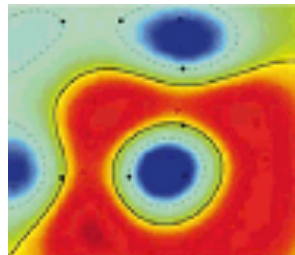




WiSe 2022/23

Deep Learning 1



Lecture 3

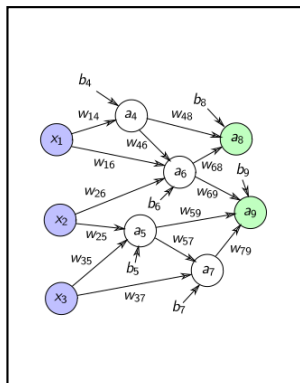
Optimization

Recap: How to Learn in a Neural Network

Observation:

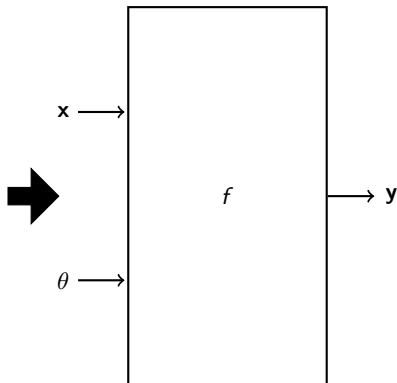
A neural network is a function of both its inputs and parameters.

graph view



$$\mathbf{x} = [x_1, x_2, x_3] \quad \mathbf{y} = [a_8, a_9]$$

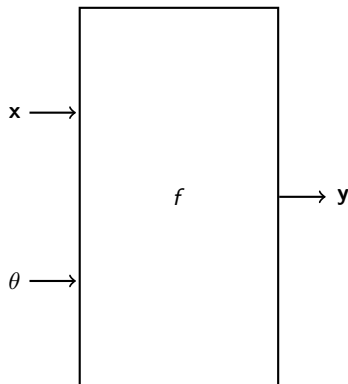
function view



$$\theta = (w_{ij})_{ij}, (b_j)_j$$

Recap: How to Learn in a Neural Network

function view



$$\theta = (w_{ij})_{ij}, (b_j)_j$$

Define an error function

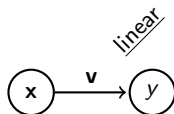
$$E(\theta) = \sum_n (f(\mathbf{x}_n; \theta) - t_n)^2$$

and minimize by gradient descend

$$\theta \leftarrow \theta - \gamma \cdot \nabla_{\theta} E(\theta)$$

- Characterizing the error function
 - Local minima and plateaus
 - Local curvature and condition number
- Improving optimization
 - Initialization
 - Choice of non-linearities
 - Momentum
- Fast Implementations

Characterizing the Error Function: One Layer



$$y = \mathbf{x}^\top \mathbf{v}$$

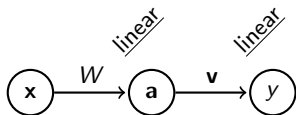
$$E = \frac{1}{N} \sum_n \|y_n - t_n\|^2 + \lambda \|\mathbf{v}\|^2 \quad (1)$$

$$= \frac{1}{N} \sum_n (\mathbf{x}_n^\top \mathbf{v} - t_n)^\top (\mathbf{x}_n^\top \mathbf{v} - t_n) + \lambda \|\mathbf{v}\|^2 \quad (2)$$

$$= \mathbf{v}^\top \underbrace{\left[\frac{1}{N} \sum_n \mathbf{x}_n \mathbf{x}_n^\top + \lambda I \right]}_{\text{positive semi-definite}} \mathbf{v} + \text{linear} + \text{constant} \quad (3)$$

 convex

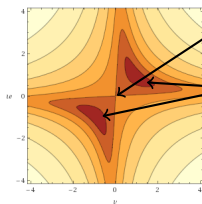
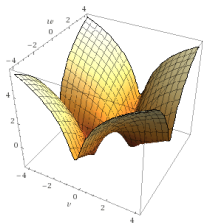
Characterizing the Error Function: Two Layers



$$y = \mathbf{v}^\top (W\mathbf{x})$$

$$E = \frac{1}{N} \sum_n \|\mathbf{v}^\top W\mathbf{x} - t_n\|^2 + \lambda(\|\mathbf{v}\|^2 + \|W\|^2) \quad (4)$$

