

# Deep Learning

*Lecture 3 – Optimization: Stochastic gradient descent and backpropagation*



UPPSALA  
UNIVERSITET

**Thomas Schön**

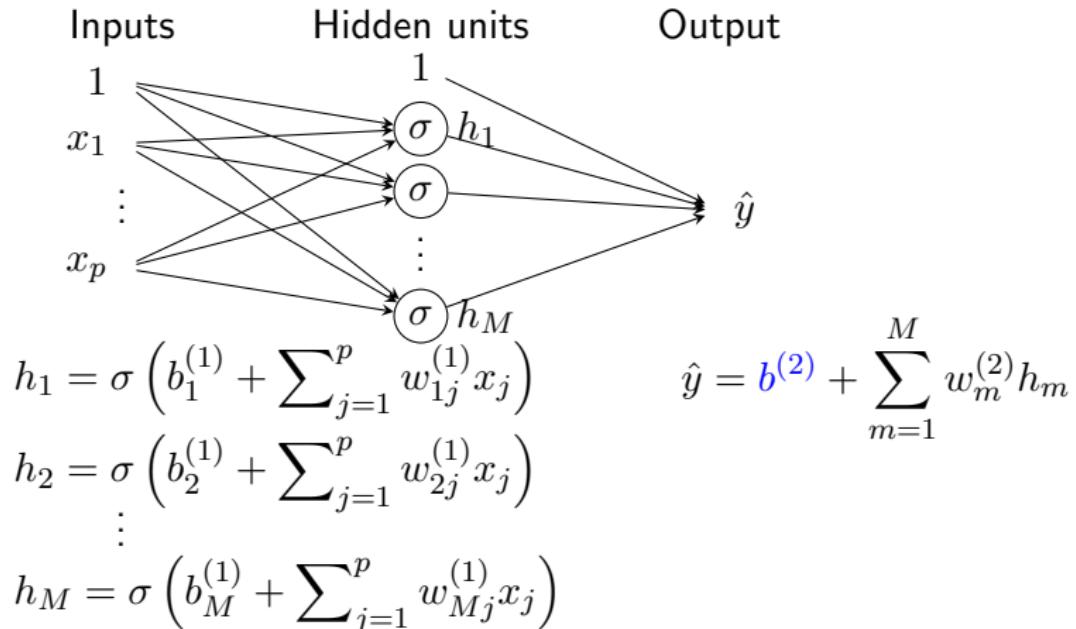
Department of Information Technology  
Uppsala University

[thomas.schon@it.uu.se](mailto:thomas.schon@it.uu.se)

[user.it.uu.se/~thosc112](http://user.it.uu.se/~thosc112)

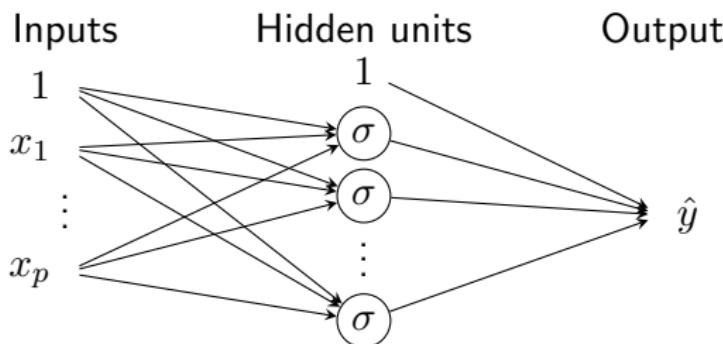
# Summary Lecture 2 (I/II) - Neural network

A neural network is a sequential construction of **several** generalized linear regression models.



# Summary Lecture 2 (I/II) - Neural network

A neural network is a sequential construction of **several** generalized linear regression models.



