# Lecture 3: Advanced Optimizations

Deep Learning @ UvA

# Previous lecture

o Machine learning paradigm for neural networks

o Backpropagation algorithm, backbone for training neural networks

o Neural network == modular architecture

o Visited different modules, saw how to implement and check them

# Lecture overview

- How to define our model and optimize it in practice

- Data preprocessing and normalization

- Optimization methods

- Regularizations

- Architectures and architectural hyper-parameters

- Learning rate

- Weight initializations

- Good practices
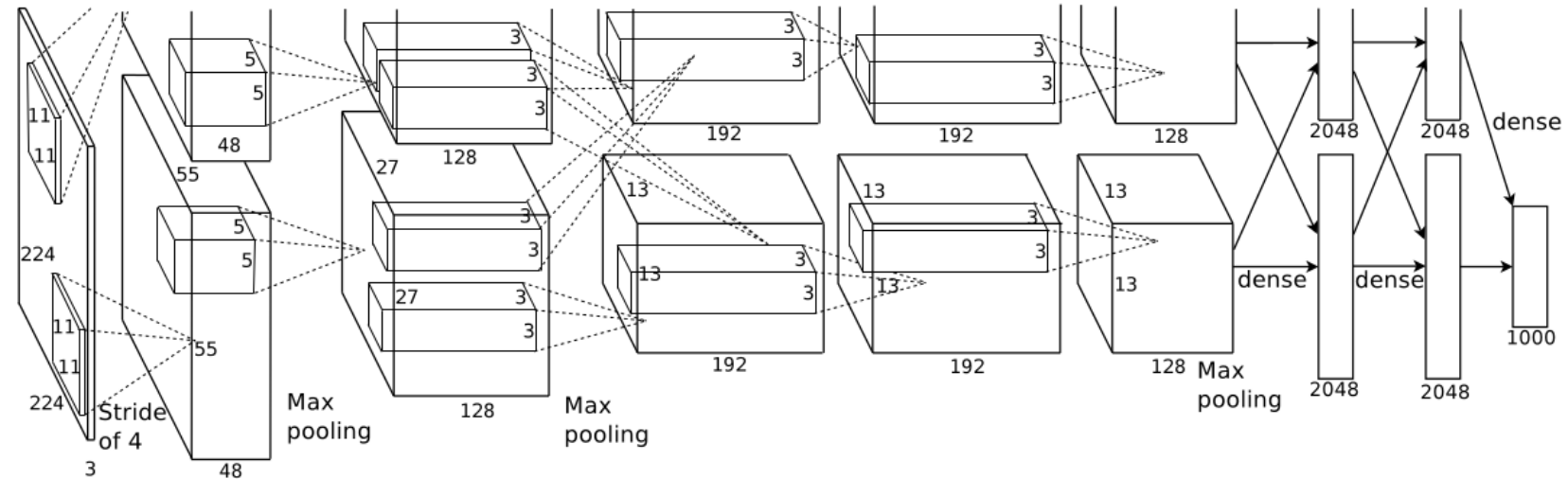
# Deeper into Neural Networks & Deep Neural Nets



Figure 2: An illustration of the architecture of our CNN, explicitly showing the delineation of responsibilities between the two GPUs. One GPU runs the layer-parts at the top of the figure while the other runs the layer-parts at the bottom. The GPUs communicate only at certain layers. The network's input is 150,528-dimensional, and the number of neurons in the network's remaining layers is given by 253,440–186,624–64,896–64,896–43,264–4096–4096–1000.

# A Neural/Deep Network in a nutshell

1. The Neural Network

$$a_L(x; \theta_{1,\ldots,L}) = h_L\left(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}(y, a_L(x; \theta_{1,\ldots,L}))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Pure Optimization vs Machine Leanring Training

o In pure optimization the end goal is finding the minimum (or maximum)
  - E.g., optimizingrailroad network in the Netherlands, the end goal is finding the optimal combination of train schedules, train availability, etc

o The goal is very direct
  - Formulate your problem mathematically as best as possible
  - Find the best possible solution of your mathematical solution

o In training Machine Learning algorithms finding the cost function is usually only a surrogate to your goal
  - You want to recognize cars from bikes (0-1 problem) in unknown images, but you optimize the classification log probabilities (continuous) in known images
  - Even the "optimal" parameters are not necessarily the best choice for your end goal

# Empirical Risk Minimization

o Differently from pure optimization which operates on the training data points, we ideally should optimize for

$$\min_{\theta} \mathrm{E}_{x,y\sim\hat{p}_{\mathrm{data}}}[\mathcal{L}(\theta; x, y)]$$

o Still, borrowing from optimization is the best way we can get satisfactory solutions to our problems by minimizing the empirical risk

$$\min_{\theta} \mathrm{E}_{x,y\sim\hat{p}_{\mathrm{data}}}[\mathcal{L}(\theta; x, y)] + \lambda\Omega(\theta) = \frac{1}{m}\Sigma_{i=1}^{m}\mathcal{L}(h(x_i;\theta),y_i)+\lambda\Omega(\theta)$$

◦ That is, the risk on the available training samples

# Optimizing with Stochastic Gradient Descent

1. The Neural Network

$$a_L\left(x; \theta_{1,...,L}\right) = h_L\left(h_{L-1}(...h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y)\subseteq(X,Y)} \mathcal{L}(y, a_L\left(x; \theta_{1,...,L}\right))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Gradient Descent

o To optimize a given loss function, most machine learning methods rely on Gradient Descent and variants

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t g^{(t)}$$

◦ $g^{(t)}: gradient$

o If the gradient is computed on the whole training set we have batch gradient descent

$$g^{(t)} = \frac{1}{m} \sum_{i=1}^{m} \nabla_\theta \mathcal{L}(\theta; x_i, y_i)$$

◦ This is only an approximation to the true gradient, as it is computed empirically from all the available training samples $(x_i, y_i)$

# Advantages of Batch Gradient Descent batch learning

o Conditions of convergence well understood

o Acceleration techniques can be applied
  ◦ Second order (Hessian based) optimizations are possible
  ◦ Measuring not only gradients, but also curvatures of the loss surface

o Simpler theoretical analysis on weight dynamics and convergence rates

# Still, optimizing with Gradient Descent is not perfect

o Often loss surfaces are
- non-quadratic
- highly non-convex
- very high-dimensional

o Datasets are typically really large to compute complete gradients

o No real guarantee that
- the final solution will be good
- we converge fast to final solution
- or that there will be convergence

# Stochastic Gradient Descent (SGD)

o The gradient equals an expectation $\mathrm{E}(\nabla_\theta \mathcal{L})$. In practice, we compute the mean from samples $\mathrm{E}(\nabla_\theta \mathcal{L}) = \frac{1}{m} \sum \nabla_\theta \mathcal{L}_i$.

o The standard error of this first approximation is given by $\frac{\sigma}{\sqrt{m}}$

   ◦ So, the error drops sublinearly with $m$. To compute <u>2x more</u> accurate gradients, we need <u>4x data</u> points

   ◦ And what's the point anyways, since our loss function is only a surrogate?

o Introduce a second approximation in computing the gradients

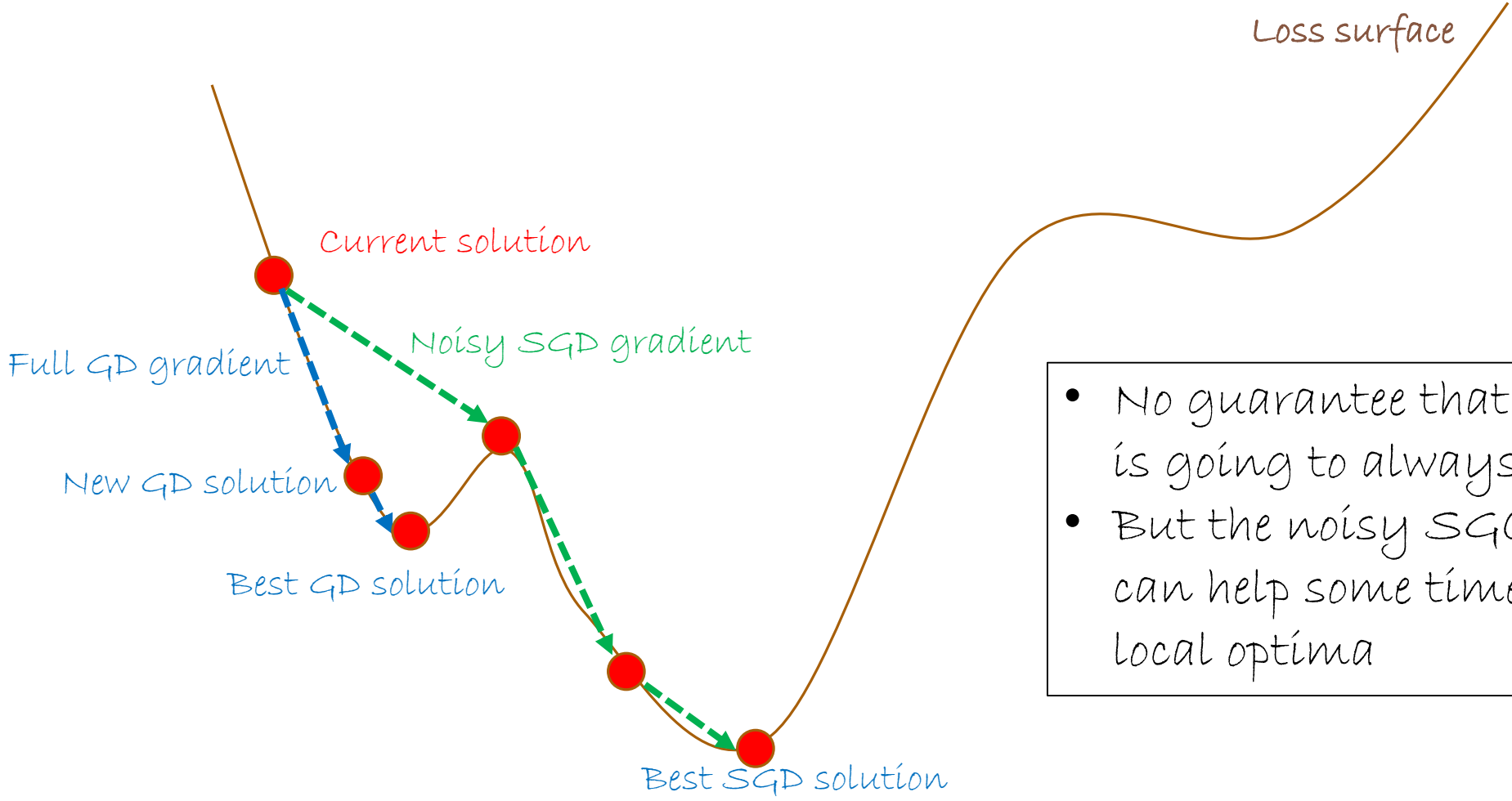   ◦ Stochastically sample "mini-training" sets ("mini-batches") from dataset $D$

$$B_j = sample(D)$$
$$\theta^{(t+1)} = \theta^{(t)} - \frac{\eta_t}{|B_j|} \sum_{i \in B_j} \nabla_\theta \mathcal{L}_i$$

o When computed from continuous streams of data (training data only seen once) SGD minimizes generalization error

   ◦ Intuitively, sampling continuously we sample from the true data distribution: $p_{\text{data}}$ not $\hat{p}_{\text{data}}$

# Some advantages of SGD

o Much faster than Gradient Descent

o Results are often better

o Also suitable for datasets that change over time

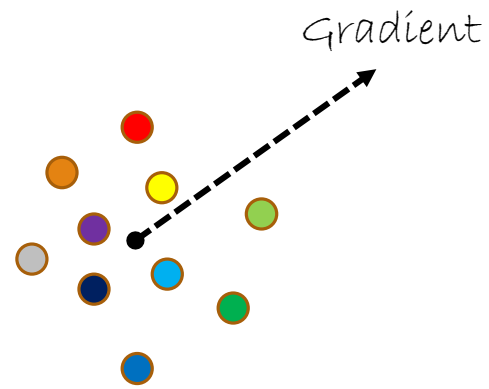o Variance of gradients increases when batch size decreases

