

**JOHANNES KEPLER
UNIVERSITY LINZ**

Deep Learning and Neural Networks I: 7. Learning methods



Günter Klambauer
LIT AI Lab & Institute for Machine Learning

JOHANNES KEPLER
UNIVERSITY LINZ
Altenberger Strasse 69
4040 Linz, Austria
jku.at

This material, no matter whether in printed or electronic form, may be used for personal and non-commercial educational use only. Any reproduction of this material, no matter whether as a whole or in parts, no matter whether in printed or in electronic form, requires explicit prior acceptance of the authors.

Current topic

- Largest and most important Machine Learning conference takes place next week: **NeurIPS**
- <https://neurips.cc/>
- Fully virtual conference; participation is at relatively low-cost (\$25 for full-time students)



- Side event on December 13:
<https://moleculediscovery.github.io/workshop2021/>

Learning methods: basic algorithms

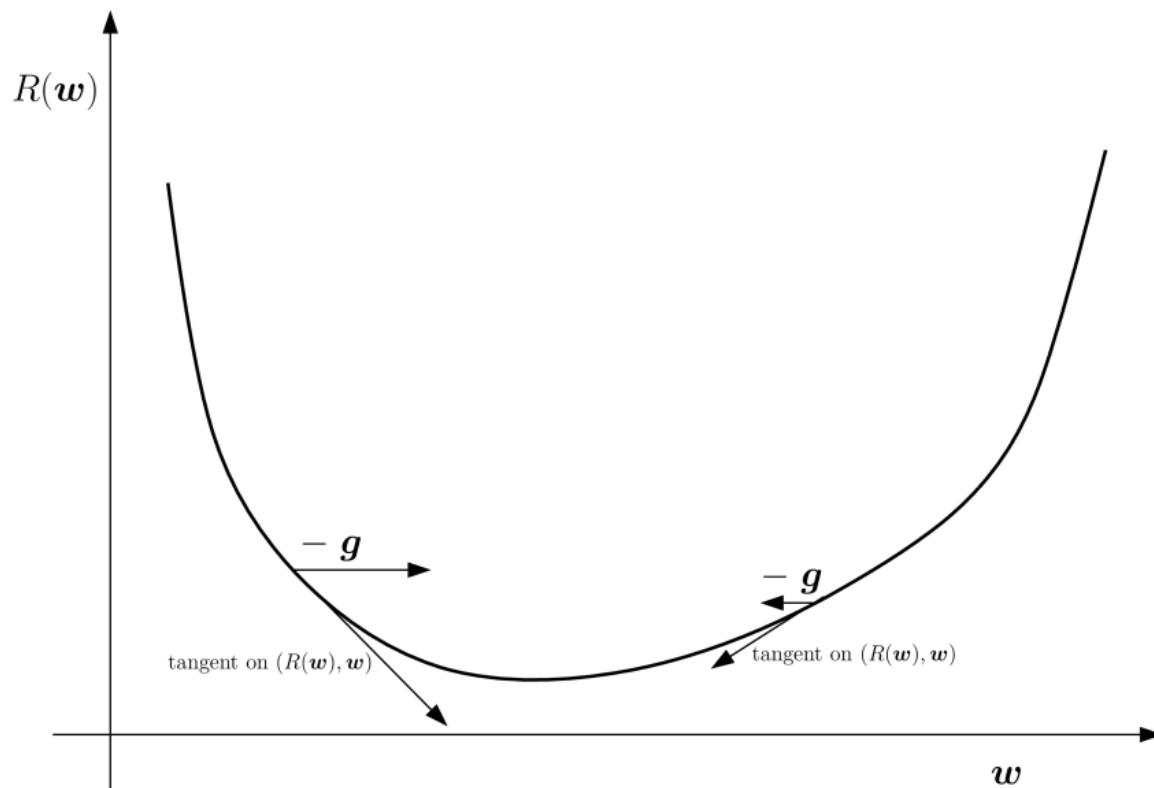
- In machine learning, we often adjust parameters w to optimize some objective function $R(w)$
- Objective functions, e.g. :
$$R(w) = \sum_{n=1}^N (y^n - w^T x^n)^2 + \|w\|^2$$
- Adapting those parameters is called *learning*
- Difference between “learning” and pure optimization:
 - e.g. we use the empirical error as surrogate to maximize accuracy.
- In ML, we often have access to the gradient of the objective function