$$\mathbf{h} = \sigma(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})$$

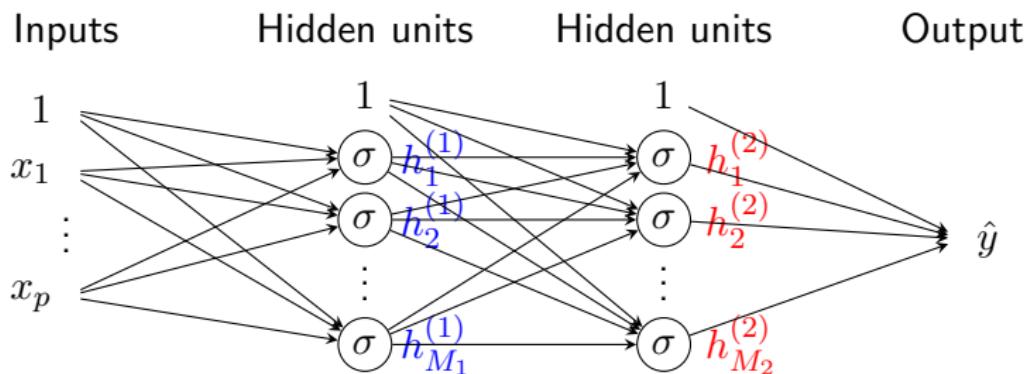
$$\hat{y} = \mathbf{W}^{(2)} \mathbf{h} + \mathbf{b}^{(2)}$$

$$\mathbf{W}^{(1)} = \begin{bmatrix} w_{11}^{(1)} & \dots & w_{1p}^{(1)} \\ \vdots & & \vdots \\ w_{M1}^{(1)} & \dots & w_{Mp}^{(1)} \end{bmatrix}, \quad \mathbf{b}^{(1)} = \begin{bmatrix} b_1^{(1)} \\ \vdots \\ b_M^{(1)} \end{bmatrix}, \quad \mathbf{h} = \begin{bmatrix} h_1 \\ \vdots \\ h_M \end{bmatrix} \quad \mathbf{b}^{(2)} = [b^{(2)}] \\ \mathbf{W}^{(2)} = [w_1^{(2)} \dots w_M^{(2)}]$$

**Weight matrix**      **Offset vector**      **Hidden units**

# Summary Lecture 2 (I/II) - Neural network

A neural network is a **sequential** construction of several generalized linear regression models.



$$\mathbf{h}^{(1)} = \sigma(\mathbf{W}^{(1)} \mathbf{x} + \mathbf{b}^{(1)})$$

$$\mathbf{h}^{(2)} = \sigma(\mathbf{W}^{(2)} \mathbf{h}^{(1)} + \mathbf{b}^{(2)})$$

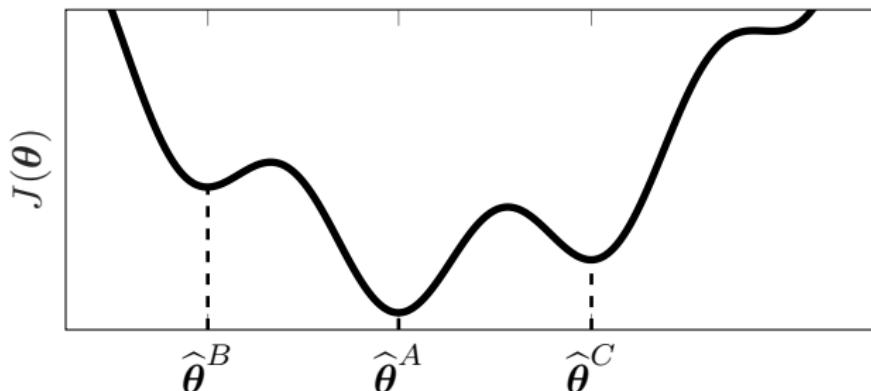
$$\hat{y} = \mathbf{W}^{(3)} \mathbf{h}^{(2)} + \mathbf{b}^{(3)}$$

The model tends to learn better using a deep network (several layers) instead of a wide and shallow network.

# Summary of Lecture 2 (II/II) - Optimization

We train our models by considering an optimization problem

$$\hat{\theta} = \arg \min_{\theta} J(\theta), \quad J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\mathbf{x}_i, \mathbf{y}_i, \theta)$$



- The best possible solution  $\hat{\theta}$  is the **global minimizer** ( $\hat{\theta}^A$ )
- The global minimizer is typically very hard to find, and we have to settle for a **local minimizer** ( $\hat{\theta}^A, \hat{\theta}^B, \hat{\theta}^C$ )

# Outline – Lecture 3

---

**Aim:** Motivate and introduce the Stochastic Gradient (SG) algorithm and derive the backpropagation algorithm to efficiently compute the gradient needed for SG.

## Outline:

1. Summary – lecture 2
2. Optimization in machine learning
3. Stochastic gradient
4. Backpropagation

# Reminder — forming the loss function

---

The model provides a prediction  $\hat{y}_\theta(x) = f(x, \theta)$  of the output  $y$ .

For a good model we would expect  $\hat{y}_\theta(x) \approx y$

Hence, it is natural to look for parameters  $\theta$  that makes  $\hat{y}_\theta(x)$  as close to  $y$  as possible.