# SGD is often better



Loss surface

Current solution

Noisy SGD gradient

Full GD gradient

New GD solution

Best GD solution

Best SGD solution

- No guarantee that this is what is going to always happen.
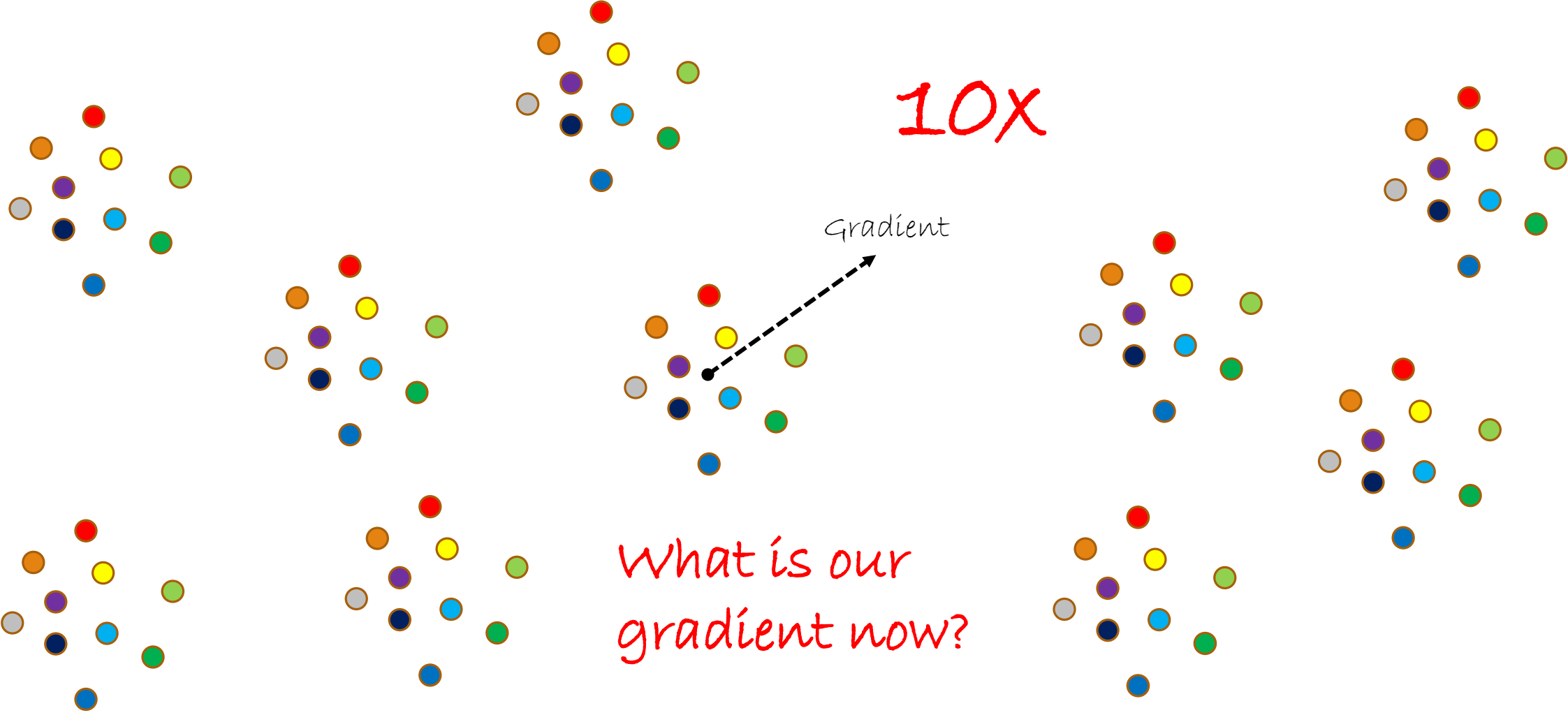- But the noisy SGC gradients can help some times escaping local optima

# SGD is often better

o (A bit) Noisy gradients act as regularization

o Gradient Descent → Complete gradients

o Complete gradients fit optimally the (arbitrary) data we have, not the distribution that generates them
  ◦ All training samples are the "absolute representative" of the input distribution
  ◦ Test data will be no different than training data
  ◦ Suitable for traditional optimization problems: "find optimal route"
  ◦ But for ML we cannot make this assumption → test data are always different

o Stochastic gradients → sampled training data sample roughly representative gradients
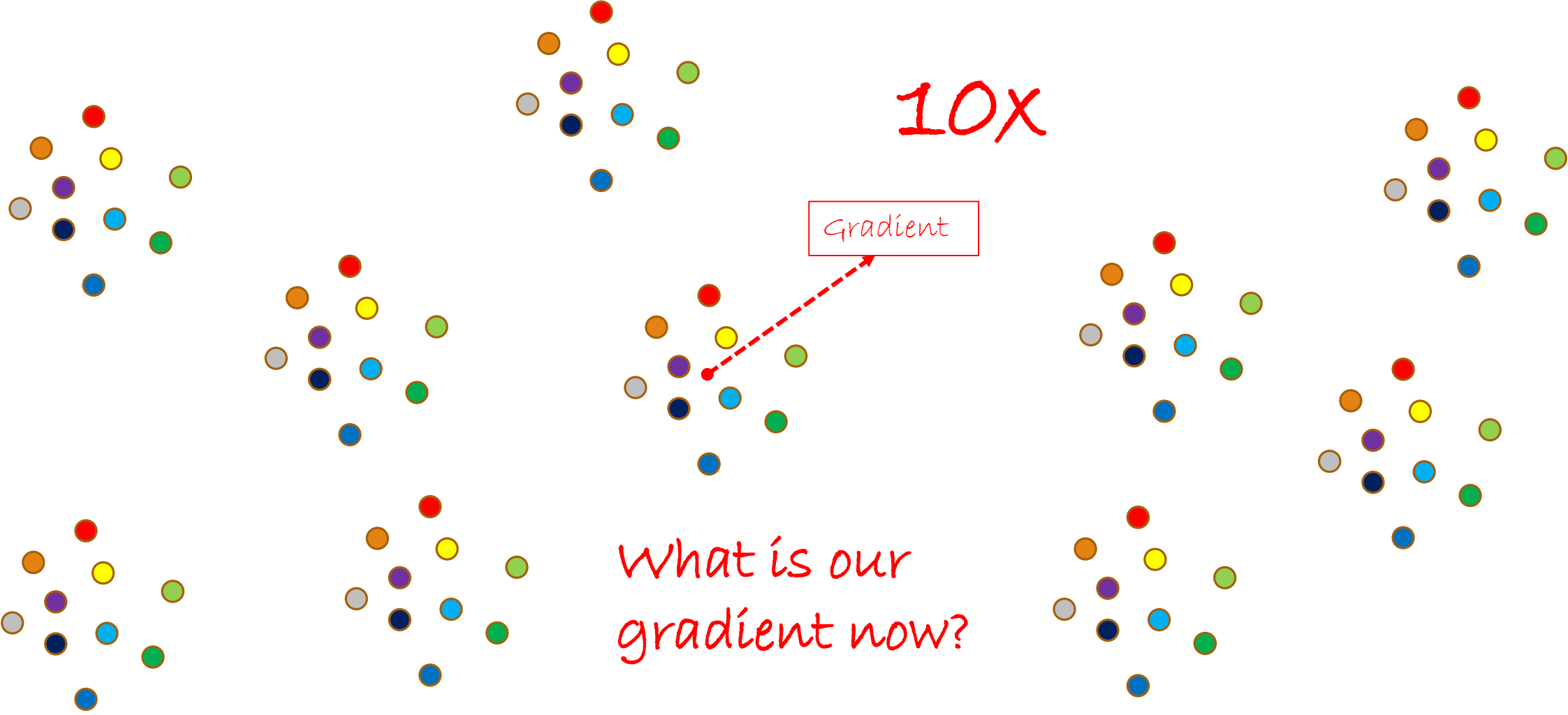  ◦ Model does not overfit to the particular training samples

# SGD is faster



Gradient

# SGD is faster



10X

Gradient

What is our
gradient now?

# SGD is faster



10x

Gradient

What is our gradient now?

# SGD is faster

o Of course in real situations data do not replicate

o However, after a sizeable amount of data there are clusters of data that are similar

o Hence, the gradient is approximately alright

o Approximate alright is great, is even better in many cases actually

# SGD for dynamically changed datasets

○ Often datasets are not "rigid"

○ Imagine Instagram
  ◦ Let's assume 1 million of **new** images uploaded per week and we want to build a "cool picture" classifier
  ◦ Should "cool pictures" from the previous year have the same as much influence?
  ◦ No, the learning machine should track these changes

○ With GD these changes go undetected, as results are averaged by the many more "past" samples
  ◦ Past "over-dominates"

○ A properly implemented SGD can track changes much better and give better models
  ◦ [LeCun2002]
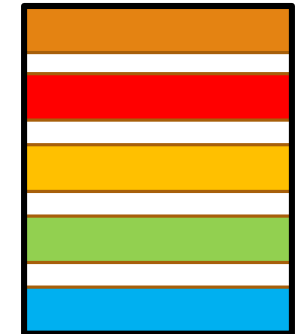

Popular today


Popular in 2014


Popular in 2010

# Shuffling examples

Dataset

o Applicable only with SGD

o Choose samples with maximum information content

o Mini-batches should contain examples from different classes
   ◦ As different as possible

Shuffling
at epoch t

o Prefer samples likely to generate larger errors
   ◦ Otherwise gradients will be small ➔ slower learning
   ◦ Check the errors from previous rounds and prefer "hard examples"
   ◦ Don't overdo it though :P, beware of outliers

Shuffling
at epoch t+1

o In practice, split your dataset into mini-batches
   ◦ Each mini-batch is as class-divergent and rich as possible
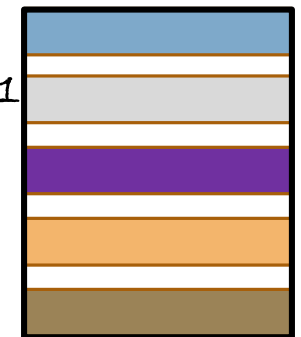   ◦ New epoch ➔ to be safe new batches & new, randomly shuffled examples

# Backpropagation again

o **Step 1.** Compute forward propagations for all layers recursively

$$a_l = h_l(x_l) \text{ and } x_{l+1} = a_l$$

o **Step 2.** Once done with forward propagation, follow the reverse path.
  ◦ Start from the last layer and for each new layer compute the gradients
  ◦ Cache computations when possible to avoid redundant operations

$$\frac{\partial \mathcal{L}}{\partial a_l} = \left(\frac{\partial a_{l+1}}{\partial x_{l+1}}\right)^T \cdot \frac{\partial \mathcal{L}}{\partial a_{l+1}}$$

$$\frac{\partial \mathcal{L}}{\partial \theta_l} = \frac{\partial a_l}{\partial \theta_l} \cdot \left(\frac{\partial \mathcal{L}}{\partial a_l}\right)^T$$

o **Step 3.** Use the gradients $\frac{\partial \mathcal{L}}{\partial \theta_l}$ with Stochastic Gradient Descent to train

# In practice

o SGD is preferred to Gradient Descent

o Training is orders faster
  ◦ In real datasets Gradient Descent is not even realistic

o Solutions generalize better
  ◦ More efficient → larger datasets
  ◦ Larger datasets → better generalization

o How many samples per mini-batch?
  ◦ Hyper-parameter, trial & error
  ◦ Usually between 32-256 samples

# Challenges in optimization

o Ill conditioning

  ◦ Let's check the 2$^{nd}$ order Taylor dynamics of optimizing the cost function

$$\mathcal{L}(\theta) = \mathcal{L}(\theta') + (\theta - \theta')^{\mathrm{T}}g + \frac{1}{2}(\theta - \theta')^{\mathrm{T}}\mathrm{H}(\theta - \theta') \quad (\mathrm{H}:\text{Hessian})$$

$$\mathcal{L}(\theta' - \varepsilon g) \approx \mathcal{L}(\theta) - \varepsilon g^{\mathrm{T}}g + \frac{1}{2}g^{T}Hg$$

  ◦ Even if the gradient $g$ is strong, if $\frac{1}{2}g^{T}Hg > \varepsilon g^{\mathrm{T}}g$ the cost will increase

o Local minima

  ◦ Non-convex optimization produces lots of equivalent, local minima

o Plateaus

o Cliffs and exploding gradients

o Long-term dependencies

# Data preprocessing & normalization

1. The Neural Network

$$a_L\left(\textcolor{red}{x}; \theta_{1,\ldots,L}\right) = h_L\left(h_{L-1}(\ldots h_1(\textcolor{red}{x}, \theta_1), \theta_{L-1}), \theta_L\right)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y)\subseteq(X,Y)} \mathcal{L}(y, a_L\left(\textcolor{red}{x}; \theta_{1,\ldots,L}\right))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Data pre-processing

o Center data to be roughly 0

　◦ Activation functions usually "centered" around 0

　◦ Convergence usually faster

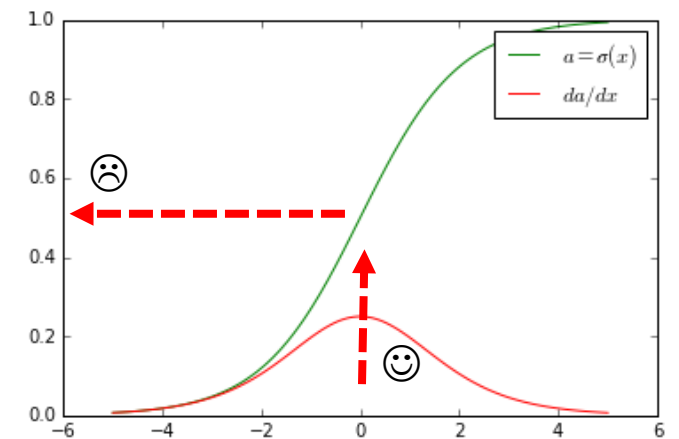　◦ Otherwise bias on gradient direction ➔ might slow down learning