Overview

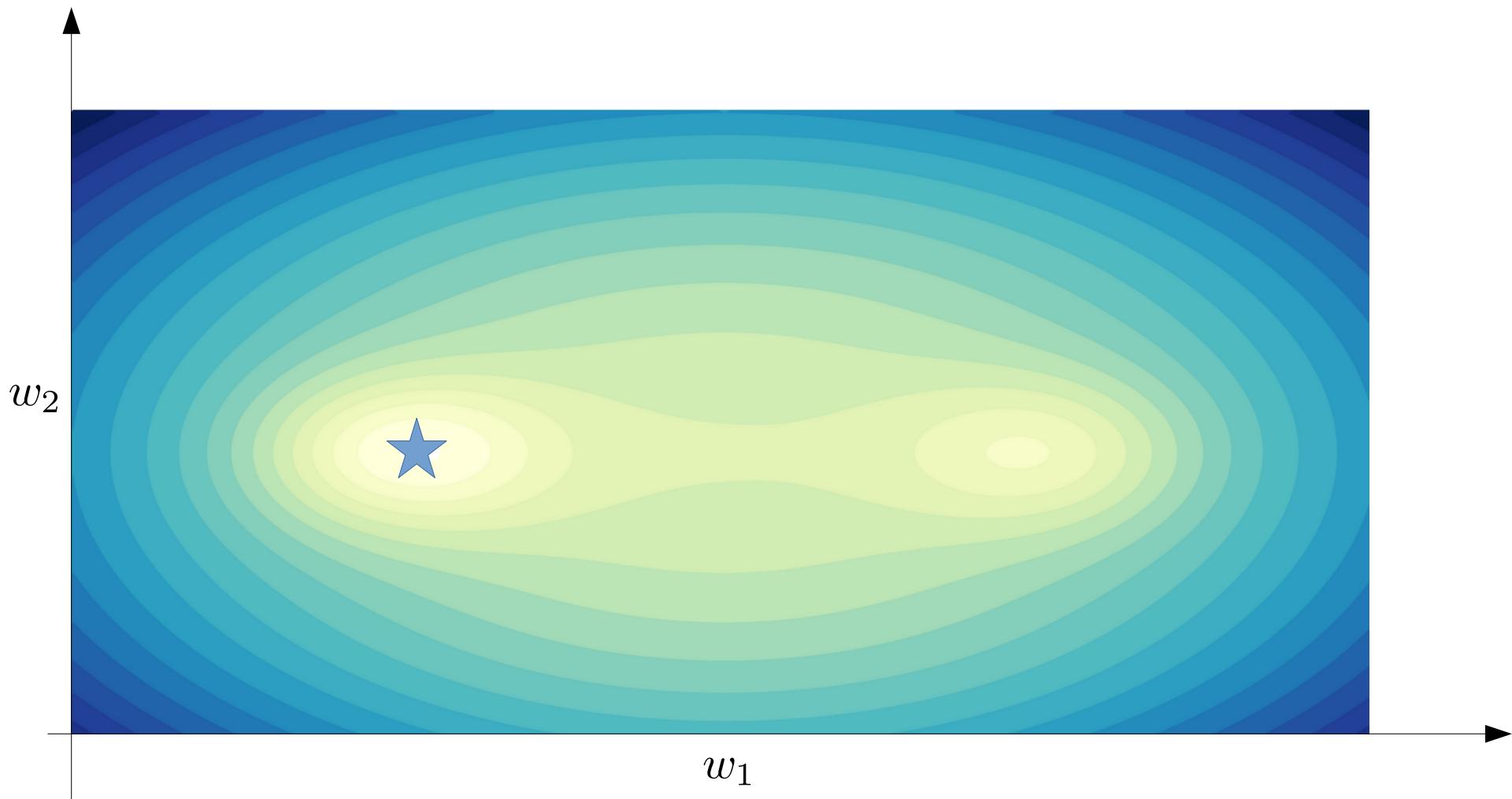
- 7.1 Gradient Descent
- 7.2 Gradient Descent with Momentum Term
- 7.3 Stochastic Gradient Descent (SGD)
- 7.4 Adaptive Learning Rate Optimizers
- 7.5 First-order Gradient Alternatives
- 7.6 Second-order Methods