Example: $N = 1$, $x_1 = 1$, $t_1 = 1$, $\lambda = 0.1$



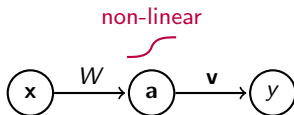
saddle point

two local minima



non-convex

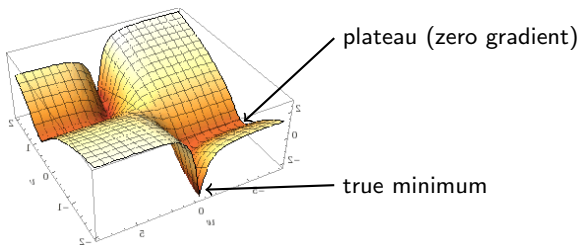
Characterizing the Error Function: Two Layers



$$y = \mathbf{v}^\top \tanh(W\mathbf{x})$$

$$E = \frac{1}{N} \sum_n \|\mathbf{v}^\top \tanh(W\mathbf{x}) - t_n\|^2 + \lambda(\|\mathbf{v}\|^2 + \|W\|^2) \quad (5)$$

Example: $N = 1$, $x_1 = 1$, $t_1 = 1$, $\lambda = 0.1$



Initializing the Neural Network

Even for the simplest two-layer neural network, the error function is already non-convex. Therefore, initialization of the neural network is important.

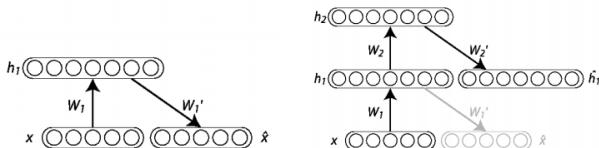
Recommendation for neural networks with tanh non-linearities: Scale parameters such that neuron outputs have variance ≈ 1 initially (LeCun'98/12 “Efficient Backprop”)

$$\theta \sim \mathcal{N}(0, \sigma^2) \qquad \sigma^2 = \frac{1}{\#\text{input neurons}} \qquad (6)$$

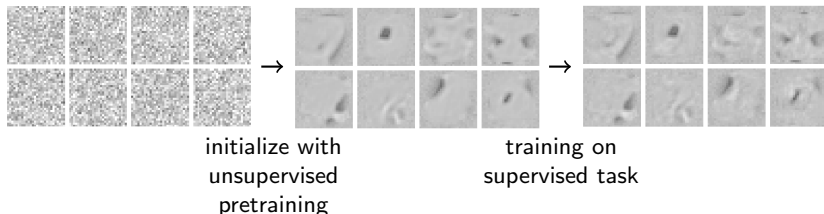
This initialization avoids saddle nodes and plateaus and also works well for ReLU non-linearities.

Initializing the Neural Network

Technique to reach good optima of the error function: layer-wise unsupervised pre-training (Hinton'06, Bengio'06, Vincent'08).



Example: Learning first-layer parameters on MNIST handwritten digits:

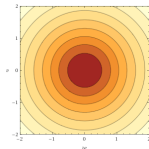
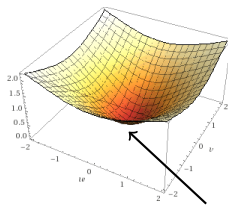


Is Convexity Sufficient?

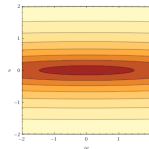
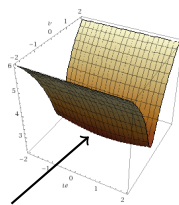
Answer: No. We must also verify that the function is well-conditioned.

Examples:

well-conditioned
error function



poorly conditioned
error function



minima of
the function

Well-conditioned functions
are easier to optimize.

Quantifying “well-conditioned”

The error function of a neural network can be approximated locally by a quadratic function.

$$E(\theta) \approx E(\theta_0) + \underbrace{\frac{\partial E}{\partial \theta} \Big|_{\theta_0}}_{\text{Gradient}} \cdot (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^\top \underbrace{\frac{\partial^2 E}{\partial \theta \partial \theta^\top} \Big|_{\theta_0}}_{\text{Hessian}} \cdot (\theta - \theta_0) \quad (7)$$

g **Gradient**
contains slope
information

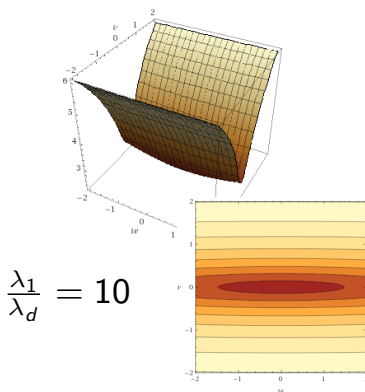
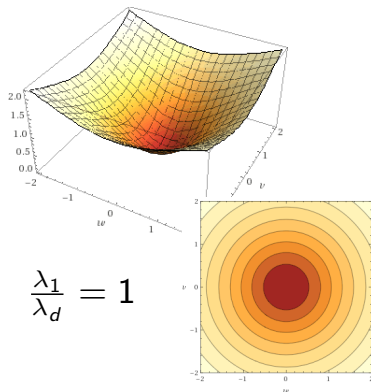
H **Hessian**
contains curva-
ture information