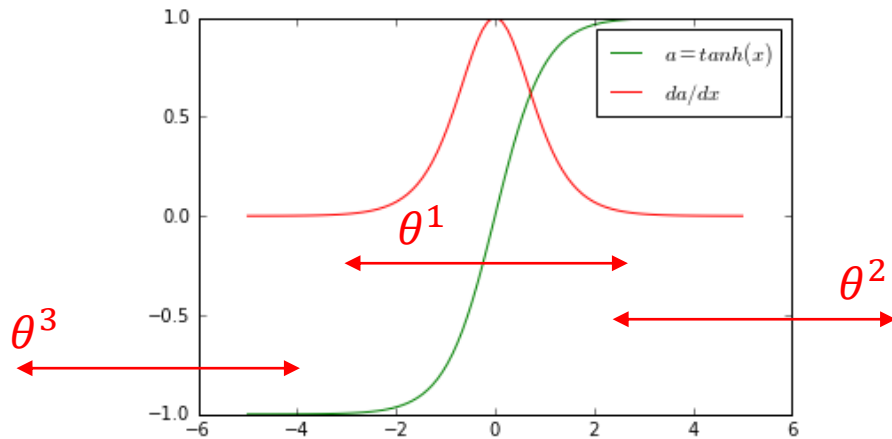ReLU ☺                           tanh(x) ☺                           σ(x) ☹

# Data pre-processing

○ Scale input variables to have similar diagonal covariances $c_i = \sum_j (x_i^{(j)})^2$

  ◦ Similar covariances $\rightarrow$ more balanced rate of learning for different weights

  ◦ Rescaling to 1 is a good choice, unless some dimensions are less important

$x = [x^1, x^2, x^3]^T, \theta = [\theta^1, \theta^2, \theta^3]^T, a = \tanh(\theta^T x)$



$x^1, x^2, x^3 \rightarrow$ much different covariances

Generated gradients $\left.\dfrac{\mathrm{d}\mathcal{L}}{d\theta}\right|_{x^1, x^2, x^3}$ : much different
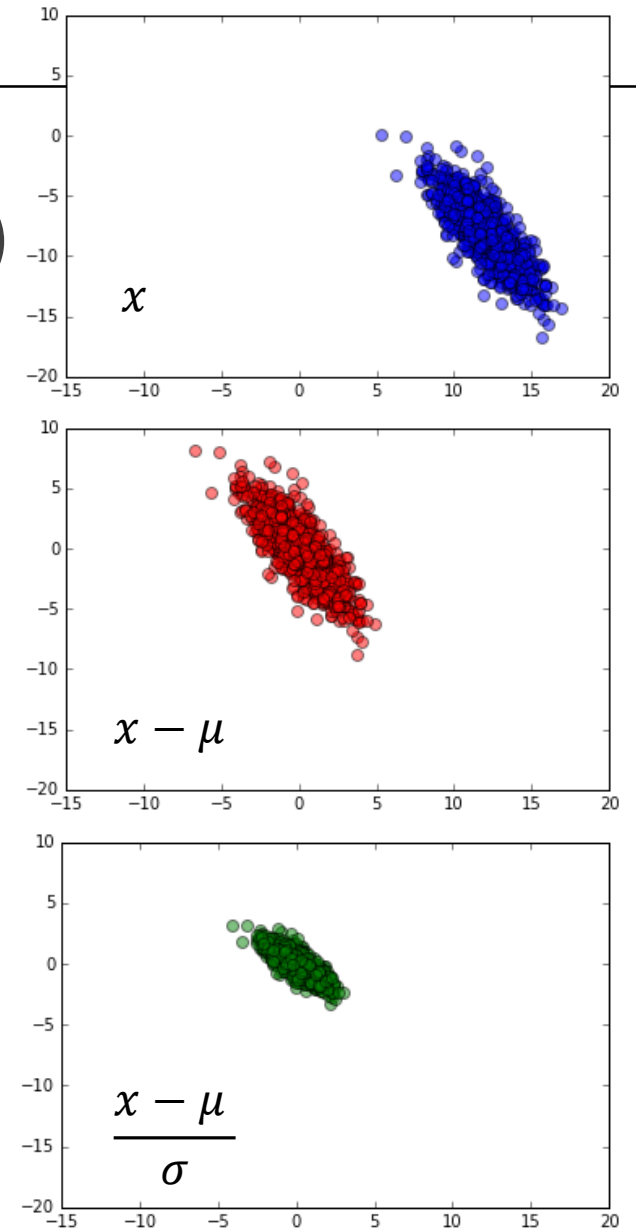
Gradient update harder: $\theta^{(t+1)} = \theta^{(t)} - \eta_t \begin{bmatrix} d\mathcal{L}/d\theta^1 \\ d\mathcal{L}/d\theta^2 \\ d\mathcal{L}/d\theta^3 \end{bmatrix}$

# Data pre-processing

o Input variables should be as decorrelated as possible

◦ Input variables are "more independent"

◦ Network is forced to find non-trivial correlations between inputs

◦ Decorrelated inputs → Better optimization

◦ Obviously not the case when inputs are by definition correlated (sequences)

o Extreme case

◦ extreme correlation (linear dependency) might cause problems [CAUTION]

# Normalization: $N(\mu, \sigma^2) = N(0, 1)$

o Input variables follow a Gaussian distribution (roughly)

o In practice:
◦ from training set compute mean and standard deviation
◦ Then subtract the mean from training samples
◦ Then divide the result by the standard deviation

# $N(\mu, \sigma^2) = N(0, 1)$ − Making things faster

o Instead of "per-dimension" ➔ all input dimensions simultaneously

o If dimensions have similar values (e.g. pixels in natural images)
  ◦ Compute one $(\mu, \sigma^2)$ instead of as many as the input variables
  ◦ Or the per color channel pixel average/variance

$$(\mu_{red}, \sigma^2_{red}), (\mu_{green}, \sigma^2_{green}), (\mu_{blue}, \sigma^2_{blue})$$

# Even simpler: Centering the input

o When input dimensions have similar ranges ...

o ... and with the right non-linearity ...

o ... centering might be enough
  ◦ e.g. in images all dimensions are pixels
  ◦ All pixels have more or less the same ranges

o Just make sure images have mean $0$ ($\mu = 0$)

# PCA Whitening

○ If $C$ the covariance matrix of your dataset, compute eigenvalues and eigenvectors with SVD
$$U, \Sigma, V^T = svd(C)$$



$X_{rot} = U^T X$
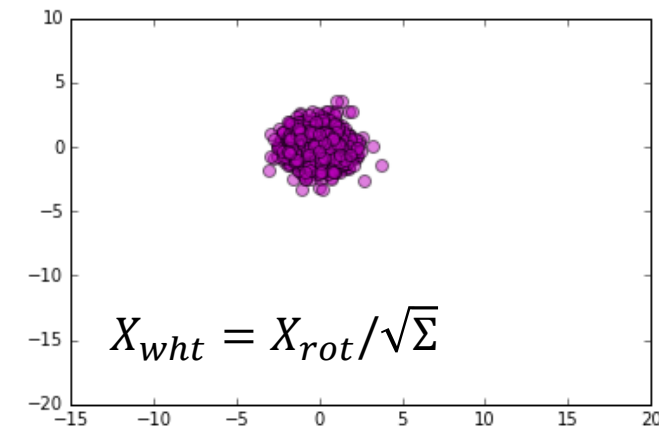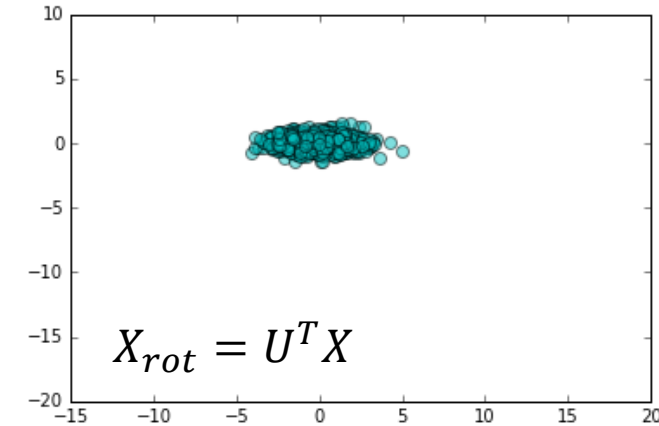
○ Decorrelate (PCA-ed) dataset by
$$X_{rot} = U^T X$$

  ◦ Subset of eigenvectors $U' = [u_1, \dots, u_q]$ to reduce data dimensions
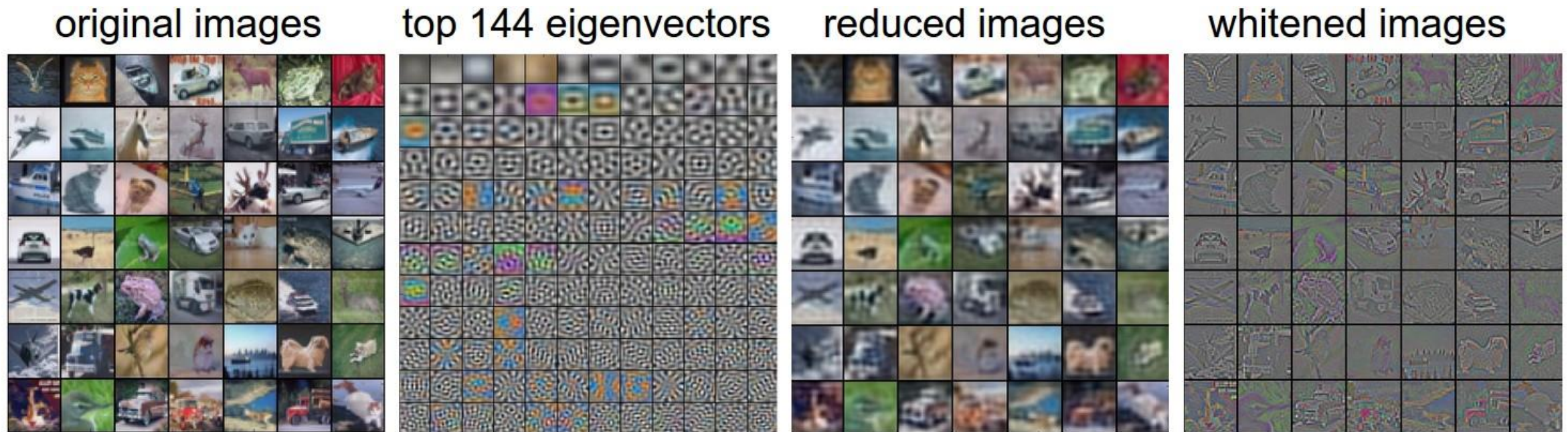
○ Scaling by square root of eigenvalues to whiten data
$$X_{wht} = X_{rot}/\sqrt{\Sigma}$$



$X_{wht} = X_{rot}/\sqrt{\Sigma}$

○ Not used much with Convolutional Neural Nets
  ◦ The zero mean normalization is more important

# Example



original images     top 144 eigenvectors     reduced images     whitened images

Images taken from A. Karpathy course website: http://cs231n.github.io/neural-networks-2/

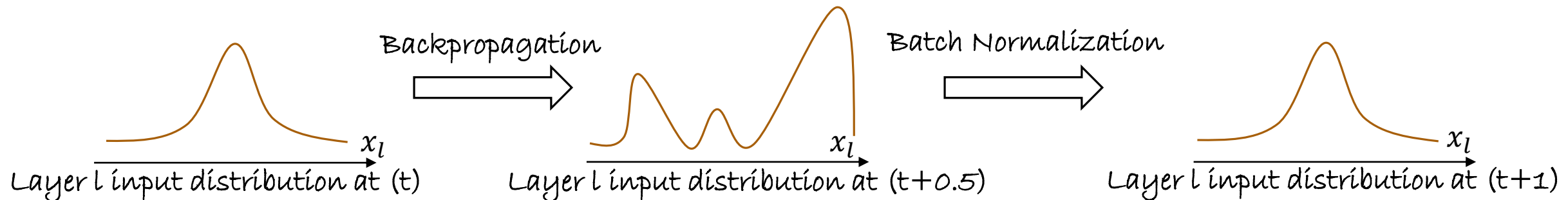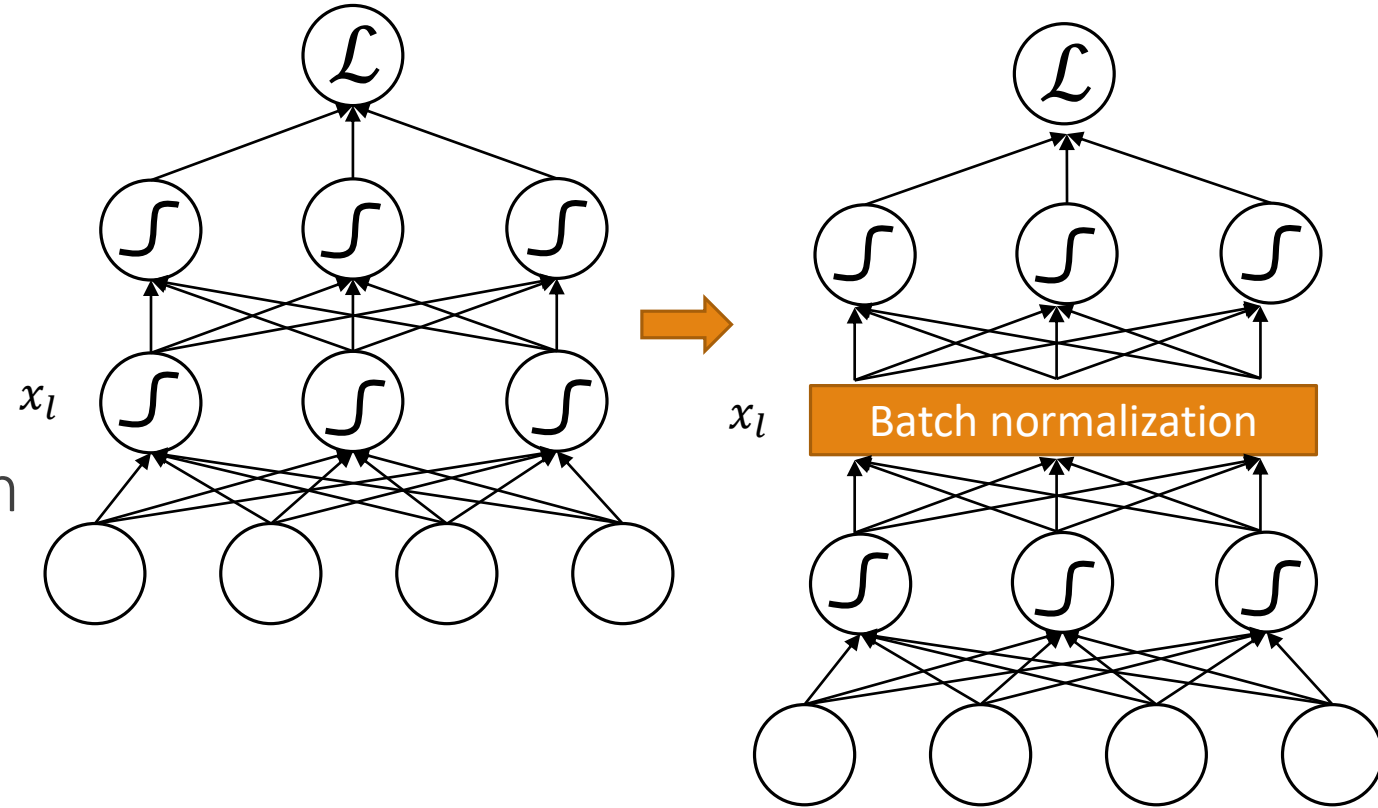# Data augmentation [Krizhevsky2012]
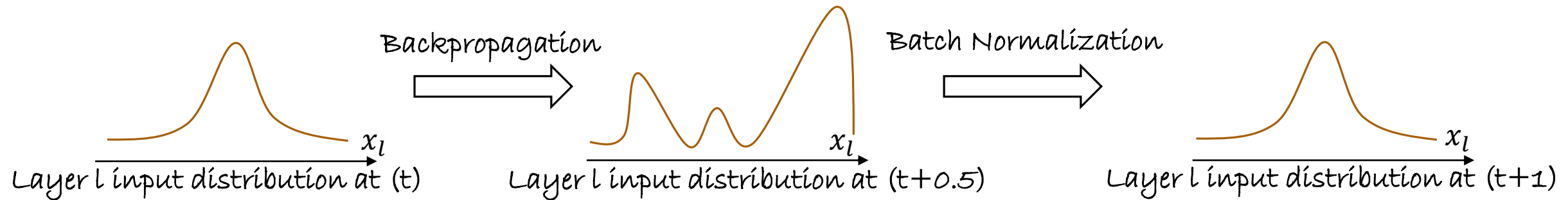


