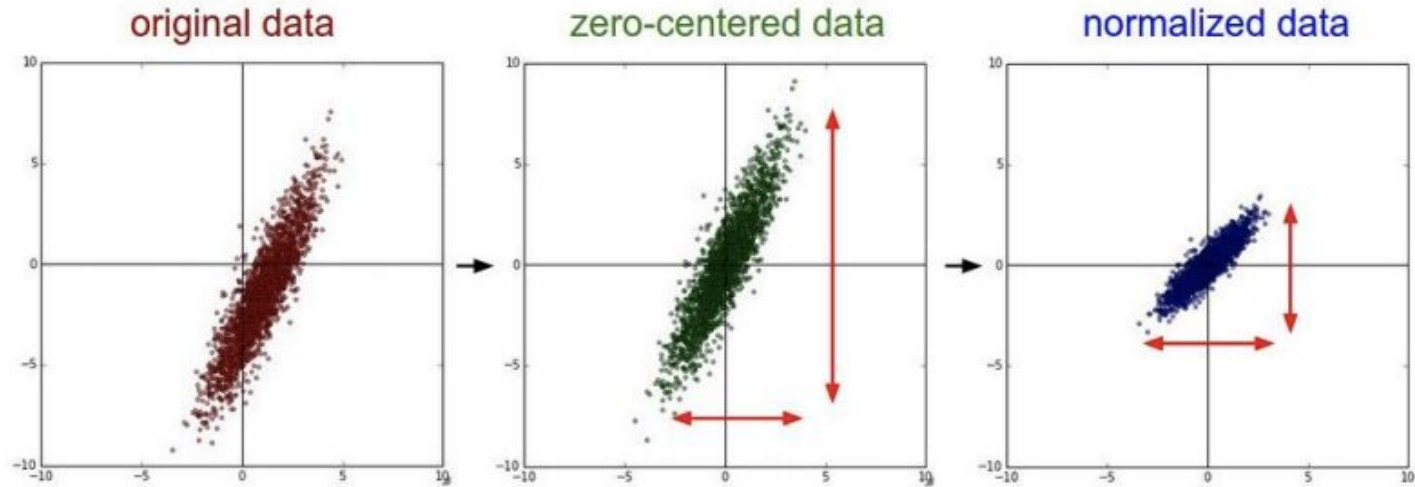- Monitor your model quality

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

# Data normalization



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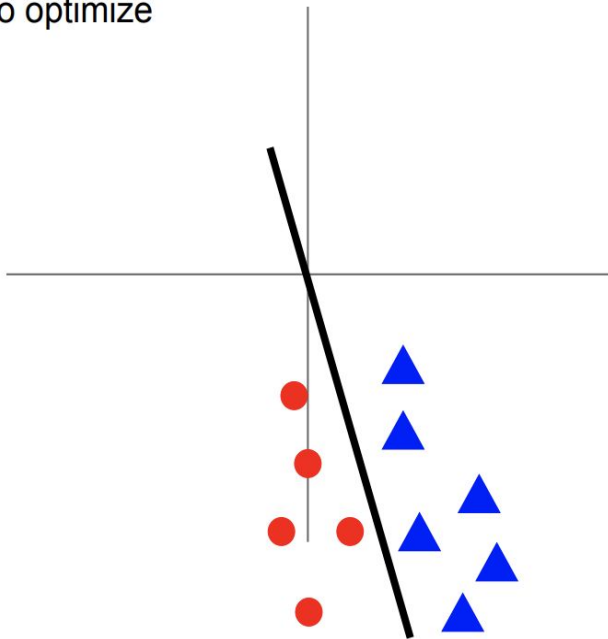
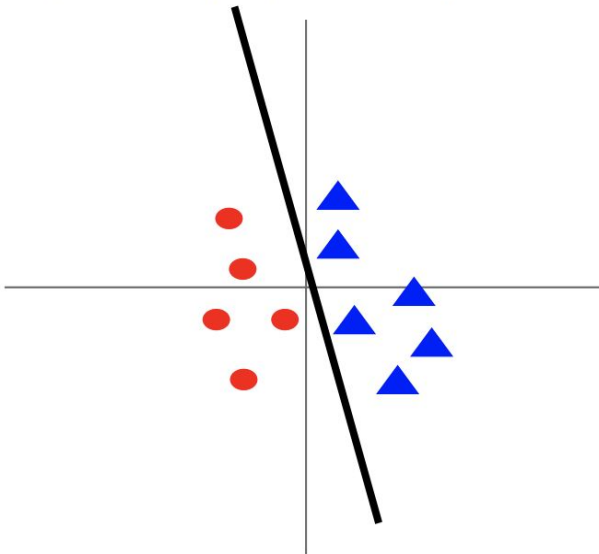# Data normalization

**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

# Weights initialization

- Pitfall: all zero initialization.

# Weights initialization

- Pitfall: all zero initialization.
- Small random numbers.

# Weights initialization

- Pitfall: all zero initialization.
- Small random numbers.
- Calibrated random numbers.

$$\text{Var}(s) = \text{Var}(\sum_i^n w_i x_i)$$

$$= \sum_i^n \text{Var}(w_i x_i)$$

$$= \sum_i^n [E(w_i)]^2 \text{Var}(x_i) + E[(x_i)]^2 \text{Var}(w_i) + \text{Var}(x_i)\text{Var}(w_i)$$

$$= \sum_i^n \text{Var}(x_i)\text{Var}(w_i)$$

$$= (n\text{Var}(w))\,\text{Var}(x)$$

# Thanks for attention!

Questions?

girafe
ai