Idea: Look at the disbalance of curvature between different directions in the input space, by computing a ratio of eigenvalues.

$$\lambda_1, \lambda_2, \dots, \lambda_d = \text{eigval}(H)$$

$$\text{condition number} = \frac{\lambda_1}{\lambda_d}$$

Quantifying “well-conditioned”



The lower the condition number, the better.

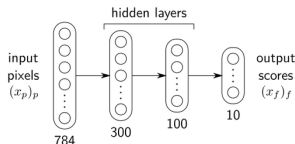
Computing the Hessian in Practice?

$$E(\theta) = E(\theta_0) + \underbrace{\frac{\partial E}{\partial \theta}}_{\text{Gradient}} \Big|_{\theta_0} \cdot (\theta - \theta_0) + \frac{1}{2} (\theta - \theta_0)^\top \underbrace{\frac{\partial^2 E}{\partial \theta \partial \theta^\top}}_{\text{Hessian}} \Big|_{\theta_0} \cdot (\theta - \theta_0) \quad (8)$$

Gradient
can be computed
with back-prop

Hessian
hard to compute and
very large for fully
connected networks

Example:

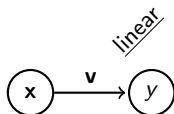


$$\begin{aligned} |\theta| &= 784 \cdot 300 + 300 \cdot 100 + 100 \cdot 10 \\ &= 266200 \\ |H| &= 266200 \cdot 266200 \\ &= 7.086 \cdot 10^{10} \sim 283GB@32bit \end{aligned}$$

For most practical tasks, we don't need to evaluate the Hessian and the condition number. We only need to apply a set of recommendations and tricks that keep the condition number low.

Improving Conditioning of the Error Function

Example: The linear model



$$y = \mathbf{x}^\top \mathbf{v}$$

$$E = \mathbf{v}^\top \underbrace{\left[\frac{1}{N} \sum_n \mathbf{x}_n \mathbf{x}_n^\top + \lambda I \right]}_H \mathbf{v} + \text{linear} + \text{constant} \quad (9)$$

H Hessian of the error function

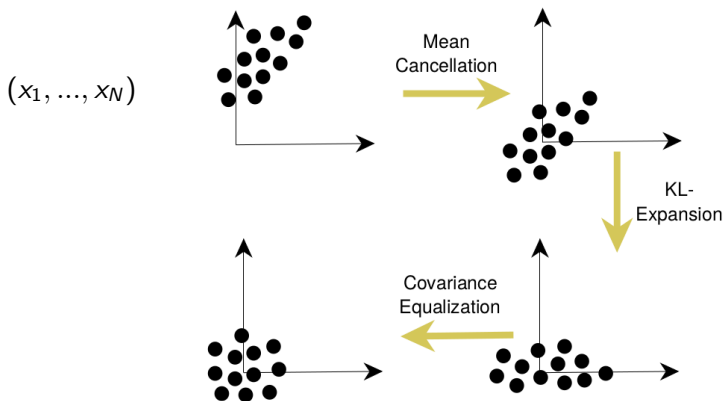
Hessian (and condition number) are influenced by the mean and covariance of the data. The closer the mean is to zero, and the closer the covariance is to the identity, the lower the condition number.