---

Results in an optimization problem of the form

$$\min_{\theta} J(\theta)$$

where

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n L(\hat{y}_\theta(x_i), y_i)$$

# Computational challenge 1 – $n$ is large

At each optimization step we need to compute the gradient

$$\mathbf{g}_t = \nabla_{\boldsymbol{\theta}} J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}).$$

**Computational challenge 1 –  $n$  large:** We typically have large-scale training data ( $n \gg 1$ ) available, which implies that it is costly to compute the gradient.

**Solution: Randomization**

# Computational challenge 1 – $n$ is large

---

At each iteration, we only use a small part of the dataset to compute an estimate of the gradient  $\mathbf{g}_k$ . This is called **minibatch stochastic gradient**.

Large scale machine learning requires **stochastic optimization**.

**Randomization** often produce reliability (in large scale problems).

Massive research on large scale optimization over the past decade.

# Computational challenge 2 – $\dim(\theta)$ large

---

At each optimization step we need to compute the gradient

$$\mathbf{g}_t = \nabla_{\theta} J(\theta) = \frac{1}{n} \sum_{i=1}^n \nabla_{\theta} L_i(\theta).$$

**Computational challenge 2 –  $\dim(\theta)$  large:** A neural network contains a lot of parameters. Computing the gradient is costly.

**Solution:** A deep NN is a composition of multiple layers. Hence, each term  $\nabla_{\theta} L_i(\theta)$  can be computed efficiently by repeatedly applying the chain rule and storing certain intermediate results for later reuse. This is called **backpropagation**.

# Stochastic gradient (SG)

# The DL optimization problem

$$J(\boldsymbol{\theta}) = \frac{1}{n} \sum_{i=1}^n L(f(\mathbf{x}_i, \boldsymbol{\theta}), y_i) = \frac{1}{n} \sum_{i=1}^n L_i(\boldsymbol{\theta})$$

Basic characteristics:

1. Non-convex
2. Large-scale (both large  $n$  and large dimension of  $\boldsymbol{\theta}$ )

Two main classes of optimization methods (for this lecture):

1. Deterministic (batch)
2. Stochastic (online)

# Problem formulation – stochastic opt.

**What?** Solve the non-convex stochastic optimization problem

$$\min_{\theta} J(\theta)$$

when we only have access to **noisy** evaluations of  $J(\theta)$  and its derivatives.

**Why?** These stochastic optimization problems are common:

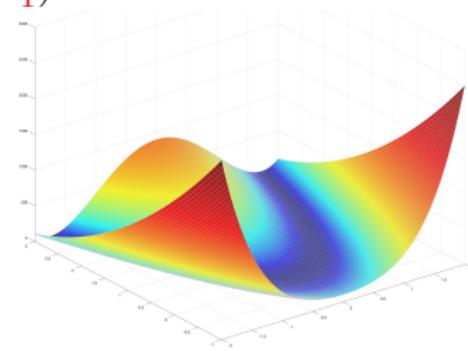
- When the cost function cannot be evaluated on the entire dataset. Large-scale ML and DL in particular.
- When numerical methods approximate  $J(\theta)$  and  $\nabla^i J(\theta)$ .
- ...

# Toy example – Rosenbrock's banana function

Let  $J(\theta) = (1 - \theta_1)^2 + 100(\theta_2 - \theta_1^2)^2$ .