Original

Flip

Random crop

Contrast

Tint

# Batch normalization [Ioffe2015]

o Weights change ➔ the distribution of the layer inputs changes per round
  ◦ Covariance shift

o Normalize the layer inputs with batch normalization
  ◦ Roughly speaking, normalize $x_l$ to $N(0, 1)$ and rescale



Layer l input distribution at (t) — Backpropagation ➔ Layer l input distribution at (t+0.5) — Batch Normalization ➔ Layer l input distribution at (t+1)

# Batch normalization - Intuitively



Backpropagation

Batch Normalization

Layer $l$ input distribution at $(t)$

Layer $l$ input distribution at $(t+0.5)$

Layer $l$ input distribution at $(t+1)$

$x_l$

# Batch normalization – The algorithm

- $\mu_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m} x_i$          [compute mini-batch mean]

- $\sigma_{\mathcal{B}} \leftarrow \frac{1}{m}\sum_{i=1}^{m}(x_i - \mu_{\mathcal{B}})^2$    [compute mini-batch variance]

- $\widehat{x_i} \leftarrow \dfrac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \varepsilon}}$          [normalize input]

- $\widehat{y_i} \leftarrow \color{red}{\gamma} x_i + \color{red}{\beta}$          [scale and shift input]

Trainable parameters

# Batch normalization - Benefits

o Gradients can be stronger → higher learning rates → faster training

◦ Otherwise maybe exploding or vanishing gradients or getting stuck to local minima

o Neurons get activated in a near optimal "regime"

o Better model regularization

◦ Neuron activations not deterministic, depend on the batch
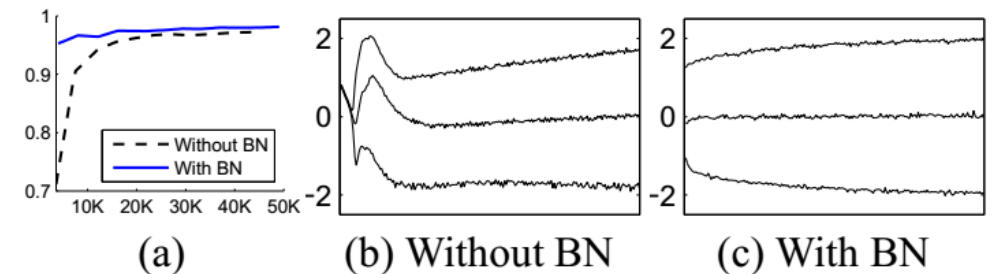
◦ Model cannot be overconfident



Figure 1: (a) *The test accuracy of the MNIST network trained with and without Batch Normalization, vs. the number of training steps. Batch Normalization helps the network train faster and achieve higher accuracy.* (b, c) *The evolution of input distributions to a typical sigmoid, over the course of training, shown as* $\{15, 50, 85\}th$ *percentiles. Batch Normalization makes the distribution more stable and reduces the internal covariate shift.*

# Regularization

1. The Neural Network

$$a_L(x; \theta_{1,\ldots,L}) = h_L(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_{1,\ldots,L}))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Regularization

o Neural networks typically have thousands, if not millions of parameters
  ◦ Usually, the dataset size smaller than the number of parameters

o Overfitting is a grave danger

o Proper weight regularization is crucial to avoid overfitting

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \ell\left(y, a_L\left(x; \theta_{1,\dots,L}\right)\right) + \textcolor{red}{\lambda\Omega(\theta)}$$

o Possible regularization methods
  ◦ $\ell_2$-regularization
  ◦ $\ell_1$-regularization
  ◦ Dropout

# $\ell_2$-regularization

- Most important (or most popular) regularization

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}(y, a_L(x; \theta_{1,\ldots,L})) + \frac{\lambda}{2} \sum_l \|\theta_l\|^2$$

- The $\ell_2$-regularization can pass inside the gradient descent update rule

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t(\nabla_\theta \mathcal{L} + \lambda\theta_l) \implies$$
$$\theta^{(t+1)} = (1 - \lambda\eta_t)\theta^{(t)} - \eta_t\nabla_\theta\mathcal{L}$$

*"Weight decay", because weights get smaller*

- $\lambda$ is usually about $10^{-1}, 10^{-2}$

# $\ell_1$-regularization

- $\ell_1$-regularization is one of the most important regularization techniques

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}\left(y, a_L\left(x; \theta_{1,\ldots,L}\right)\right) + \frac{\lambda}{2} \sum_l \|\theta_l\|$$

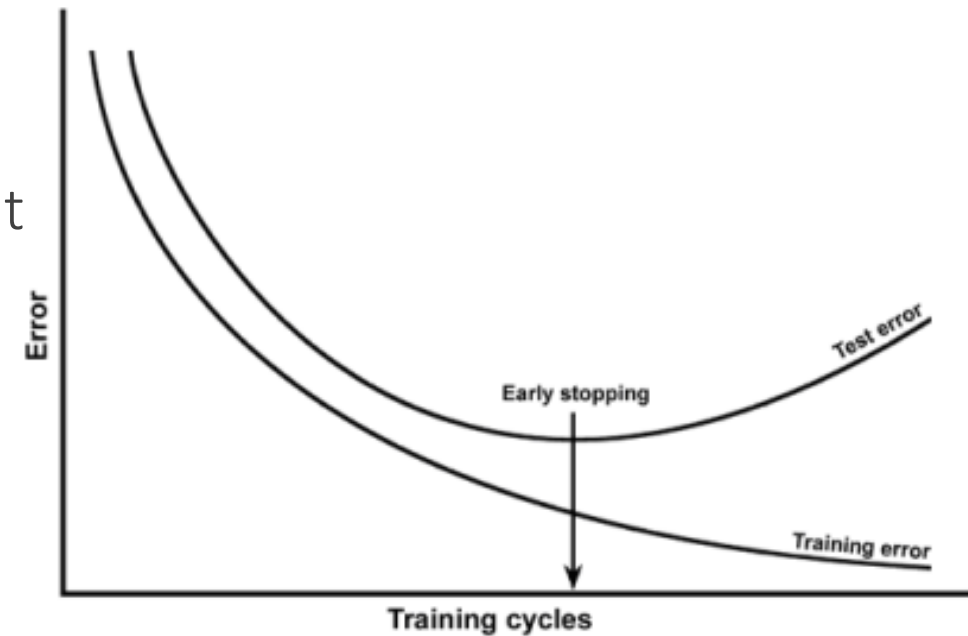- Also $\ell_1$-regularization passes inside the gradient descent update rule

$$\theta^{(t+1)} = \theta^{(t)} - \lambda\eta_t \frac{\theta^{(t)}}{|\theta^{(t)}|} - \eta_t \nabla_\theta \mathcal{L}$$

Sign function

- $\ell_1$-regularization $\rightarrow$ sparse weights
  - $\lambda \nearrow$ $\rightarrow$ more weights become 0

# Early stopping

o To tackle overfitting another popular technique is early stopping

o Monitor performance on a separate validation set

o Training the network will decrease training error, as well validation error (although with a slower rate usually)

o Stop when validation error starts increasing
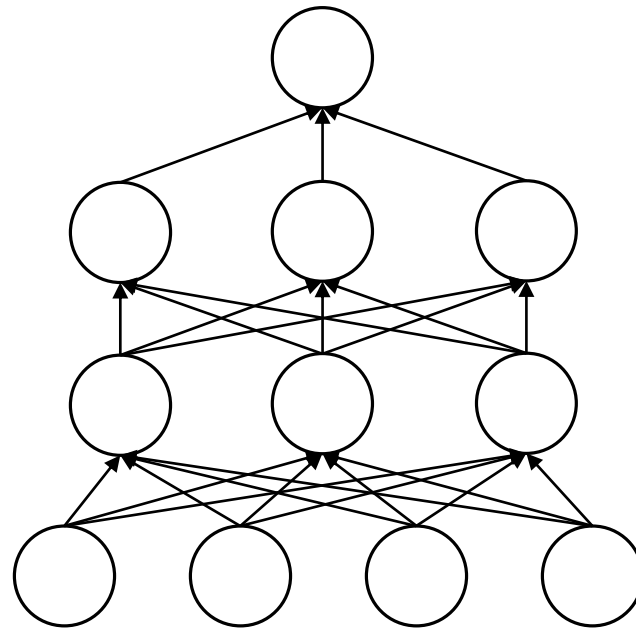  ◦ This quite likely means the network starts to overfit

# Dropout [Srivastava2014]

o During training setting activations randomly to 0
  ◦ Neurons sampled at random from a Bernoulli distribution with $p = 0.5$

o At test time all neurons are used
  ◦ Neuron activations reweighted by $p$

o Benefits
  ◦ Reduces complex co-adaptations or co-dependencies between neurons
  ◦ No "free-rider" neurons that rely on others
  ◦ Every neuron becomes more robust
  ◦ Decreases significantly overfitting
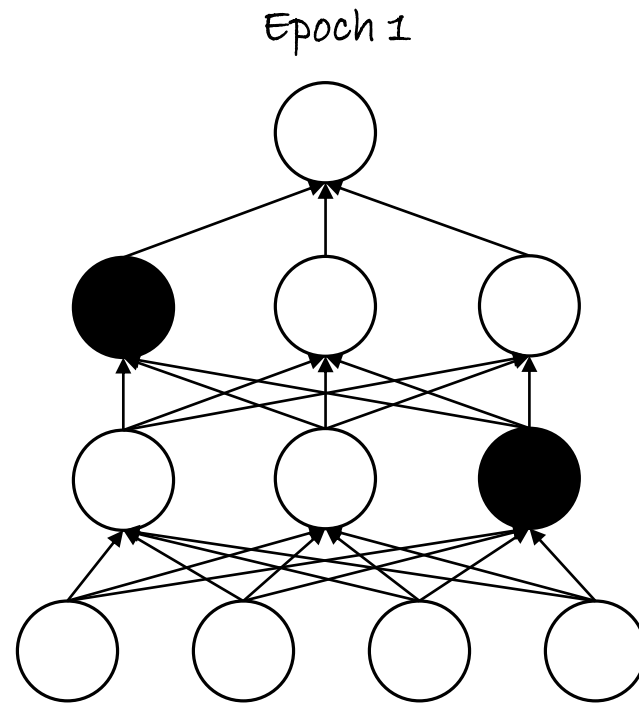  ◦ Improves significantly training speed