Trick: Normalize the data

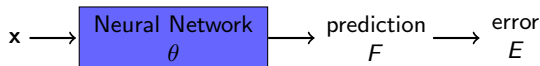
Data Normalization to Improve Conditioning

Data pre-processing *before* training:



(image from LeCun'98/12)

Approximation of the Hessian



General formula for the Hessian of a neural network (size: $|\theta| \times |\theta|$)

$$H = \frac{\partial^2 E}{\partial \theta^2} = \frac{\partial F^\top}{\partial \theta} \frac{\partial^2 E}{\partial F^2} \frac{\partial F}{\partial \theta} + \frac{\partial E}{\partial F} \frac{\partial^2 F}{\partial \theta^2}$$

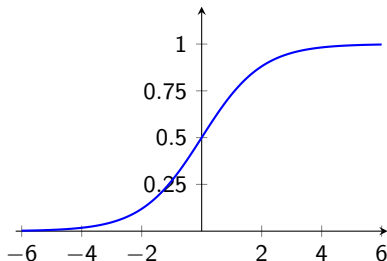
Hessian between weights of a single neuron (mean square error case):

$$[H_k]_{jj'} = \frac{\partial^2 E}{\partial w_{jk} \partial w_{j'k}} = \mathbb{E} \left[\underbrace{a_j a_{j'} \delta_k^2}_{\substack{\text{similar to} \\ \text{the simple} \\ \text{linear model}}} \right] + \mathbb{E} \left[\underbrace{a_j \cdot \frac{\partial \delta_k}{\partial w_{j'k}} \cdot (y - t)}_{\text{complicated}} \right]$$

Improving Conditioning of Higher-Layers

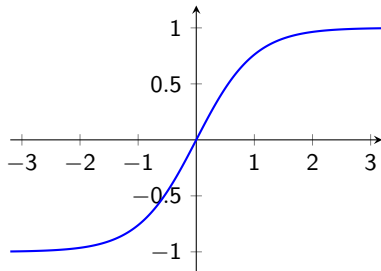
To improve conditioning, not only the input data should be normalized, but also the representations built from this data at each layer. This can be done by carefully choosing the activation function.

logistic sigmoid



activations are
not centered
→ **high condition number**

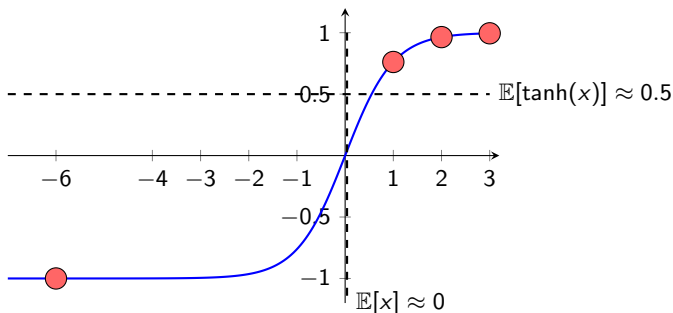
hyperbolic tangent



activations ap-
prox. centered at 0
→ **low condition number**

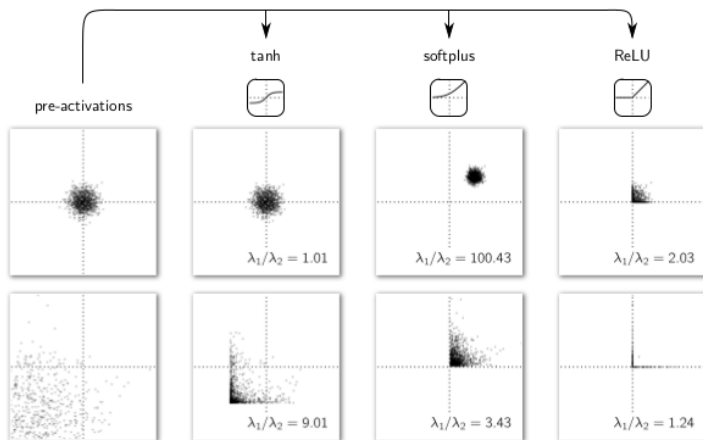
Limitation of tanh

The tanh non-linearity works well initially, but after some training steps, it might no longer work as expected as the input distribution will drift to negative or positive values.



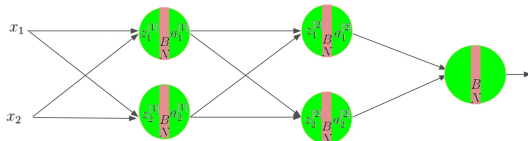
Remark: If the input of tanh is centered but skewed, the output of tanh will not be centered. This happens a lot in practice, e.g. when the problem representation needs to be sparse.