## Deterministic problem

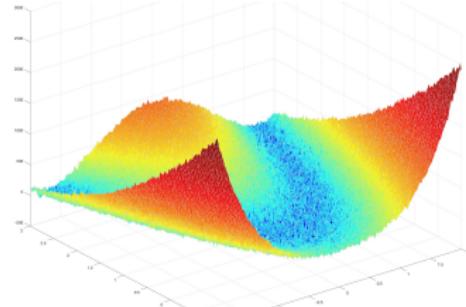
$$\min_{\theta} J(\theta)$$



## Stochastic problem

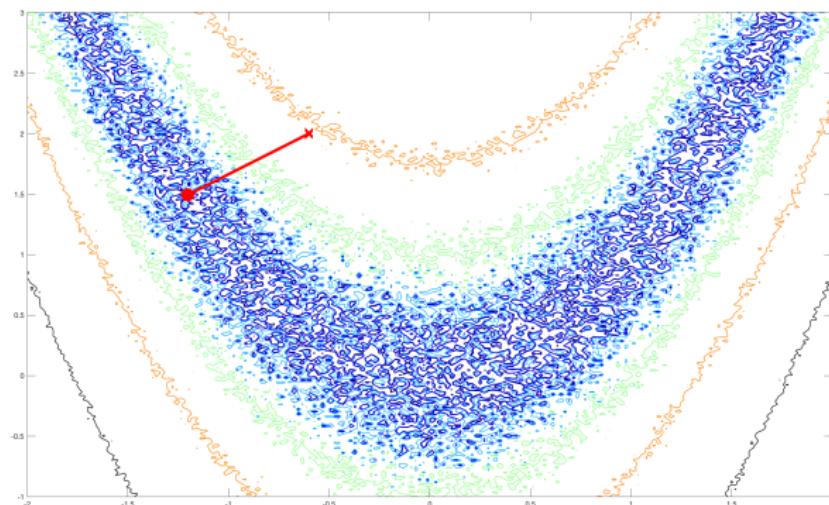
$$\min_{\theta} J(\theta)$$

when we only have access to noisy versions of the cost function  
 $(\tilde{J}(\theta) = J(\theta) + e, e = \mathcal{N}(0, 30^2))$   
and its noisy gradients.



# Std. deterministic optimization at work

---



Terminates at the wrong solution after 3 iterations.

The true solution is  $(1, 1)$ .

# SG is a Markov chain!

---

The SG idea was first proposed almost 70 years ago in

Robbins, H. and Monro, S. **A stochastic approximation method.** *The Annals of Mathematical Statistics.* 22(3):400–407, 1951.

They introduce SG as a particular **Markov chain**, which is since then 'forgotten'.

# Motivating the use of SG

---

Advantages of the batch approach:

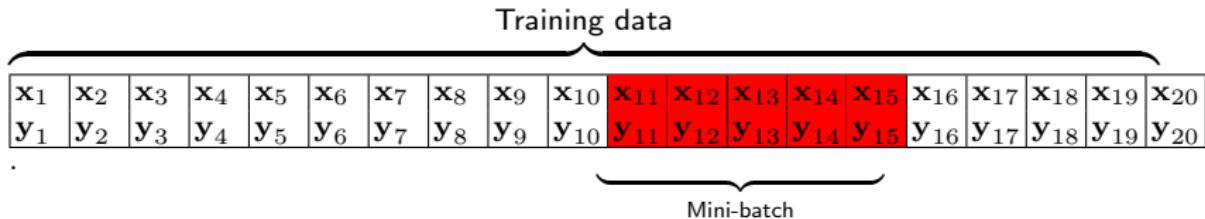
1. There are many good optimization methods available for the deterministic problem.
2. The sum-structure of  $J$  and  $\nabla L$  allows for computation in a distributed manner.

Why is the stochastic gradient (SG) algorithm still the standard for large-scale machine learning (including DL)?

1. Intuitive/practical motivation
2. Theoretical motivation

We will mainly use a middle ground (combining the best properties of batch and stochastic algorithms): **minibatch SG**.

# Minibatch SG



$$\theta_3 = \theta_2 - \gamma \frac{1}{5} \sum_{i=11}^{15} \nabla_{\theta} L(\mathbf{x}_i, \mathbf{y}_i, \theta_2)$$

Recall:

- The extreme version is to use only one data point at each training step (called **online learning**)
- We typically do something in between (not one data point, and not all data). We use a smaller set called **minibatch**.

One pass through the training data is called an **epoch**.

# Minibatch SG

x <sub>7</sub>	x <sub>10</sub>	x <sub>3</sub>	x <sub>20</sub>	x <sub>16</sub>	x <sub>2</sub>	x <sub>1</sub>	x <sub>18</sub>	x <sub>19</sub>	x <sub>12</sub>	x <sub>6</sub>	x <sub>11</sub>	x <sub>17</sub>	x <sub>15</sub>	x <sub>5</sub>	x <sub>14</sub>	x <sub>4</sub>	x <sub>9</sub>	x <sub>13</sub>	x <sub>8</sub>
y <sub>7</sub>	y <sub>10</sub>	y <sub>3</sub>	y <sub>20</sub>	y <sub>16</sub>	y <sub>2</sub>	y <sub>1</sub>	y <sub>18</sub>	y <sub>19</sub>	y <sub>12</sub>	y <sub>6</sub>	y <sub>11</sub>	y <sub>17</sub>	y <sub>15</sub>	y <sub>5</sub>	y <sub>14</sub>	y <sub>4</sub>	y <sub>9</sub>	y <sub>13</sub>	y <sub>8</sub>

**Iteration:** 3

**Epoch:** 1

- If we pick the minibatches in order, they might be unbalanced and not representative for the whole data set.
- Therefore, we pick data points **at random** from the training data to form a minibatch.
- One implementation is to randomly reshuffle the data before dividing it into minibatches.
- After each epoch we do another reshuffling and another pass through the data set.