# Dropout

o Effectively, a different architecture at every training epoch
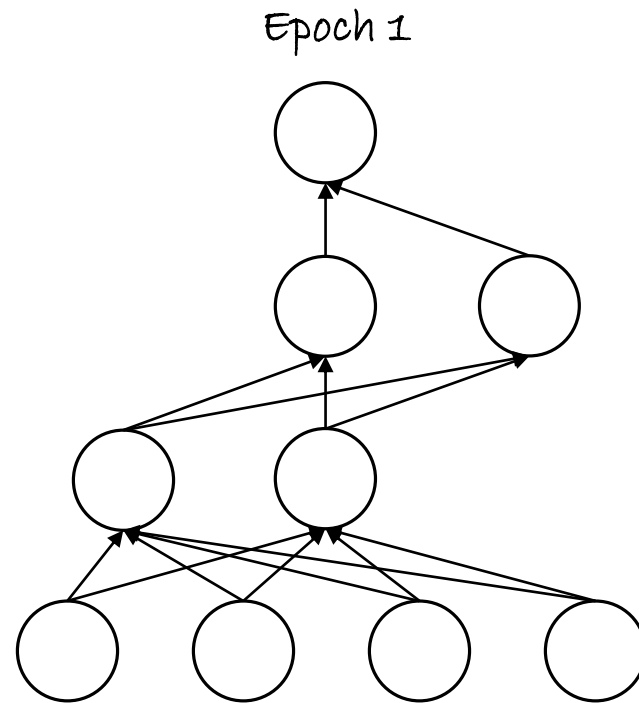  ◦ Similar to model ensembles

Original model

# Dropout

o Effectively, a different architecture at every training epoch
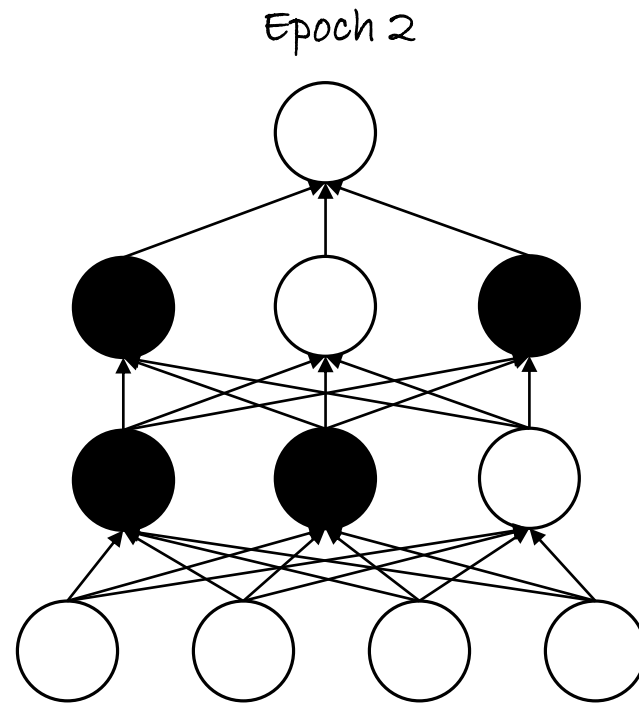  ◦ Similar to model ensembles

Epoch 1

# Dropout

o Effectively, a different architecture at every training epoch
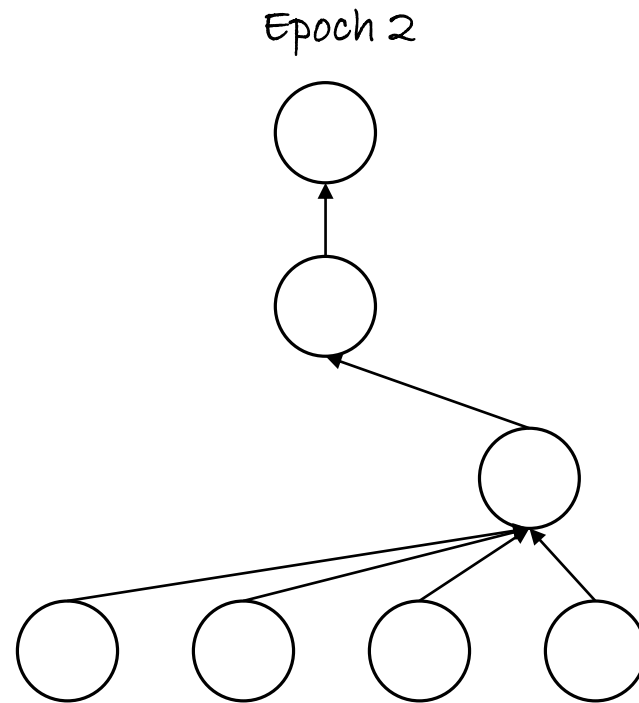  ◦ Similar to model ensembles



Epoch 1

# Dropout

o Effectively, a different architecture at every training epoch
  ◦ Similar to model ensembles



Epoch 2

# Dropout

o Effectively, a different architecture at every training epoch
  ◦ Similar to model ensembles



Epoch 2

# Architectural details

1. The Neural Network

$$a_L\left(x; \theta_{1,\ldots,L}\right) = h_L\left(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$$

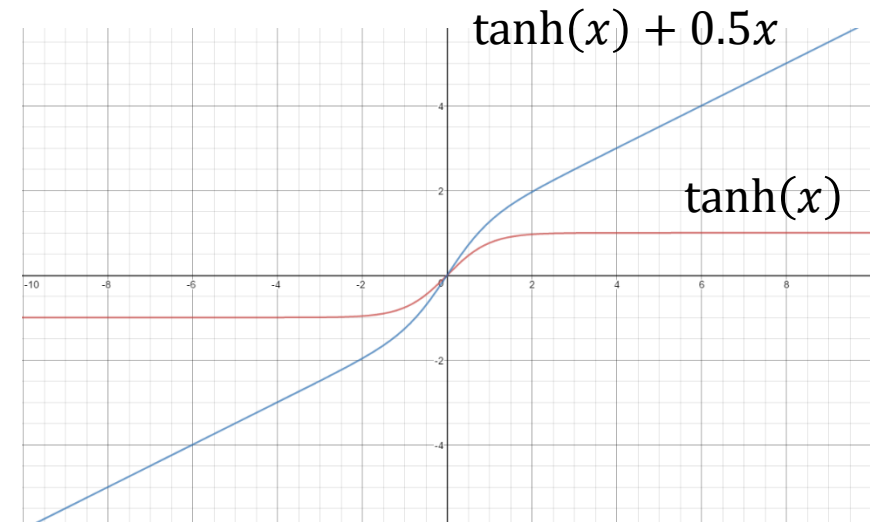2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y)\subseteq(X,Y)} \mathcal{L}(y, a_L\left(x; \theta_{1,\ldots,L}\right))$$

3. Optimizing with Gradient Descent based methods

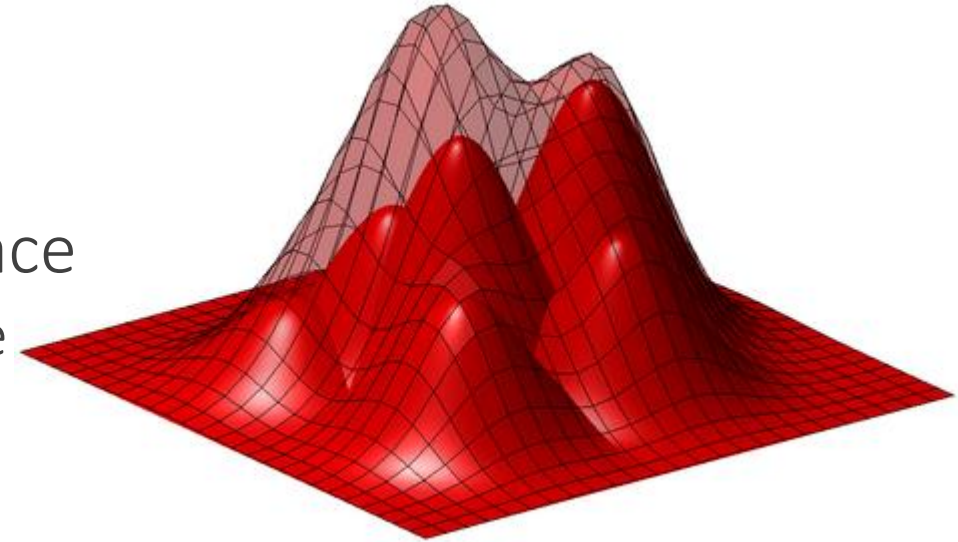$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Sigmoid-like activation functions

o Straightforward sigmoids not a very good idea

o Symmetric sigmoids converge faster
   ◦ E.g. tanh, returns a(x=0)=0
   ◦ Recommended sigmoid: $a = h(x) = 1.7159 \tanh(\frac{2}{3}x)$

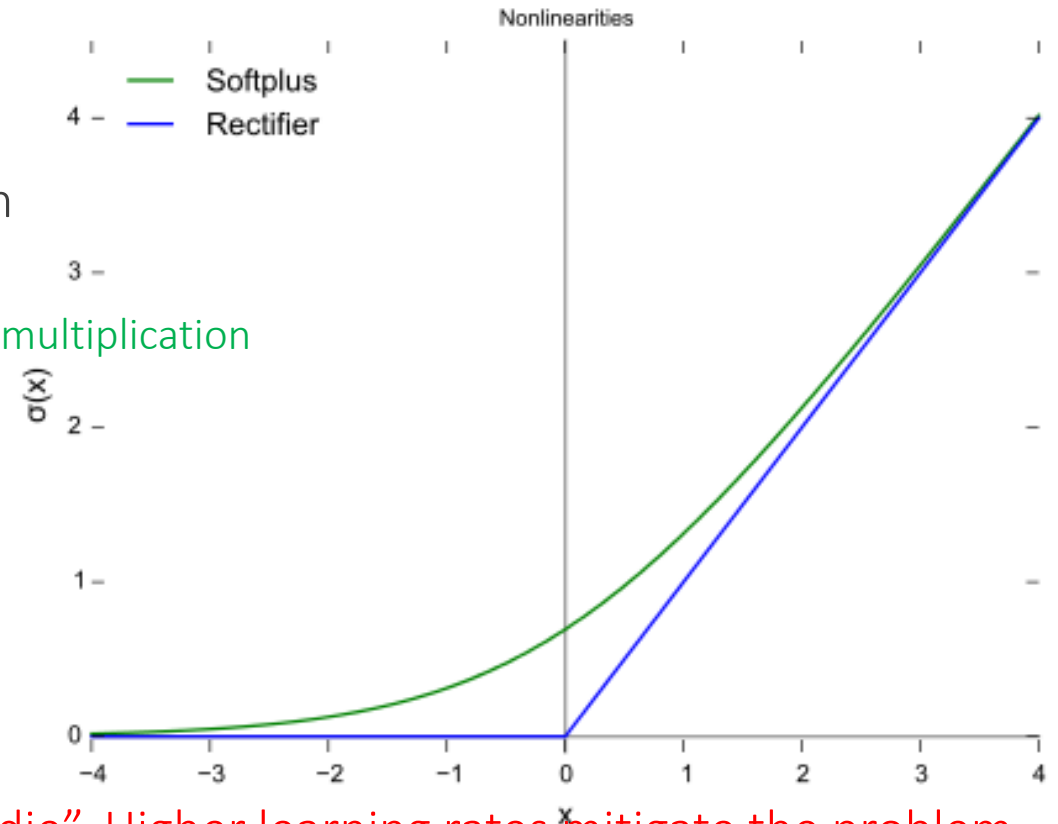o You can add a linear term to avoid flat areas
$$a = h(x) = \tanh(x) + \beta x$$