Comparing Non-Linearities



Further Improving the Hessian

Batch Normalization [Ioffe'15]



$$\begin{aligned}
 z^{[l]} &= W^{[l]} a^{[l-1]} \longrightarrow \mu^{[l]} = \frac{1}{m} \sum_i z^{[l](i)} \\
 \sigma^{[l]2} &= \frac{1}{m} \sum_i (z^{[l](i)} - \mu^{[l]})^2 \\
 z_{norm}^{[l](i)} &= \frac{z^{[l](i)} - \mu^{[l]}}{\sqrt{\sigma^{[l]2} + \epsilon}} \\
 \tilde{z}^{[l](i)} &= \gamma^{[l]} z_{norm}^{[l](i)} + \beta^{[l]} \longrightarrow a^{[l]} = g^{[l]}(\tilde{z}^{[l]})
 \end{aligned}$$

Advantage: Decorrelates gradients between different layers. \rightarrow Improve cross-layer parts of the Hessian.

Observation: Useful to train very deep networks efficiently.



(Stochastic) Gradient Descent

Objective to Minimize:

$$E(\theta) = \frac{1}{N} \sum_{n=1}^N E^{(n)}(\theta)$$

Batch GD:

while True:

$$\theta \leftarrow \theta - \gamma \frac{\partial}{\partial \theta} \frac{1}{N} \sum_{n=1}^N E^{(n)}$$

$O(N)$

Stochastic GD (SGD):

while True:

$$n \leftarrow \text{random}(1, N)$$

$$\theta \leftarrow \theta - \gamma \frac{\partial E^{(n)}}{\partial \theta}$$

$O(1)$

Batch GD:

while True:

$$\theta \leftarrow \theta - \gamma \frac{\partial}{\partial \theta} \frac{1}{N} \sum_{n=1}^N E^{(n)}$$

Stochastic GD (SGD):

while True:

$$n \leftarrow \text{random}(1, N)$$

$$\theta \leftarrow \theta - \gamma \frac{\partial E^{(n)}}{\partial \theta}$$



use a decreasing schedule

$$\gamma^{(1)}, \gamma^{(2)}, \dots, \gamma^{(T)}$$

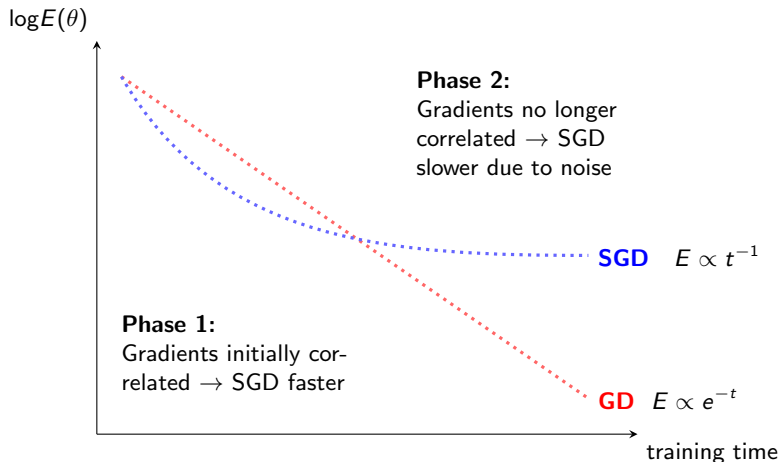
Conditions for SGD convergence:

$$\sum_{t=1}^{\infty} \gamma^{(t)} \rightarrow \infty \quad \lim_{t \rightarrow \infty} \gamma^{(t)} = 0$$

Which SGD learning schedule to choose?

	$\gamma^{(t)} = 1$	$\gamma^{(t)} = t^{(-1)}$	$\gamma^{(t)} = e^{-t}$
$\sum_{t=1}^{\infty} \gamma^{(t)} \rightarrow \infty$	✓	✓	✗
$\lim_{t \rightarrow \infty} \gamma^{(t)} = 0$	✗	✓	✓

GD vs. SGD Convergence



Insight: Phase 2 is not relevant, because the model already starts overfitting before reaching it. \rightarrow SGD is the method of choice for most practical purposes.

Momentum

Idea: Choose the update direction as a weighted average of previous updates.

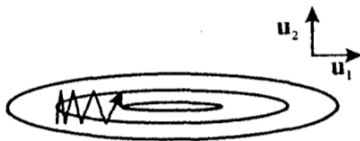


Image from Bishop'95

Accelerates convergence along direction of low curvature.

Momentum can help to overcome a poorly conditioned neural network.

Momentum

Update the direction of descent as:

$$\Delta^{(t)} = \mu \cdot \Delta^{(t-1)} + \gamma \cdot \nabla E(\theta^{(t)})$$

Diagram illustrating the update of the direction of descent $\Delta^{(t)}$ using momentum and the error gradient.