# Minibatch SG

---

1. Initialize  $\boldsymbol{\theta}_0$ , set  $t \leftarrow 1$ , choose batch size  $n_b$  and number of epochs  $E$ .
  2. For  $i = 1$  to  $E$ 
    - (a) Randomly shuffle the training data  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=1}^n$ .
    - (b) For  $j = 1$  to  $\frac{n}{n_b}$ 
      - (i) Approximate the gradient of the loss function using the current minibatch  $\{\mathbf{x}_i, \mathbf{y}_i\}_{i=(j-1)n_b+1}^{jn_b}$ 
$$\hat{\mathbf{g}}_t = \frac{1}{n_b} \sum_{i=(j-1)n_b+1}^{jn_b} \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t} \approx \mathbf{g}_t.$$
      - (ii) Take a step in the gradient direction  $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \gamma \hat{\mathbf{g}}_t$ .
      - (iii) Update the iteration index  $t \leftarrow t + 1$ .
- 

At each iteration we get a stochastic approximation  $\hat{\mathbf{g}}_t$  of the true gradient

$$\mathbf{g}_t = \frac{1}{n} \sum_{i=1}^n \nabla_{\boldsymbol{\theta}} L_i(\boldsymbol{\theta}) \Big|_{\boldsymbol{\theta}=\boldsymbol{\theta}_t}$$

# Learning rate within minibatch SG

It is necessary to select a learning rate that decrease over time.  
A sufficient condition to guarantee convergence is that

$$\sum_{t=1}^{\infty} \gamma_t = \infty \quad \text{and} \quad \sum_{t=1}^{\infty} \gamma_t^2 < \infty.$$

Practically common approach: decay the learning rate linearly

$$\gamma_t = (1 - \alpha)\gamma_0 + \alpha\gamma_\tau,$$

with

$$\alpha = \frac{t}{\tau}$$

until iteration  $\tau$  and then leave  $\gamma$  constant after that.

# Can we do better?

---

The answer to the title question is **Yes!**

The next question is then: **How?**

By using more information in computing:

- the search direction  $d_t$
- and the learning rate  $\gamma_t$

$$\theta_{t+1} = \theta_t + \gamma_t d_t.$$

# Quasi-Newton – A non-standard take

Our problem is of the form

$$\min_{\theta} J(\theta)$$

**Idea underlying (quasi-)Newton methods:** Learn a local quadratic model  $q(\theta_t, \delta)$  of  $J(\theta)$  around the current iterate  $\theta_t$

$$q(\theta_t, \delta) = f(\theta_t) + g(\theta_t)^T \delta + \frac{1}{2} \delta^T H(\theta_t) \delta$$

$$g(\theta_t) = \nabla J(\theta)|_{\theta=\theta_t}, \quad H(\theta_t) = \nabla^2 J(\theta)|_{\theta=\theta_t}, \quad \delta = \theta - \theta_k.$$

We have measurements of

- the cost function  $J_t = J(\theta_t)$ ,
- and its gradient  $g_t = g(\theta_t)$ .

**Question:** How do we update the Hessian model?

# Useful basic facts

Line segment connecting two adjacent iterates  $\theta_t$  and  $\theta_{t+1}$ :

$$r_t(\tau) = \theta_t + \tau(\theta_{t+1} - \theta_t), \quad \tau \in [0, 1].$$

1. The **fundamental theorem of calculus** states that

$$\int_0^1 \frac{\partial}{\partial \tau} \nabla f(r_t(\tau)) d\tau = \nabla J(r_t(1)) - \nabla f(r_t(0)) = \underbrace{\nabla J(\theta_{t+1})}_{g_{t+1}} - \underbrace{\nabla J(\theta_t)}_{g_t}$$

2. The **chain rule** tells us that

$$\frac{\partial}{\partial \tau} \nabla J(r_t(\tau)) = \nabla^2 J(r_t(\tau)) \frac{\partial r_t(\tau)}{\partial \tau} = \nabla^2 J(r_t(\tau)) (\theta_{t+1} - \theta_t).$$

$$\underbrace{g_{t+1} - g_t}_{=y_t} = \int_0^1 \frac{\partial}{\partial \tau} \nabla J(r_t(\tau)) d\tau = \int_0^1 \nabla^2 J(r_t(\tau)) d\tau \underbrace{(\theta_{t+1} - \theta_t)}_{s_t}.$$

# Result – the quasi-Newton integral

With the definitions  $y_t \triangleq g_{t+1} - g_t$  and  $s_t \triangleq \theta_{t+1} - \theta_t$  we have

$$y_t = \int_0^1 \nabla^2 J(r_t(\tau)) d\tau s_t.$$

**Interpretation:** The difference between two consecutive gradients ( $y_t$ ) constitute a **line integral observation of the Hessian**.