7.1 Gradient descent

The negative gradient $-g$ gives the direction of the *steepest decent*, depicted by the tangent on $(R(\mathbf{w}), \mathbf{w})$:



Optimization problems



7.1 Gradient descent

Assume that the objective function $R(g(.; \mathbf{w}))$ is a differentiable function w. r. t \mathbf{w} .

The gradient is

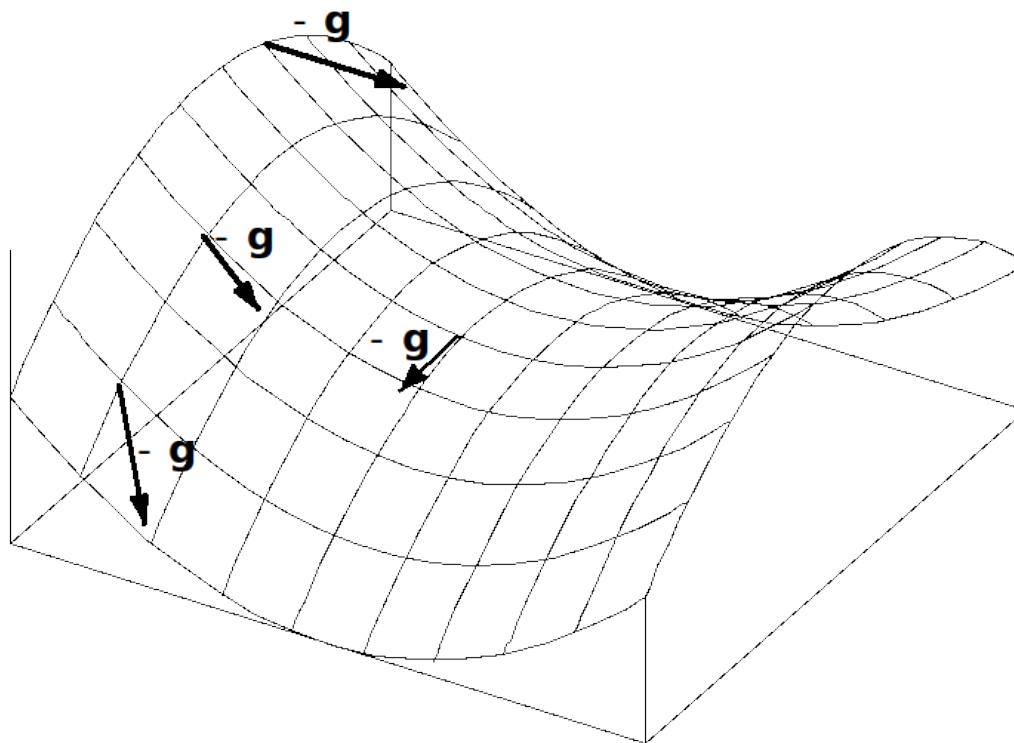
$$\frac{\partial R(\mathbf{w})}{\partial \mathbf{w}} = (\nabla_{\mathbf{w}} R(\mathbf{w}))^T = \left(\frac{\partial R(\mathbf{w})}{\partial w_1}, \dots, \frac{\partial R(\mathbf{w})}{\partial w_W} \right).$$