$\tanh(x) + 0.5x$

$\tanh(x)$

# RBFs vs "Sigmoids"

o RBF: $a = h(x) = \sum_j u_j \exp\left(-\beta_j (x - w_j)^2\right)$

o Sigmoid: $a = h(x) = \sigma(x) = \frac{1}{1+e^{-x}}$

o Sigmoids can cover the full feature space

o RBF's are much **more local** in the feature space
  ◦ Can be faster to train but with a more limited range
  ◦ Can give better set of basis functions
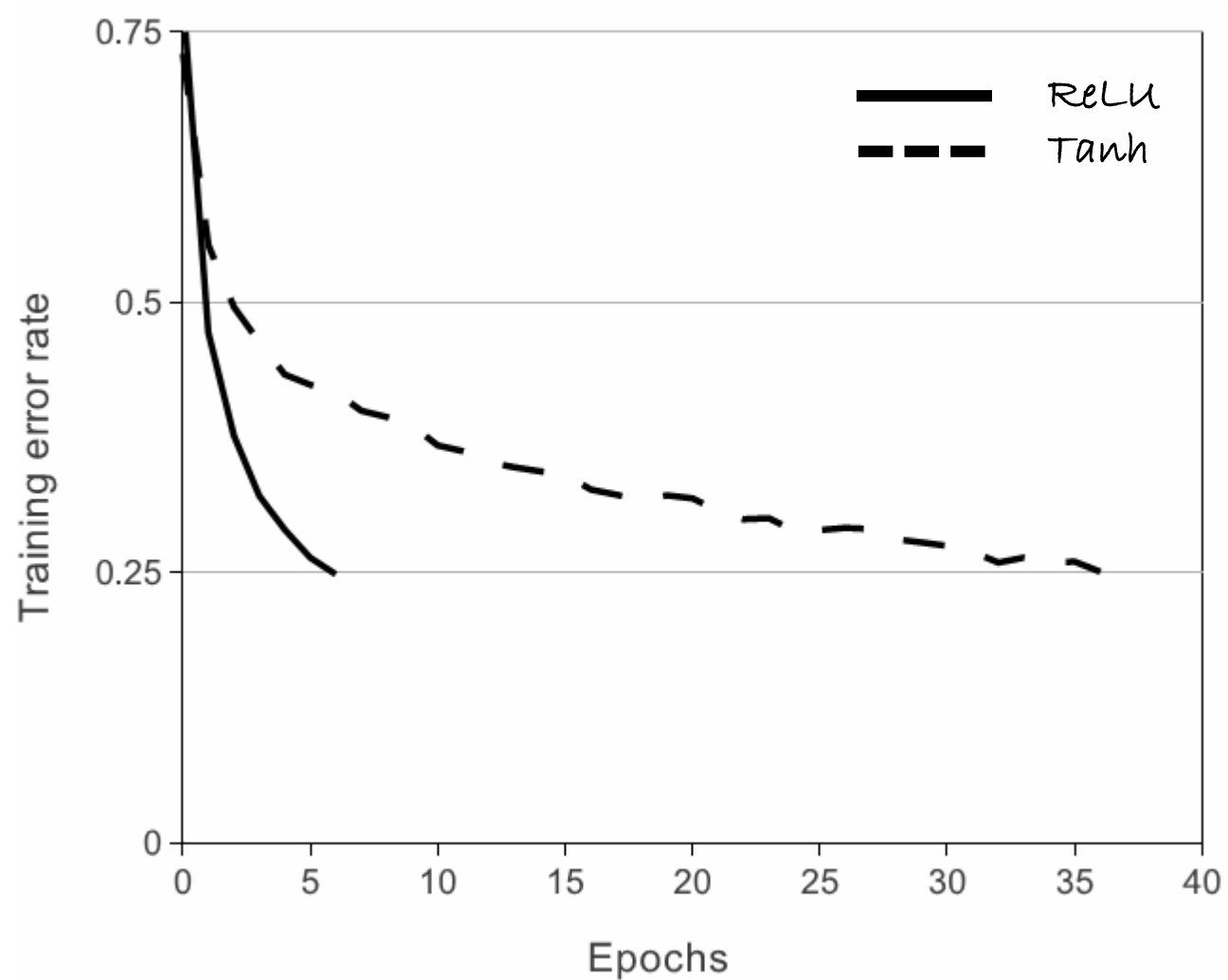  ◦ Preferred in lower dimensional spaces

# Rectified Linear Unit (ReLU) module [Krizhevsky2012]

o Activation function $a = h(x) = \max(0, x)$

o Gradient wrt the input $\dfrac{\partial a}{\partial x} = \begin{cases} 0, & if\ x \leq 0 \\ 1, & if\ x > 0 \end{cases}$

o Very popular in computer vision and speech recognition

o Much faster computations, gradients
   ◦ No vanishing or exploding problems, only comparison, addition, multiplication

o People claim biological plausibility

o Sparse activations

o No saturation

o Non-symmetric

o Non-differentiable at 0

o A large gradient during training can cause a neuron to "die". Higher learning rates mitigate the problem
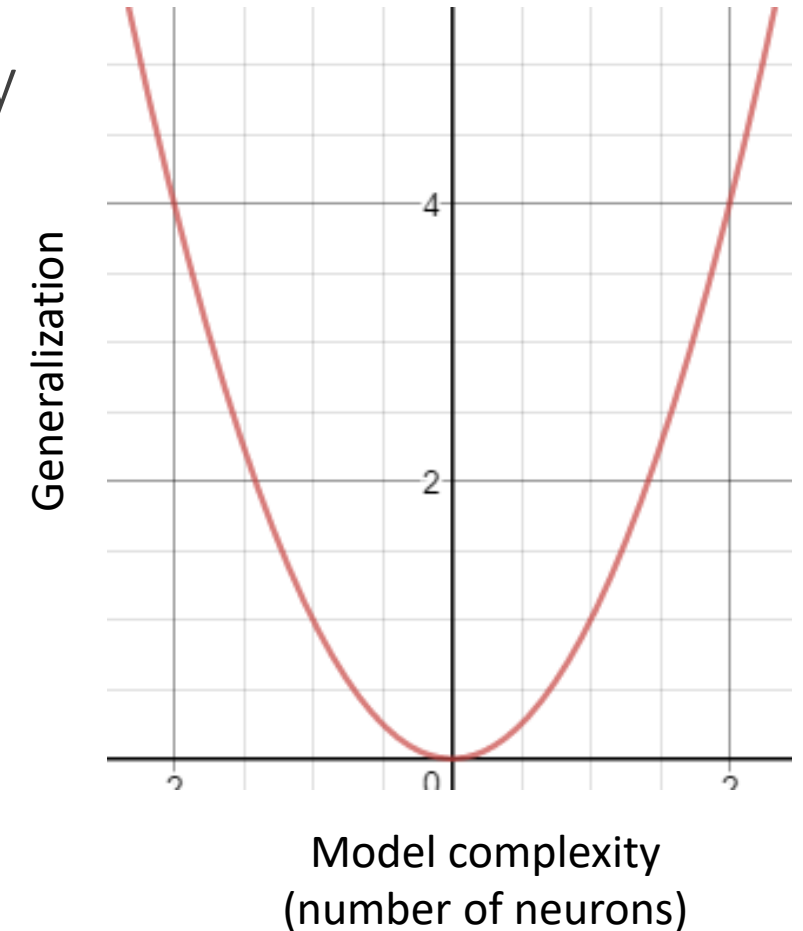
# ReLU convergence rate

# Architectural hyper-parameters

- Number of hidden layers

- Number of neuron in each hidden layer

- Type of activation functions

- Type and amount of regularization

# Number of neurons, number of hidden layers

o Dataset dependent hyperparameters

o Tip: **Start small** → increase complexity gradually
- e.g. start with a 2-3 hidden layers
- Add more layers → does performance improve?
- Add more neurons → does performance improve?

o Regularization is very important, use $\ell_2$
- Even if with very deep or wide network
- With strong $\ell_2$-regularization we avoid overfitting



Generalization

Model complexity
(number of neurons)

# Learning rate

1. The Neural Network

$$a_L\left(x; \theta_{1,\ldots,L}\right) = h_L\left(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y)\subseteq(X,Y)} \mathcal{L}(y, a_L\left(x; \theta_{1,\ldots,L}\right))$$

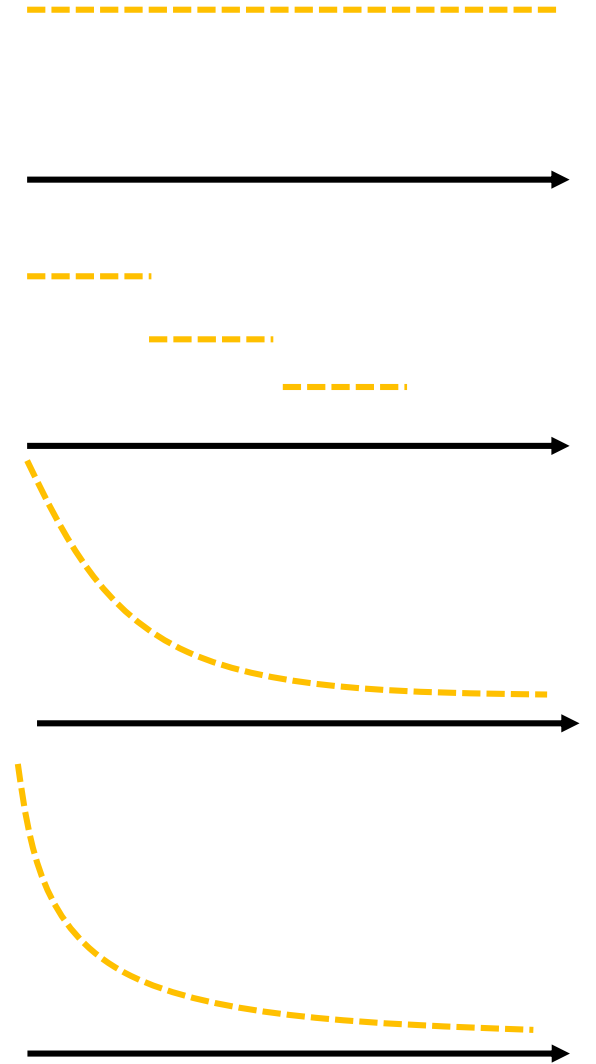3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Learning rate

o The right learning rate $\eta_t$ very important for fast convergence
  ◦ Too strong → gradients overshoot and bounce
  ◦ Too weak, → too small gradients → slow training

o Learning rate per weight is often advantageous
  ◦ Some weights are near convergence, others not

o Rule of thumb
  ◦ Learning rate of (shared) weights prop. to square root of share weight connections

o Adaptive learning rates are also possible, based on the errors observed
  ◦ [Sompolinsky1995]

# Learning rate schedules

o Constant
  ◦ Learning rate remains the same for all epochs

o Step decay
  ◦ Decrease (e.g. $\eta_t/T$ or $\eta_t/T$) every T number of epochs

o Inverse decay $\eta_t = \dfrac{\eta_0}{1+\varepsilon t}$

o Exponential decay $\eta_t = \eta_0 e^{-\varepsilon t}$

o Often step decay preferred
  ◦ simple, intuitive, works well and only a single extra hyper-parameter $T$ ($T$ =2, 10)

# Learning rate in practice

o Try several log-spaced values $10^{-1}, 10^{-2}, 10^{-3}, \ldots$ on a smaller set
  ◦ Then, you can narrow it down from there around where you get the lowest error

o You can decrease the learning rate every 10 (or some other value) full training set epochs
  ◦ Although this highly depends on your data

# Weight initialization

1. The Neural Network

$$a_L(x; \theta_{1,\ldots,L}) = h_L(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

2. Learning by minimizing empirical error
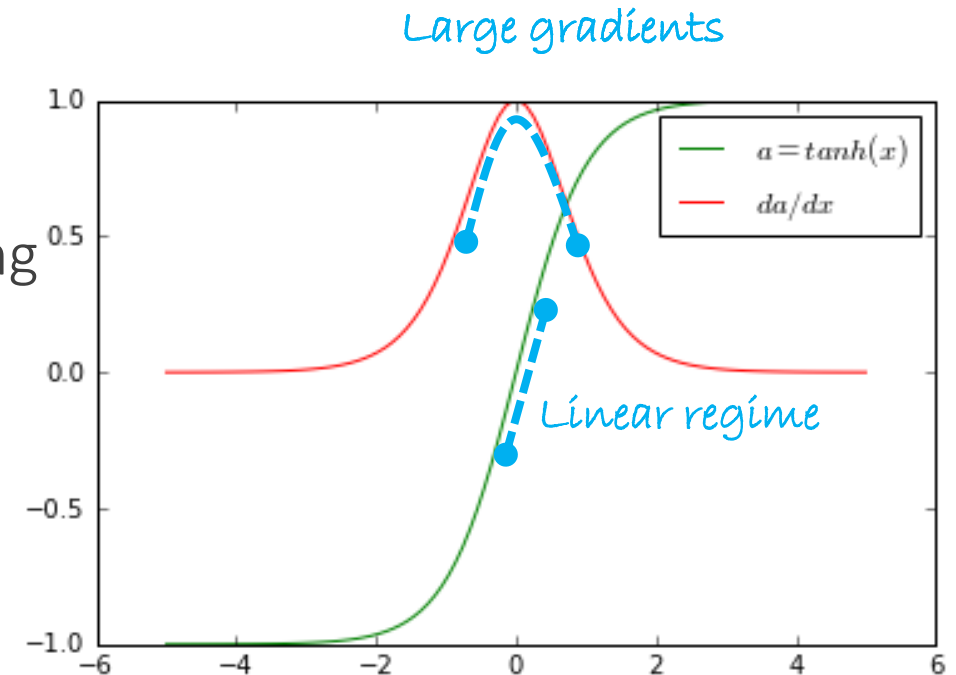
$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \ell(y, a_L(x; \theta_{1,\ldots,L}))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Weight initialization

○ There are few contradictory requirements

○ Weights need to be small enough
  ◦ around origin ($\vec{\mathbf{0}}$) for symmetric functions (tanh, sigmoid)
  ◦ When training starts better stimulate activation functions near their linear regime
  ◦ larger gradients → faster training

○ Weights need to be large enough
  ◦ Otherwise signal is too weak for any serious learning

Large gradients

Linear regime

# Weight initialization

- Weights must be initialized to preserve the variance of the activations during the forward and backward computations
  - Especially for deep learning
  - All neurons operate in their full capacity

  Question: Why similar input/output variance?