**Problem:** Since the Hessian is unknown there is no functional form available for it.

# Solution 1 – existing quasi-Newton algorithms

---

Existing quasi-Newton algorithms (e.g. BFGS, DFP, Broyden's method) assume the Hessian to be constant

$$\nabla^2 J(r_t(\tau)) \approx H_{t+1}, \quad \tau \in [0, 1],$$

implying the following approximation of the integral (**secant condition**)

$$y_t = H_{t+1} s_t.$$

---

Find  $H_{t+1}$  by **regularizing**  $H$ :

$$\begin{aligned} H_{t+1} &= \min_H \|H - H_t\|_W^2, \\ \text{s.t. } H &= H^\top, \quad H s_t = y_t, \end{aligned}$$

Equivalently, the existing quasi-Newton methods can be interpreted as **particular instances of Bayesian linear regression**.

Philipp Hennig. Probabilistic interpretation of linear solvers, *SIAM Journal on Optimization*, 25(1):234–260, 2015.

## Solution 2 – use a flexible nonlinear model

---

The approach used here is fundamentally different.

Recall that the problem is **stochastic** and **nonlinear**.

Hence, we need a model that can deal with such a problem.

**Idea:** Represent the Hessian using a **Gaussian process** learned from data.

Adrian Wills and Thomas B. Schön. **Stochastic quasi-Newton with line-search regularization.** *Automatica*, 127:109503, May 2021.

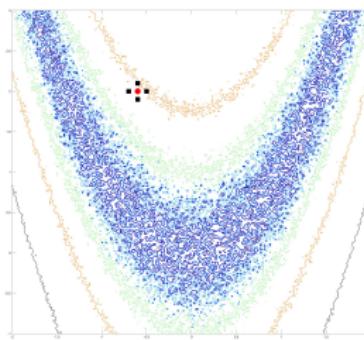
# Adam at work

---

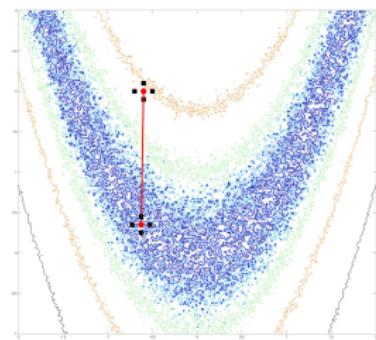
By not using the curvature information we expose ourselves to the "banana-problem".

# 2-order alg. at work — overall result

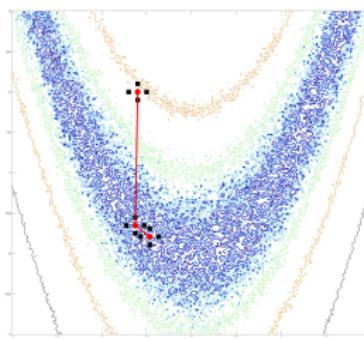
Initial value



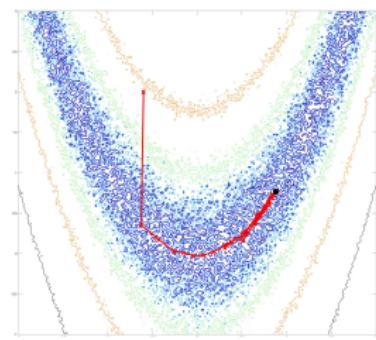
Iteration 1



Iteration 2



Iteration 50



# Adaptive (momentum) methods

Adaptive (momentum) methods **accelerates SG** especially when

- high curvature,
- small but consistent gradients,
- or noisy gradients.

**Idea:** Consider different ways of using the earlier gradients to adapt learning rates

$$\gamma_t = \gamma(\nabla J_t, \nabla J_{t-1}, \dots, \nabla J_0)$$

and search directions

$$d_t = d(\nabla J_t, \nabla J_{t-1}, \dots, \nabla J_0).$$

Update parameters as

$$\theta_{t+1} = \theta_t - \gamma_t d_t.$$

# Adaptive (momentum) methods

---

The most popular methods use an **exponential moving average**, where recent gradients have higher weights than older gradients.