- For the gradient of $R(\mathbf{w})$ we use the *column vector*

$$\mathbf{g} = \nabla_{\mathbf{w}} R(\mathbf{w}).$$

7.1 Gradient descent

The negative gradient $-g$ attached at different positions on a two dimensional error surface ($R(w)$, w):



7.1 Gradient descent

- The gradient descent update is

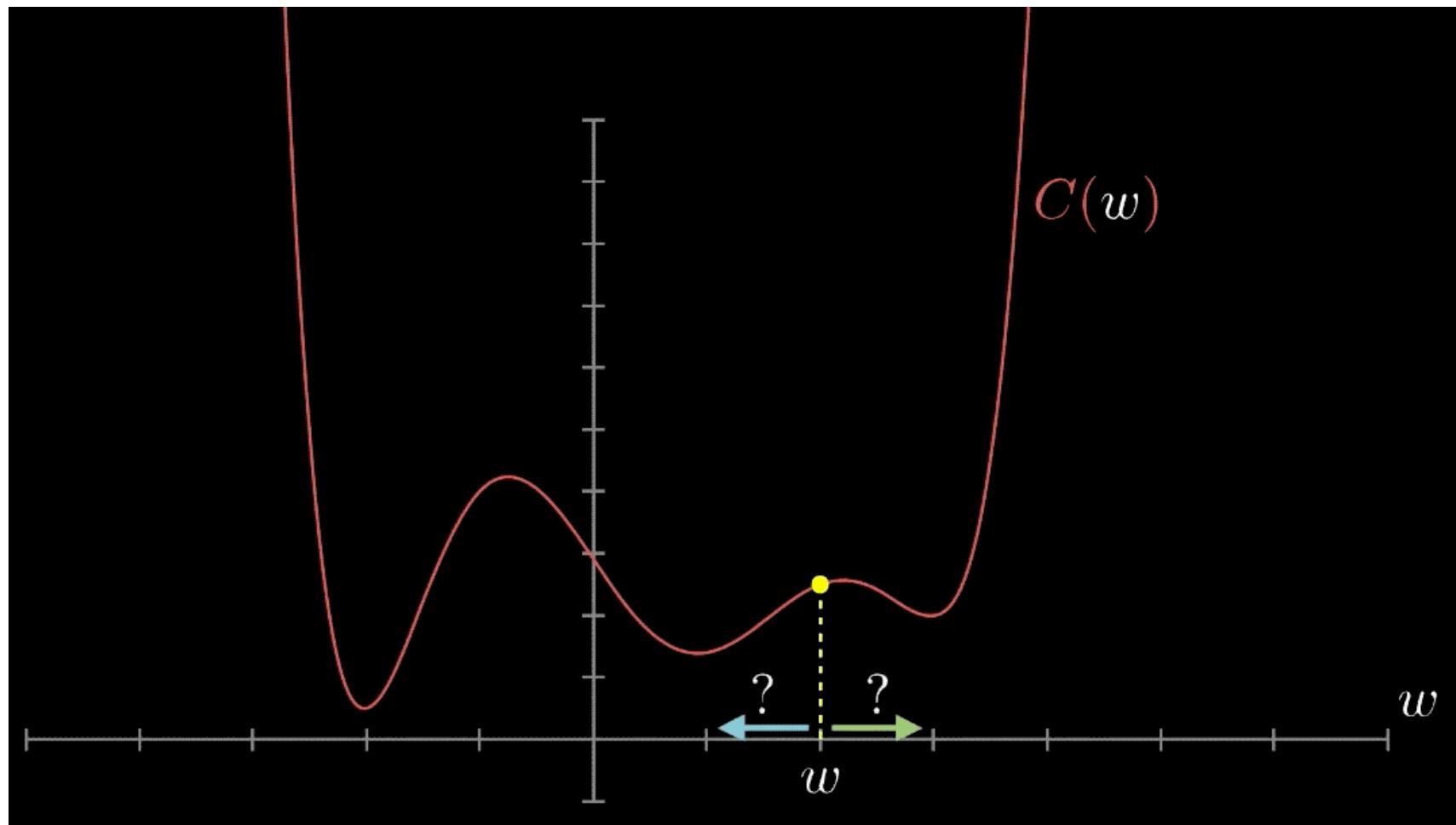
$$\begin{aligned}\Delta \mathbf{w}^{(t)} &= -\eta \nabla_{\mathbf{w}} R(\mathbf{w})|_{\mathbf{w}^{(t)}} \\ \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} + \Delta \mathbf{w}^{(t)}\end{aligned}$$