- Good practice: initialize weights to be asymmetric
  - Don't give save values to all weights (like all $\vec{\mathbf{0}}$)
  - In that case all neurons generate same gradient → no learning

- Generally speaking initialization depends on
  - non-linearities
  - data normalization

# Weight initialization

○ Weights must be initialized to preserve the variance of the activations during the forward and backward computations
- Especially for deep learning
- All neurons operate in their full capacity

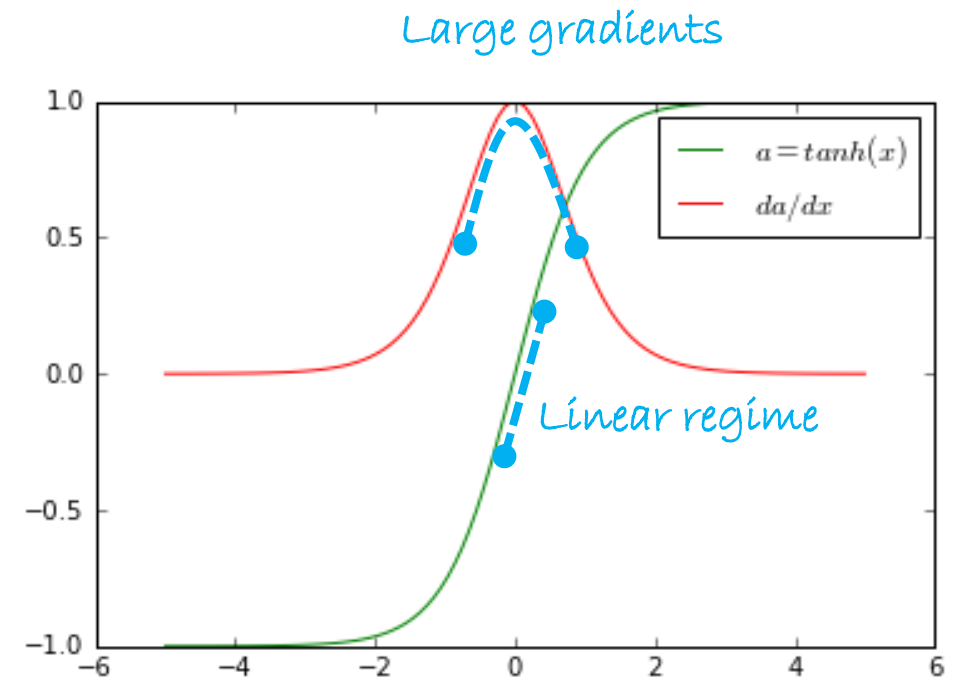Question: Why similar input/output variance?

Answer: Because the output of one module is the input to another

○ Good practice: initialize weights to be asymmetric
- Don't give save values to all weights (like all $\vec{\mathbf{0}}$)
- In that case all neurons generate same gradient → no learning

○ Generally speaking initialization depends on
- non-linearities
- data normalization

# One way of initializing sigmoid-like neurons

○ For tanh initialize weights from $\left[-\sqrt{\dfrac{6}{d_{l-1}+d_l}}, \sqrt{\dfrac{6}{d_{l-1}+d_l}}\right]$

　◦ $d_{l-1}$ is the number of input variables to the tanh layer and $d_l$ is the number of the output variables

○ For a sigmoid $\left[-4 \cdot \sqrt{\dfrac{6}{d_{l-1}+d_l}}, 4 \cdot \sqrt{\dfrac{6}{d_{l-1}+d_l}}\right]$

# Xavier initialization [Glorot2010]

o For $a = \theta x$ the variance is
$$var(a) = E[x]^2 var(\theta) + \mathrm{E}[\theta]^2 var(x) + var(x)var(\theta)$$

o Since $E[x] = E[\theta] = 0$
$$var(a) = var(x)var(\theta) \approx d \cdot var(x^i)var(\theta^i)$$

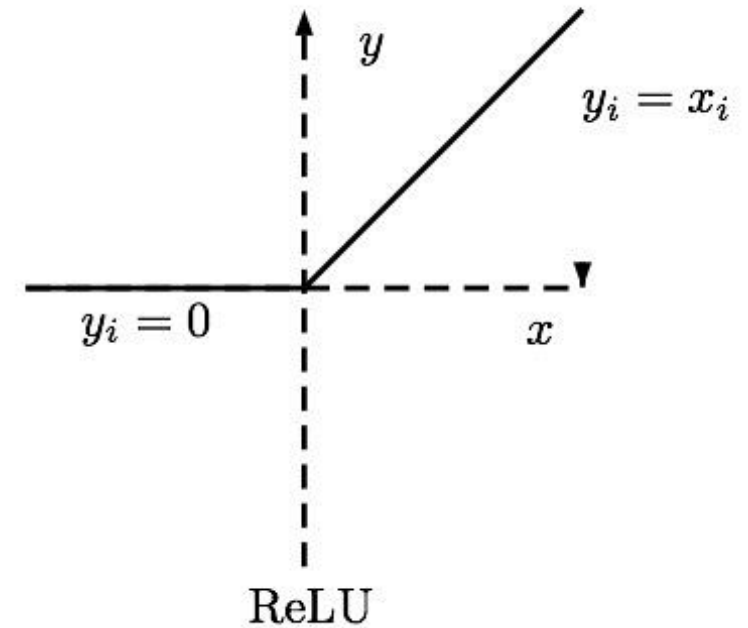o For $var(a) = var(x) \Rightarrow var(\theta^i) = \frac{1}{d}$

o Draw random weights from
$$\theta \sim N\left(0, \sqrt{1/d}\right)$$

where $d$ is the number of neurons in the input

# [He2015] initialization for ReLUs

o Unlike sigmoids, ReLUs ground to 0 the linear activations half the time

o Double weight variance
  ◦ Compensate for the zero flat-area →
  ◦ Input and output maintain same variance
  ◦ Very similar to Xavier initialization

o Draw random weights from $w \sim N\left(0, \sqrt{2/d}\right)$

   where $d$ is the number of neurons in the input



ReLU

# Loss functions

1. The Neural Network

$$a_L(x; \theta_{1,\dots,L}) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}(y, a_L(x; \theta_{1,\dots,L}))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Multi-class classification

- Our samples contains only one class
  - There is only one correct answer per sample

- Negative log-likelihood (cross entropy) + Softmax

$$\mathcal{L}(\theta; x, y) = -\sum_{c=1}^{C} y_c \log a_L^c \quad \text{for all classes } c = 1, \dots, C$$

- Hierarchical softmax when C is very large

*Is it a cat? Is it a horse? …*

- Hinge loss (aka SVM loss)

$$\mathcal{L}(\theta; x, y) = \sum_{\substack{c=1 \\ c \neq y}}^{C} \max(0, a_L^c - a_L^y + 1)$$

- Squared hinge loss

# Multi-class, multi-label classification

o  Each sample can have many correct answers

o  Hinge loss and the likes
  ◦ Also sigmoids would also work



o  Each output neuron is independent
  ◦ "Does this contain a car, <u>yes</u> or no?"
  ◦ "Does this contain a person, <u>yes</u> or no?"
  ◦ "Does this contain a motorbike, <u>yes</u> or no?"
  ◦ "Does this contain a horse, yes or <u>no</u>?"

o  Instead of "Is this a car, motorbike or person?"
  ◦ $p(car|x) = 0.55, p(m/bike|x) = 0.25, p(person|x) = 0.15, p(horse|x) = 0.05$
  ◦ $p(car|x) + p(m/bike|x) + p(person|x) + p(horse|x) = 1.0$

# Regression

o The good old Euclidean Loss

$$\mathcal{L}(\theta; x, y) = \frac{1}{2}|y - a_L|_2^2$$

o Or RBF on top of Euclidean loss

$$\mathcal{L}(\theta; x, y) = \sum_j u_j \exp(-\beta_j(y - a_L)^2)$$

o Or $\ell_1$ distance

$$\mathcal{L}(\theta; x, y) = \sum_j |y_j - a_L^j|$$

# Even better optimizations

1. The Neural Network

$$a_L(x; \theta_{1,\dots,L}) = h_L(h_{L-1}(\dots h_1(x, \theta_1), \theta_{L-1}), \theta_L)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}(y, a_L(x; \theta_{1,\dots,L}))$$

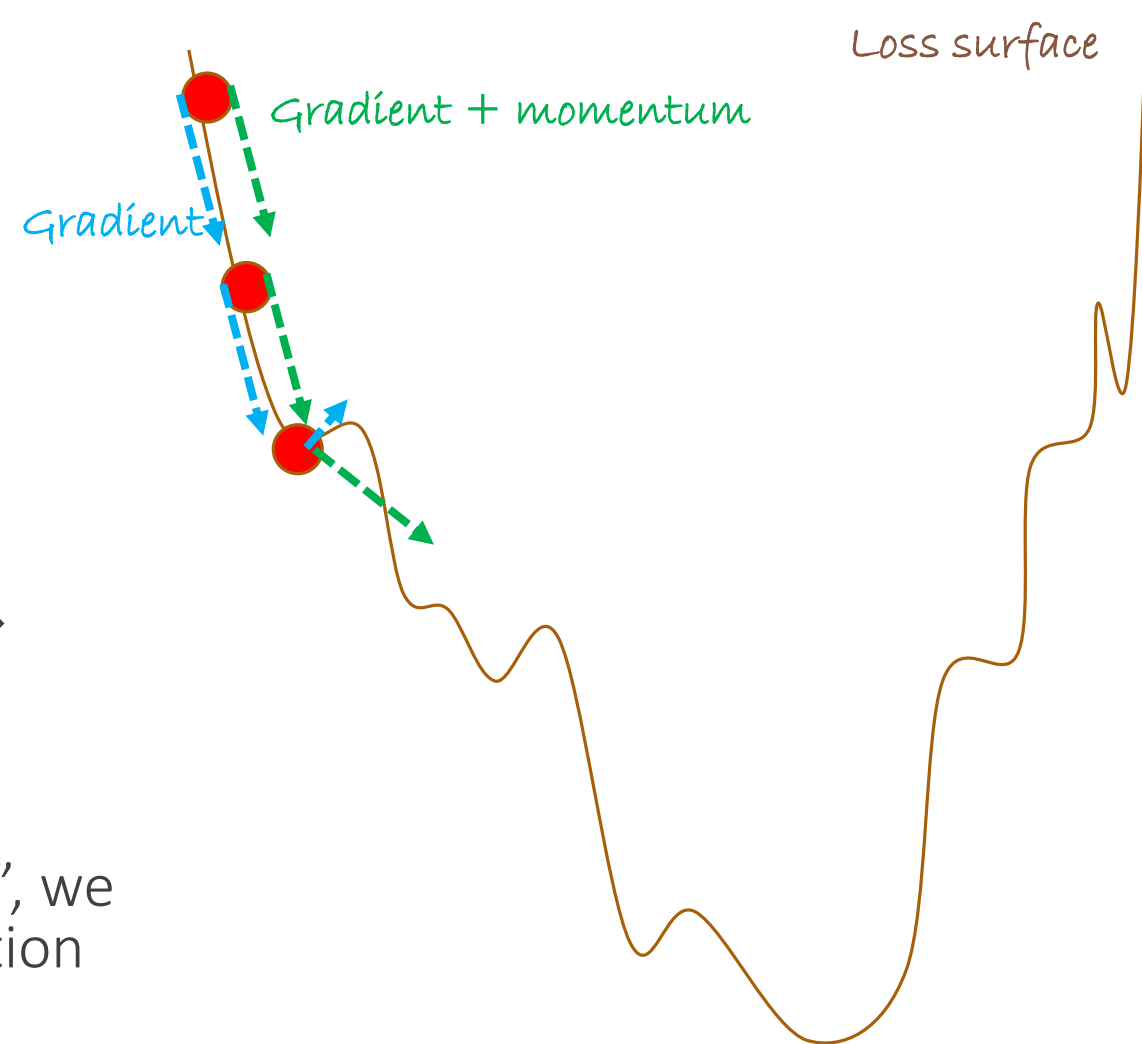3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Momentum

o Don't switch gradients all the time

o Maintain "momentum" from previous parameters

$$u^{(t+1)} = \gamma u^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$
$$\theta^{(t+1)} = \theta^{(t)} + u^{(t+1)}$$

o More robust gradients and learning → faster convergence

o Nice "physics"-based interpretation
   ◦ Instead of updating the position of the "ball", we update the velocity, which updates the position