Adam updates according to (typical values  $\beta_1 = 0.9, \beta_2 = 0.999$ )

$$\begin{aligned} d_t &= \beta_1 d_{t-1} + (1 - \beta_1) \nabla J_t, \\ \gamma_t^2 &= \beta_2 \gamma_{t-1}^2 + (1 - \beta_2) \|\nabla J_t\|^2. \end{aligned}$$

# Numerical opt. – Want to know more?

---

Great pedagogical introduction to the use of stochastic optimization for large-scale ML:

Bottou, L., Curtis, F. E. and Nocedal, J. **Optimization methods for large-scale machine learning**, *SIAM Review*, 60(2):223–311, 2018.

Solid and well-structured introduction to numerical optimization:

Nocedal, J. and Wright, S. **Numerical Optimization**, 2006.

---

To follow the research on optimization methods for ML, see the NeurIPS and ICML conferences.

[neurips.cc](http://neurips.cc)

[icml.cc](http://icml.cc)

# Backpropagation

# Bookkeeping – chain rule

---

Just as we can compute gradients of  $z$  w.r.t. scalars and vectors, we can compute gradients of  $z$  w.r.t. matrices and tensors.

Let  $\mathbf{X}$  denote a tensor. If  $z = f(\mathbf{Y})$  and  $\mathbf{Y} = g(\mathbf{X})$ , then

$$\nabla_{\mathbf{X}} z = \sum_j (\nabla_{\mathbf{X}} \mathbf{Y}_j) \frac{\partial z}{\partial \mathbf{Y}_j}$$

Note that the index variable  $j$  now encodes the complete tuple of indices, e.g. for a matrix, it contains two coordinates.

# Forward propagation and backpropagation

---

We only consider feedforward networks to focus on the key ideas. Generalizing this to a general network is then mainly a bookkeeping exercise.

**Forward propagation:** The input  $x$  is propagated through the network populating all hidden units with values and finally we reach the output  $\hat{y}$ . Based on this we can then compute the cost  $J(\theta)$ .

**Backpropagation:** Refers to the method for computing the gradient. It amounts to repeated use of the chain rule coupled with a memory that allows us to reuse computations.

# Ex. scalar regression

---

$$J(\theta) = \frac{1}{n} \sum_{i=1}^n (y_i - z^{(L)}(x_i, \theta))^2$$

To keep the notation simple we only consider one data point, i.e.

$$J(\theta) = (y - z^{(L)}(x, \theta))^2$$

# Forward propagation – fully connected MLP

---

Weight matrices:  $\mathbf{W}^{(i)}, i \in \{i = 1, \dots, L\}$ .

Offset terms:  $\mathbf{b}^{(i)}, i \in \{i = 1, \dots, L\}$ .

Input:  $\mathbf{x}$ , output:  $\mathbf{y}$

Network depth:  $L$

1.  $\mathbf{q}^{(0)} = \mathbf{x}$
2. **For**  $l = 1$  **to**  $L$  **do**
  - (a)  $\mathbf{z}^{(l)} = \mathbf{W}^{(l)}\mathbf{q}^{(l-1)} + \mathbf{b}^{(l)}$
  - (b)  $\mathbf{q}^{(l)} = h(\mathbf{z}^{(l)})$
3. **end for**
4.  $\mathbf{z}^L = \mathbf{W}^{(L)}\mathbf{q}^{(L-1)} + \mathbf{b}^{(L)}$
5.  $J = (y - z^{(L)})^2$

# Backpropagation

---

$$d\mathbf{z}^{(l)} \triangleq \nabla_{\mathbf{z}^{(l)}} \mathbf{J}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial z_1^{(l)}} \\ \vdots \\ \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial z_{U^{(l)}}^{(l)}} \end{bmatrix}, \quad d\mathbf{q}^{(l)} \triangleq \nabla_{\mathbf{q}^{(l)}} \mathbf{J}(\boldsymbol{\theta}) = \begin{bmatrix} \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial q_1^{(l)}} \\ \vdots \\ \frac{\partial \mathbf{J}(\boldsymbol{\theta})}{\partial q_{U^{(l)}}^{(l)}} \end{bmatrix}.$$

We need the gradients  $d\mathbf{z}^{(l)}$  and  $d\mathbf{q}^{(l)}$  for all layers.