- We gradually change by update steps

$$\mathbf{w}^{(0)} \rightarrow \mathbf{w}^{(1)} \rightarrow \dots \rightarrow \mathbf{w}^{(t)}$$

```
while True:  
    dw = get_gradient(w)  
    w  = w - eta*dw
```

7.1 Gradient Descent



7.1 Gradient Descent: Convergence

Theorem 7.1. Suppose the function $R : \mathbb{R}^W \mapsto \mathbb{R}$ is convex and differentiable, and that its gradient is Lipschitz continuous with constant $L > 0$, i.e. $\|\nabla R(\mathbf{w}) - \nabla R(\mathbf{v})\|_2 \leq L\|\mathbf{w} - \mathbf{v}\|_2$ for any \mathbf{w}, \mathbf{v} . Then gradient descent with a fixed step size $\eta \leq 1/L$ will yield a solution $\mathbf{w}^{(t)}$ which satisfies

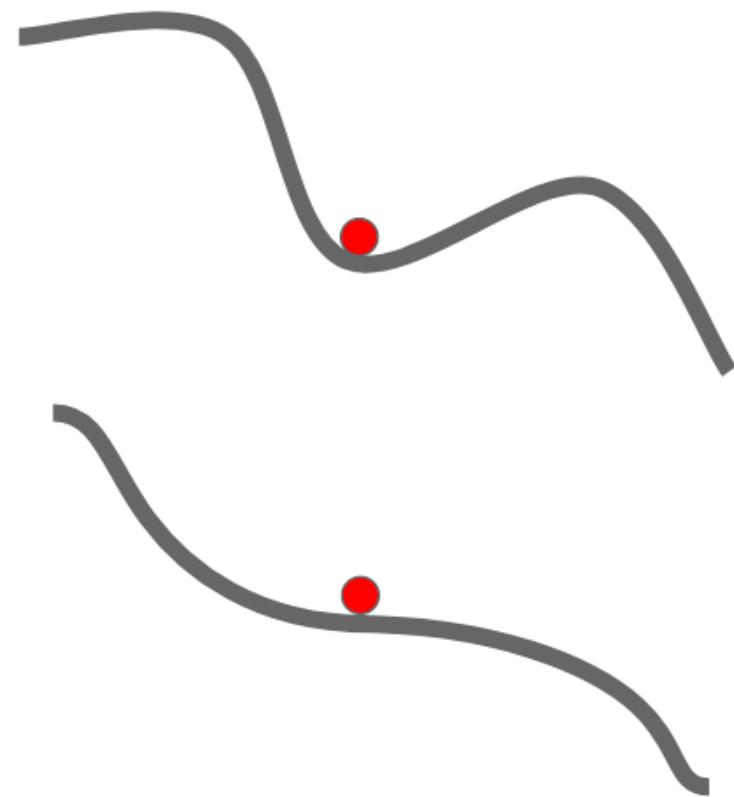
$$R(\mathbf{w}^{(t)}) - R(\mathbf{w}^*) \leq \frac{\|\mathbf{w}^{(0)} - \mathbf{w}^*\|_2^2}{2\eta t}, \quad (7.5)$$

where $R(\mathbf{w}^*)$ is the optimal value.