The equation is: $\Delta^{(t)} = \mu \cdot \Delta^{(t-1)} + \gamma \cdot \nabla E(\theta^{(t)})$

Labels and arrows indicating the components:

- $\Delta^{(t)}$: new update step
- μ : momentum
- $\Delta^{(t-1)}$: previous update step
- γ : learning rate
- $\nabla E(\theta^{(t)})$: error gradient

and update the neural network parameters following this direction:

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

The Adam Algorithm

Algorithm 1: *Adam*, our proposed algorithm for stochastic optimization. See section 2 for details, and for a slightly more efficient (but less clear) order of computation. g_t^2 indicates the elementwise square $g_t \odot g_t$. Good default settings for the tested machine learning problems are $\alpha = 0.001$, $\beta_1 = 0.9$, $\beta_2 = 0.999$ and $\epsilon = 10^{-8}$. All operations on vectors are element-wise. With β_1^t and β_2^t we denote β_1 and β_2 to the power t .

Require: α : Stepsize

Require: $\beta_1, \beta_2 \in [0, 1)$: Exponential decay rates for the moment estimates

Require: $f(\theta)$: Stochastic objective function with parameters θ

Require: θ_0 : Initial parameter vector

$m_0 \leftarrow 0$ (Initialize 1st moment vector)

$v_0 \leftarrow 0$ (Initialize 2nd moment vector)

$t \leftarrow 0$ (Initialize timestep)

while θ_t not converged **do**

$t \leftarrow t + 1$

$g_t \leftarrow \nabla_{\theta} f_t(\theta_{t-1})$ (Get gradients w.r.t. stochastic objective at timestep t)

$m_t \leftarrow \beta_1 \cdot m_{t-1} + (1 - \beta_1) \cdot g_t$ (Update biased first moment estimate)

$v_t \leftarrow \beta_2 \cdot v_{t-1} + (1 - \beta_2) \cdot g_t^2$ (Update biased second raw moment estimate)

$\hat{m}_t \leftarrow m_t / (1 - \beta_1^t)$ (Compute bias-corrected first moment estimate)

$\hat{v}_t \leftarrow v_t / (1 - \beta_2^t)$ (Compute bias-corrected second raw moment estimate)

$\theta_t \leftarrow \theta_{t-1} - \alpha \cdot \hat{m}_t / (\sqrt{\hat{v}_t} + \epsilon)$ (Update parameters)

end while

return θ_t (Resulting parameters)

from Kingma'15

Neural Network Training Time

SGD

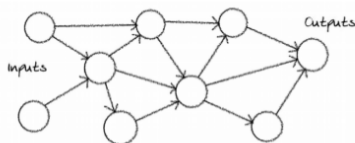


$$O(\#iterations \times \#edges)$$



Backprop

This network: 13 connections



can probably be trained with a few
hundreds iterations.

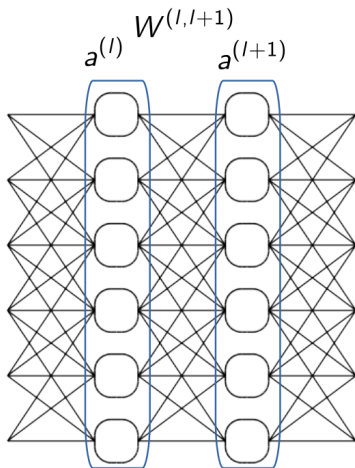
Googlenet: $> 10^9$ connections



need millions of iterations
to be trained.

Part 2: Scaling Deep Nets to Big Data

Step 1: Systemize Computations



Per-neuron forward computations

$$\forall_j : a_j = g\left(\sum_i a_i w_{ij} + b_j\right)$$

Whole-layer computation

$$a^{(l+1)} = g\left(\underbrace{W^{(l,l+1)} \cdot a^{(l)}}_{\text{matrix-vector products (e.g. numpy.dot)}} + b^{(l+1)}\right)$$

matrix-vector
products (e.g.
numpy.dot)

element-wise
application of
nonlinearity

Step 2: Mini-Batches

Idea: Take advantage of fast matrix-matrix multiplications by feeding several examples at a time to the neural network.