---

Back to our example scalar regression problem

$$J(\boldsymbol{\theta}) = (y - z^{(L)}(x, \boldsymbol{\theta}))^2$$

$$dz^{(L)} = \nabla_{\mathbf{z}^{(L)}} \mathbf{J}(\boldsymbol{\theta}) = -2(y - z^{(L)})$$

# Backpropagation

---

Applying the chain rule to

$$\begin{aligned} \mathbf{q}^{(0)} &= \mathbf{x}, \\ \begin{cases} \mathbf{z}^{(l)} = \mathbf{W}^{(l)} \mathbf{q}^{(l-1)} + \mathbf{b}^{(l)}, \\ \mathbf{q}^{(l)} = h(\mathbf{z}^{(l)}), \end{cases} &\quad \text{for } l = 1, \dots, L-1 \\ \mathbf{z}^{(L)} &= \mathbf{W}^{(L)} \mathbf{q}^{(L-1)} + \mathbf{b}^{(L)}, \end{aligned}$$

results in

$$\begin{aligned} d\mathbf{z}^{(l)} &= d\mathbf{q}^{(l)} \odot h'(\mathbf{z}^{(l)}), \\ d\mathbf{q}^{(l-1)} &= \mathbf{W}^{(l)\top} d\mathbf{z}^{(l)}, \end{aligned}$$

where  $\odot$  denotes the element-wise product.

# Backpropagation

---

Finally, we can make use of  $d\mathbf{z}^{(l)}$  to compute the gradients of the weight matrices and offset vectors as

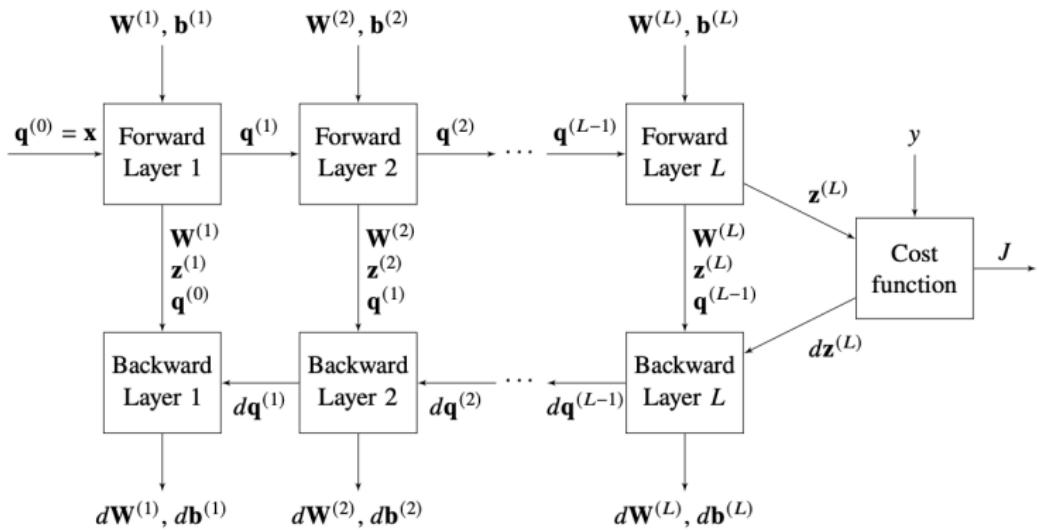
$$\begin{aligned} d\mathbf{W}^{(l)} &= d\mathbf{z}^{(l)} \mathbf{q}^{(l-1)\top}, \\ d\mathbf{b}^{(l)} &= d\mathbf{z}^{(l)}. \end{aligned}$$

The equations on this and the previous slide are derived using the **chain rule**, see Appendix 6.A for detailed derivation.

---

This can be generalized (essentially a bookkeeping exercise) into any network structure. The **computational graph** is a useful tool.

# Computational graph – backpropagation



Summarized in Algorithm 6.1.

# A few concepts to summarize lecture 3

---

**Gradient descent:** An iterative optimization algorithm where we at each iteration take a step proportional to the negative gradient.

**Learning rate: (a.k.a step length):** A scalar tuning parameter deciding the length of each gradient step in gradient descent.

**Stochastic gradient (SG):** A version of gradient descent where we at each iteration only use a small part of the training data (a mini-batch).

**Minibatch:** The group of training data that we use at each iteration in SG.

**Batch size:** The number of data points in one minibatch.

**Epoch:** One pass through the training data.

**Backpropagation:** A method based on the chain rule to compute the gradient by reusing information.