- Intuitively, Gradient Descent converges when the error function is convex and sufficiently smooth
 - The learning rate has to be sufficiently small
- (proof on blackboard/whiteboard)

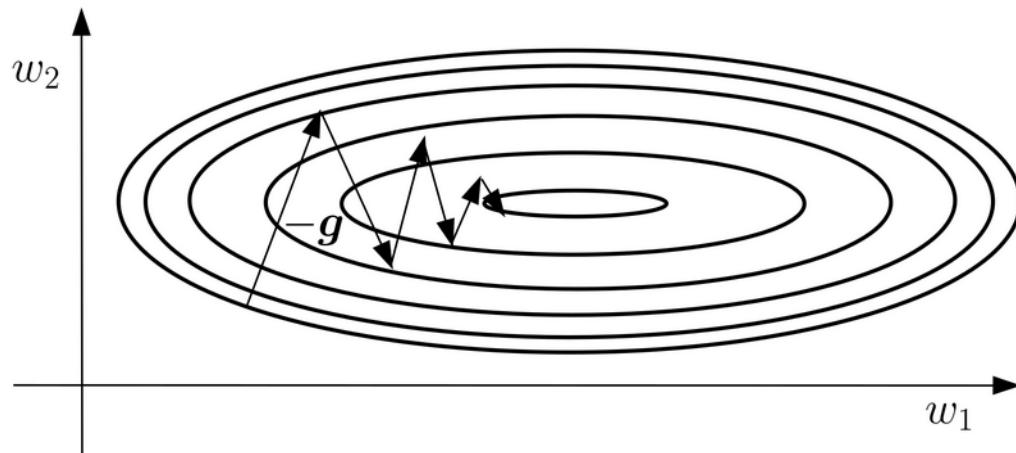
7.1 Gradient Descent: problems

- Local minima in loss function?
- Saddle points in loss function?
- Interplay with dimension of parameter space