Loss surface

Gradient + momentum

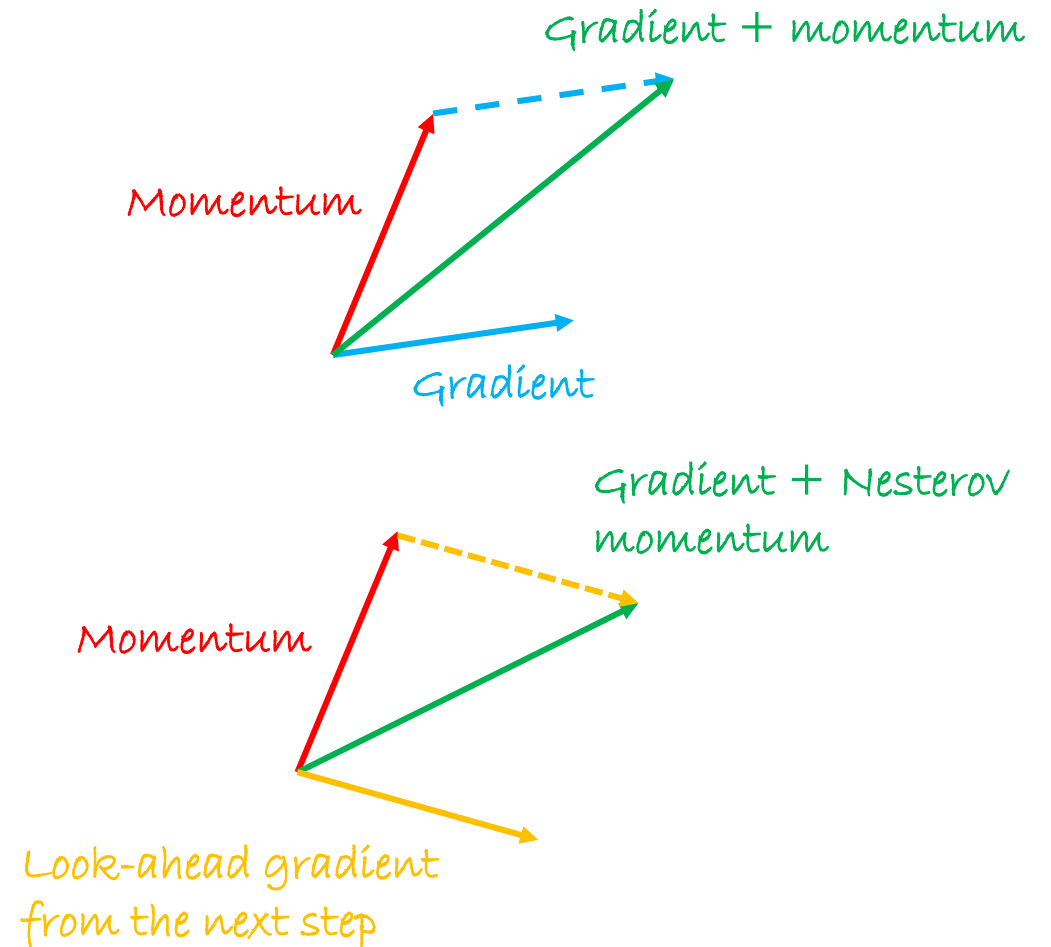Gradient

# Nesterov Momentum [Sutskever2013]

- Use the future gradient instead of the current gradient

$$\theta^{(t+0.5)} = \theta^{(t)} + \gamma u^{(t)}$$

$$u^{(t+1)} = \gamma u^{(t)} - \eta_t \nabla_{\theta^{(t+0.5)}} \mathcal{L}$$

$$\theta^{(t+1)} = \theta^{(t)} + u^{(t+1)}$$

- Better theoretical convergence

- Generally works better with Convolutional Neural Networks

# Second order optimization

o Normally all weights updated with same "aggressiveness"

  ◦ Often some parameters could enjoy more "teaching"

  ◦ While others are already about there

o Adapt learning per parameter

$$\theta^{(t+1)} = \theta^{(t)} - H_{\mathcal{L}}^{-1} \eta_t \nabla_\theta \mathcal{L}$$

o $H_{\mathcal{L}}$ is the Hessian matrix of $\mathcal{L}$: second-order derivatives

$$H_{\mathcal{L}}^{ij} = \frac{\partial \mathcal{L}}{\partial \theta_i \partial \theta_j}$$

# Second order optimization methods in practice

o Inverse of Hessian usually very expensive
  ◦ Too many parameters

o Approximating the Hessian, e.g. with the L-BFGS algorithm
  ◦ Keeps memory of gradients to approximate the inverse Hessian

o L-BFGS works alright with Gradient Descent. What about SGD?

o In practice SGD with some good momentum works just fine

# Other per-parameter adaptive optimizations

- Adagrad [Duchi2011]
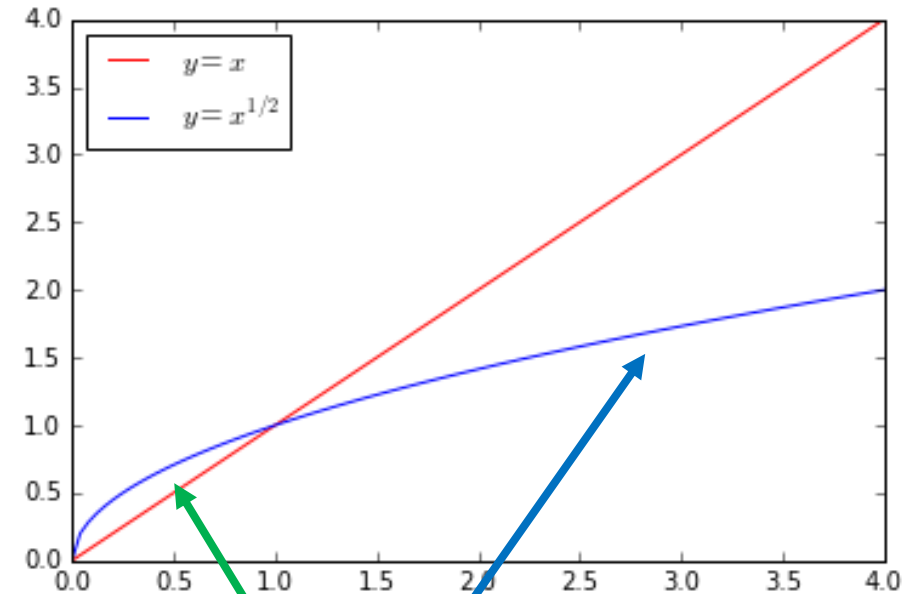
- RMSprop

- Adam [Kingma2014]

# Adagrad [Duchi2011]

○ Schedule

  ◦ $r_j = \sum_\tau (\nabla_\theta \mathcal{L}_j)^2 \implies \theta^{(t+1)} = \theta^{(t)} - \eta_t \frac{\nabla_\theta \mathcal{L}}{\sqrt{r} + \varepsilon}$

  ◦ $\varepsilon$ is a small number to avoid division with 0

  ◦ Gradients become gradually smaller and smaller

# RMSprop

Schedule

$$r = \alpha \sum_{\tau=1}^{t-1} (\nabla_\theta^{(t)} \mathcal{L}_j)^2 + (1-\alpha)\nabla_\theta^{(t)} \mathcal{L}_j \implies$$

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \frac{\nabla_\theta \mathcal{L}}{\sqrt{r}+\varepsilon}$$

*Decay hyper-parameter*

Moving average of the squared gradients
  ◦ Compared to Adagrad

**Large gradients**, e.g. too "noisy" loss surface
  ◦ Updates are tamed

**Small gradients**, e.g. stuck in flat loss surface ravine
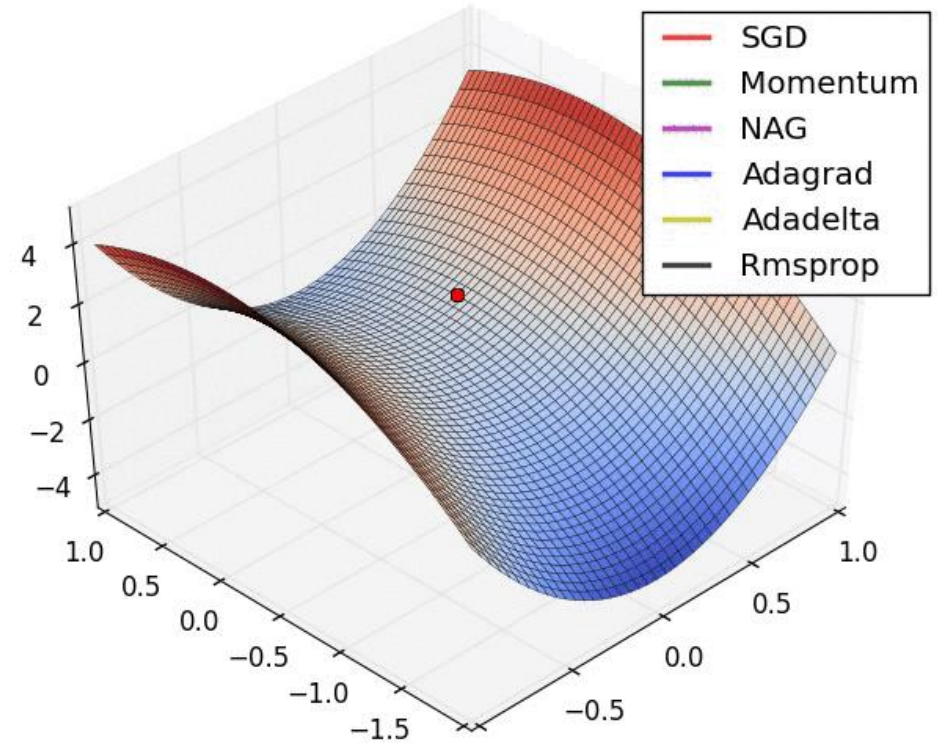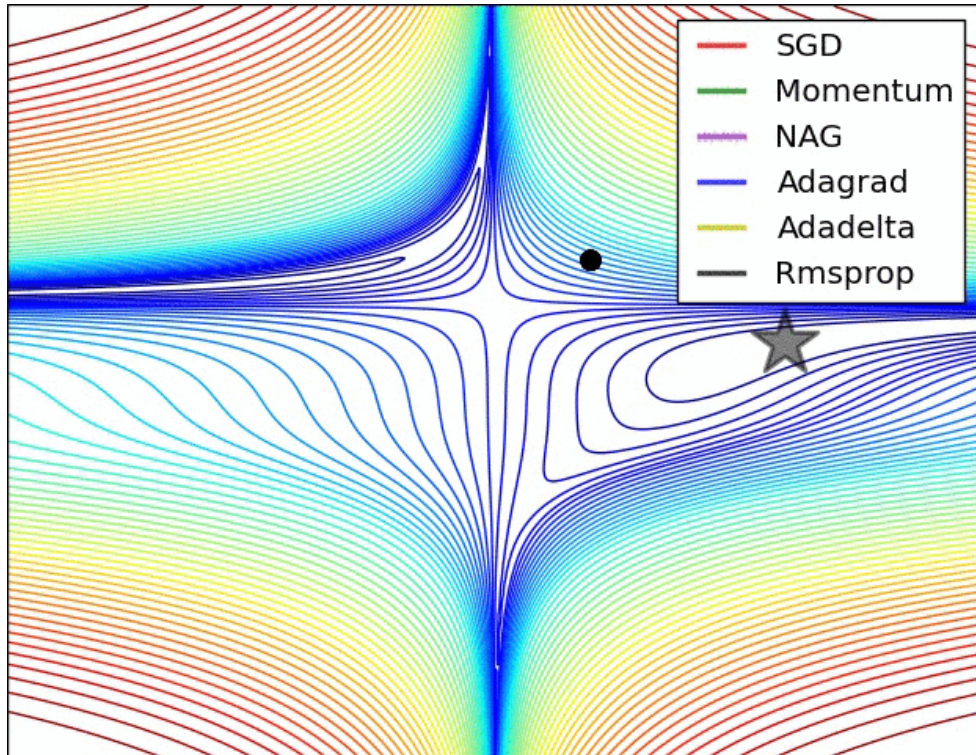  ◦ Updates become more aggressive



*Square rooting boosts small values while suppresses large values*

# Adam [Kingma2014]

- One of the most popular learning algorithms

$$g_t = \nabla_\theta \mathcal{L}$$

$$m_t = \beta_1 m_{t-1} + (1 - \beta_1) g_t$$

$$v_t = \beta_2 v_{t-1} + (1 - \beta_2) g_t^2$$

$$\widehat{m}_t = \frac{m_t}{1 - \beta_1^t}, \widehat{v}_t = \frac{v_t}{1 - \beta_2^t}$$

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \frac{\widehat{m}^{(t)}}{\sqrt{\widehat{v}^{(t)}} + \varepsilon}$$

  ◦ Recommended values: $\beta_1 = 0.9, \beta_2 = 0.999, \varepsilon = 10^{-8}$

- Similar to RMSprop, but with momentum & correction bias

# Visual overview



Picture credit: *Alec Radford*

# Another visualization

- https://habrahabr.ru/post/318970/

# Learning –not computing– the gradients

- Learning to learn by gradient descent by gradient descent
  - [Andrychowicz2016]

- $\theta^{(t+1)} = \theta^{(t)} + g_t(\nabla_\theta \mathcal{L}, \varphi)$

- $g_t$ is an "optimizer" with its own parameters $\varphi$
  - Implemented as a recurrent network

# Good practice

o Preprocess the data to at least have 0 mean

o Initialize weights based on activations functions
  ◦ For ReLU Xavier or HeICCV2015 initialization

o Always use $\ell_2$-regularization and dropout

o Use batch normalization

# Babysitting Deep Nets

1. The Neural Network

$$a_L\left(x; \theta_{1,\ldots,L}\right) = h_L\left(h_{L-1}(\ldots h_1(x, \theta_1), \theta_{L-1}), \theta_L\right)$$

2. Learning by minimizing empirical error

$$\theta^* \leftarrow \arg\min_\theta \sum_{(x,y) \subseteq (X,Y)} \mathcal{L}(y, a_L\left(x; \theta_{1,\ldots,L}\right))$$

3. Optimizing with Gradient Descent based methods

$$\theta^{(t+1)} = \theta^{(t)} - \eta_t \nabla_\theta \mathcal{L}$$

# Babysitting Deep Nets

o  Always check your gradients if not computed automatically

o  Check that in the first round you get a random loss

o  Check network with few samples
  ◦ Turn off regularization. You should predictably overfit and have a 0 loss
  ◦ Turn or regularization. The loss should increase

o  Have a separate validation set
  ◦ Compare the curve between training and validation sets
  ◦ There should be a gap, but not too large

# Summary

o How to define our model and optimize it in practice

o Data preprocessing and normalization

o Optimization methods

o Regularizations

o Architectures and architectural hyper-parameters

o Learning rate

o Weight initializations

o Good practices

# Reading material & references

- Chapter 8, 11

# Next lecture

- What are the Convolutional Neural Networks?

- Why are they important in Computer Vision?

- Differences from standard Neural Networks

- How to train a Convolutional Neural Network?