Example for layer with weight matrix: $W^{(l)} \in \mathbb{R}^{d \times d}$

Pure SGD (propagate 1 data point)

$$\mathbf{a}^{(l)}, \mathbf{a}^{(l+1)} \in \mathbb{R}^d \quad \mathbf{a}^{(l+1)} = g(W^{(l)} \mathbf{a}^{(l)} + \mathbf{b}^{(l)})$$

$O(d^2)$ computations

Mini-batch SGD (propagate m data points)

$$A^{(l)}, A^{(l+1)} \in \mathbb{R}^{d \times m} \quad A^{(l+1)} = g(W^{(l)} A^{(l)} + b^{(l)})$$

$$m \approx d \Rightarrow O(d^{2.4}) < O(md^2) \text{ computations}$$

Step 2: How to Choose Mini-Batch Size?


$$m \approx d$$

Advantages:

- (1) Largest speed up in terms of matrix multiplications

Disadvantages:

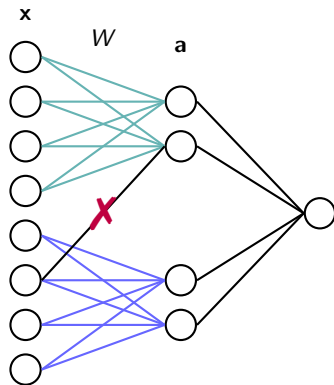
- (1) Multiplication might not fit in memory
- (2) If gradients are correlated, same direction as pure SGD
- (3) Layers have different size

Step 2: How to Choose Mini-Batch Size?

	pure SGD	medium mini-batch	largest batch ($M = N$) = GD
avoid redundant gradient computations	✓	✓	✗
speed-up from matrix-matrix multiplications	✗	✓	✓

Step 3: Prune Irrelevant Computations

Example: Prune long-range interactions in images or text, or other types of sequential data.



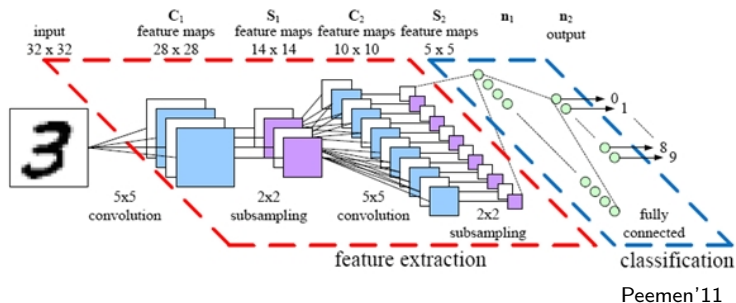
$$\mathbf{a} = \underbrace{W^T \mathbf{x}} = \underbrace{W_A^T \mathbf{x}_A + W_B^T \mathbf{x}_B}$$

$$8 \times 4 = 32 \\ \text{computations}$$

$$2 \times (2 \times 4) = 16 \\ \text{computations}$$

$$W = \begin{bmatrix} W_A & 0 \\ 0 & W_B \end{bmatrix}$$

Step 3: Avoid Computational Bottlenecks



- Lower layers detect simple features at exact locations.
- Higher layers detect complex features at approximate locations.
- Layers progressively replace spatial information with semantic information
→ keep dimensionality and number of connections low at each layer.

Step 3: GoogLeNet Example

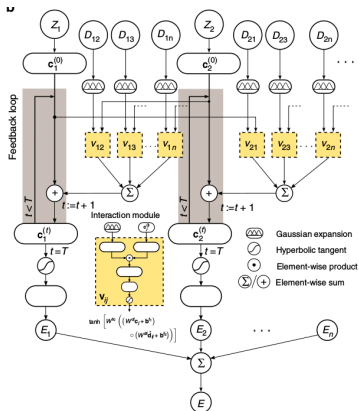
type	patch size/ stride	output size	output filters	depth	#1×1	#3×3 reduce	#3×3	#5×5 reduce	#5×5	pool proj	params	ops
convolution	7×7/2	112×112	64	1							2.7K	34M
max pool	3×3/2	56×56	64	0								
convolution	3×3/1	56×56	192	2		64	192				112K	360M
max pool	3×3/2	28×28	192	0								
inception (3a)		28×28	256	2	64	96	128	16	32	32	159K	128M
inception (3b)		28×28	480	2	128	128	192	32	96	64	380K	304M
max pool	3×3/2	14×14	480	0								
inception (4a)		14×14	512	2	192	96	208	16	48	64	364K	73M
inception (4b)		14×14	512	2	160	112	224	24	64	64	437K	88M
inception (4c)		14×14	512	2	128	128	256	24	64	64	463K	100M
inception (4d)		14×14	528	2	112	144	288	32	64	64	580K	119M
inception (4e)		14×14	832	2	256	160	320	32	128	128	840K	170M
max pool	3×3/2	7×7	832	0								
inception (5a)		7×7	832	2	256	160	320	32	128	128	1072K	54M
inception (5b)		7×7	1024	2	384	192	384	48	128	128	1388K	71M
avg pool	7×7/1	1×1	1024	0								
dropout (40%)		1×1	1024	0								
linear		1×1	1000	1							1000K	1M
softmax		1×1	1000	0								

Szegedy '14

GoogLeNet incarnation of the Inception architecture

Step 3: SchNet Example

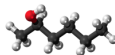
SchNet: neural network that predicts molecular properties, and where each layer models an exchange of local information in the molecular graph [Schütt'17].