7.1 Gradient Descent: problems

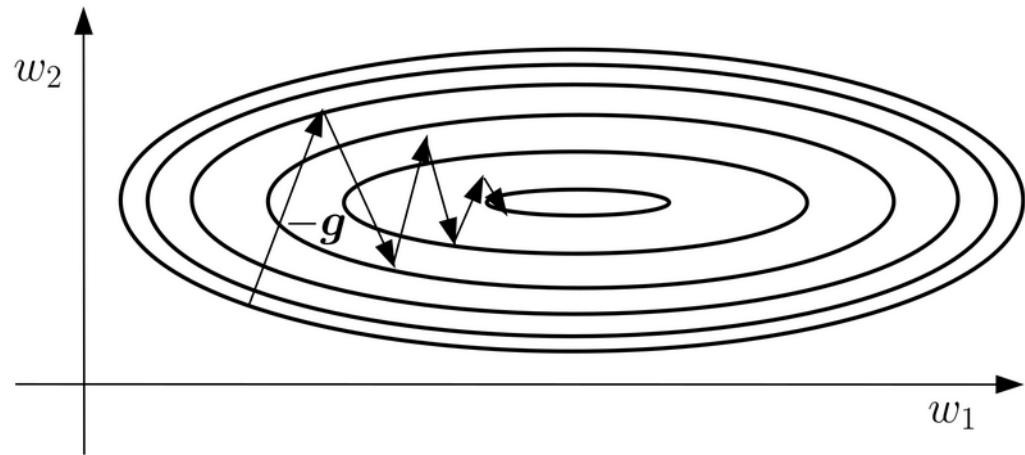
- Oscillation



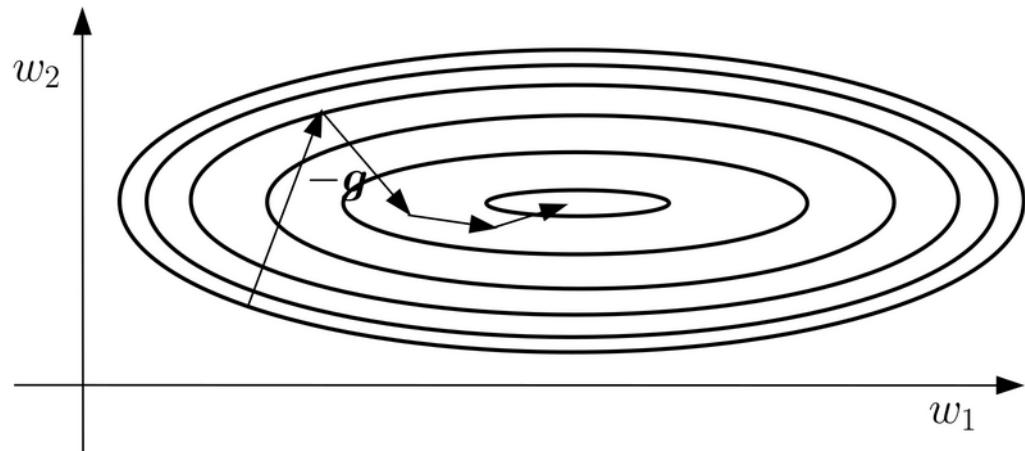
7.2 Gradient Descent with Momentum Term

Oscillation of the negative gradient $-g$

- without momentum:



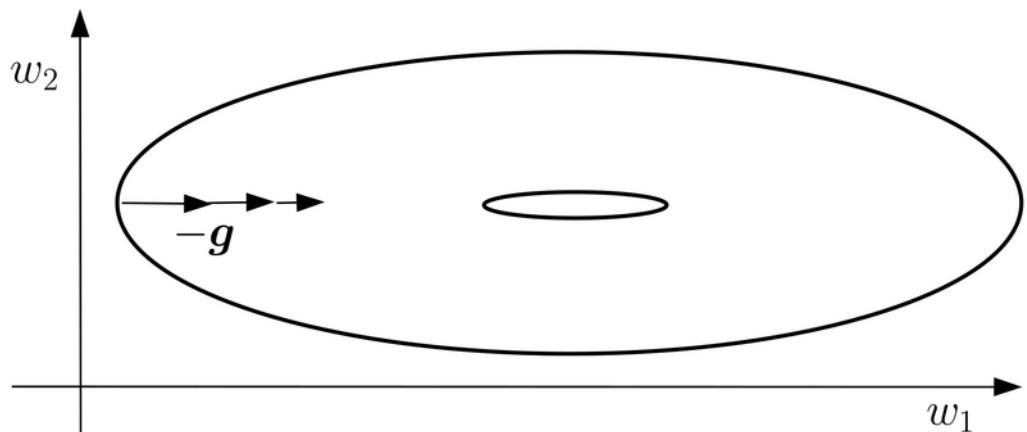
- with momentum:



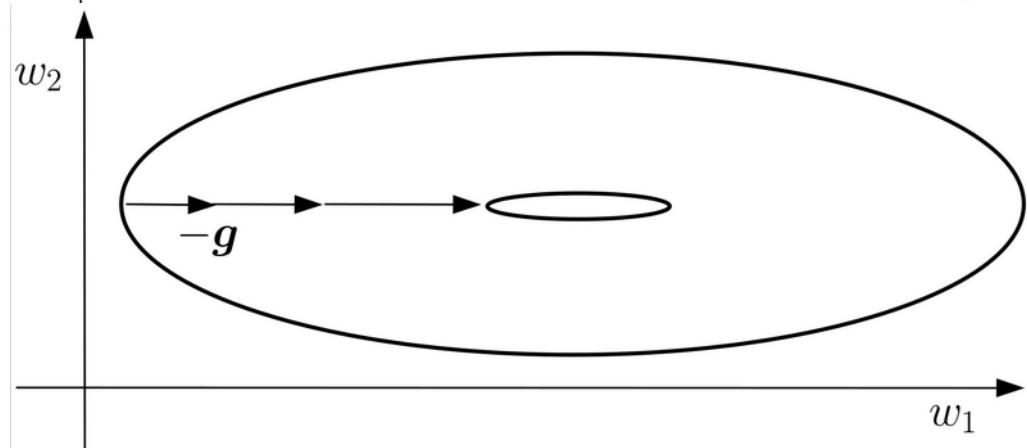
7.2 Gradient Descent with Momentum Term

Another effect of the momentum term is that in flat regions the gradients pointing in the same directions are accumulated and the learning rate is implicitly increased.