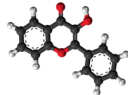
GDB-3:
 - max 3 heavy atoms
 - few of them



GDB-7:
 - max 7 heavy atoms
 - thousands of them

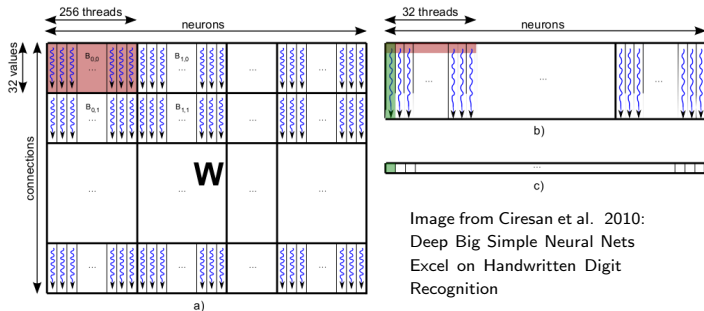


GDB-17:
 - max 17 heavy atoms
 - billions of them



SchNet handles arbitrary large molecules by modeling local atomic interaction in a small neighborhood, thereby keeping the dimensionality low.

Step 4: Map Neural Network to Hardware

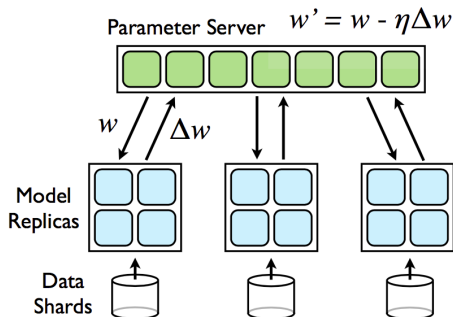


In order for the training procedure to match the hardware specifications (e.g. CPU cache, GPU block size) optimally, neural network computations (e.g. matrix multiplications) must be decomposed into blocks of appropriate size.

These hardware-specific optimizations are already built in most fast neural network libraries (e.g. PyTorch, Tensorflow, cuDNN, ...).

Step 5: Distributed Training

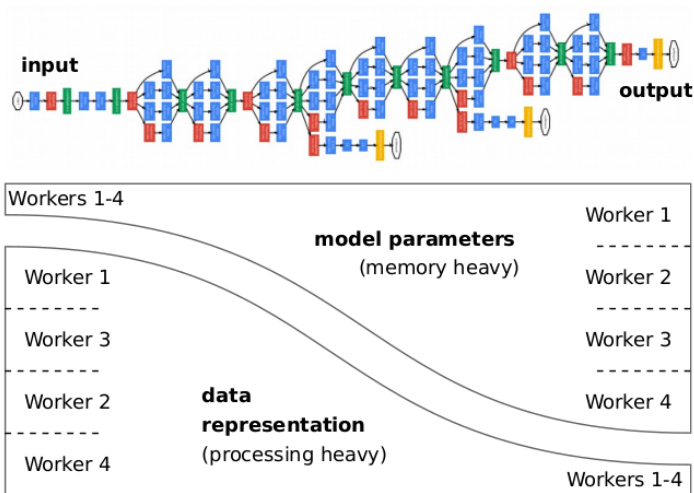
Example: Google's DistBelief Architecture [Dean'12]



Each model replica trains on its own data, and synchronizes the model parameters it has learned with other replica via a dedicated parameter server.

Step 5: Distributed Training

Combining data-parallelism and model-parallelism



see also Krizhevsky'14: One weird trick for parallelizing convolutional neural networks

Summary

- Optimization of neural networks is harder than linear models, because of (1) non-convexity and (2) conditioning issues.
- Therefore, we must carefully choose the initialization, the structure of the model (e.g. non-linearities), and if necessary add momentum to gradient descent.
- Neural networks principal goal is to enable the transformation of large datasets into complex highly predictive models.
- To achieve this, it is important to make sure they can be trained as quickly as possible (e.g., mini-batches, layered structure, avoiding bottlenecks, matching the hardware, distributed architectures).