- without momentum:



- with momentum:



7.2 Gradient Descent with Momentum Term

The gradient descent update with momentum term is

$$\begin{aligned}\boldsymbol{m}^{(t)} &= \mu \boldsymbol{m}^{(t-1)} + (1 - \mu) \boldsymbol{g}^{(t)} \\ \boldsymbol{w}^{(t+1)} &= \boldsymbol{w}^{(t)} - \eta \boldsymbol{m}^{(t)}\end{aligned}$$

where $0 \leq \mu \leq 1$ is the *momentum parameter* or *momentum factor*, which is usually close to one.

```
vm = 0
while True:
    dw = get_gradient(w)
    vm = mu*vm + (1-mu)*dw
    w = w - eta*vm
```

7.2 Gradient Descent with Momentum Term Interpretation as physical quantity

Physical view of momentum

Learning with momentum has a particular interpretation as the velocity m of a physical object:

$$\begin{aligned} m^{(t)} &= \mu m^{(t-1)} + (1 - \mu)g^{(t)} = \\ &= \mu^2 m^{(t-2)} + \mu(1 - \mu)g^{(t-1)} + (1 - \mu)g^{(t)} = \\ &= \mu^3 m^{(t-3)} + \mu^2(1 - \mu)g^{(t-2)} + \mu(1 - \mu)g^{(t-1)} + (1 - \mu)g^{(t)} = \\ &= \mu^t m^{(0)} + (1 - \mu) \sum_{i=1}^t \mu^{t-i} g^{(i)} \end{aligned}$$

7.2 Gradient Descent with Momentum Term Interpretation as physical quantity

From this and $m^{(0)} = 0$ follows that

$$\Delta w^{(t)} = - \sum_{i=1}^t \mu^{t-i} (\eta(1-\mu)g^{(i)}) .$$

Which means that we could also write the algorithm as

$$\begin{aligned} \mathbf{m}^{(t)} &= \mu \mathbf{m}^{(t-1)} - \eta' \mathbf{g} \\ \Delta \mathbf{w}^{(t)} &= \mathbf{m}^{(t)} \end{aligned} , \quad \eta' = -\eta(1-\mu).$$

We can interpret the right hand side as force times a time interval (which is equal to one).

$$\begin{aligned} \mathbf{m}^{(t)} - \mathbf{m}^{(t-1)} &= \left(-(1-\mu)\mathbf{m}^{(t-1)} - \eta' \mathbf{g}^{(t)} \right) \Delta t \\ \frac{d\mathbf{m}(t)}{dt} &= \underbrace{-(1-\mu)\mathbf{m}(t)}_{\text{friction}} \underbrace{-\eta' \mathbf{g}(t)}_{\text{Lossforce}} \end{aligned}$$

7.2 Gradient Descent with Momentum Term Interpretation as physical quantity

- First we look at the case when $\mathbf{g}(t) = 0$:

$$\mathbf{m}(t) = \mathbf{m}(t_0) e^{-(1-\mu)t}.$$

- Now we take a look at the parameter update in continuous time:

$$\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)} = \mathbf{m}^{(t)} \Delta t$$

$$\frac{\mathbf{w}^{(t)} - \mathbf{w}^{(t-1)}}{\Delta t} = \mathbf{m}^{(t)}.$$

In the limit $\Delta t \rightarrow 0$ we get:

$$\frac{d\mathbf{w}(t)}{dt} = \mathbf{m}(t).$$

- Differential equation which describes the motion of $\mathbf{w}(t)$:

$$\frac{d^2\mathbf{w}(t)}{dt^2} = -(1 - \mu) \frac{d\mathbf{w}(t)}{dt} - \eta' \mathbf{g}(t).$$

7.3 Stochastic Gradient Descent (SGD)

Objective function for neural networks:

$$R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N L(y^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w}))$$

- **Full-batch-gradient:**

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \frac{1}{N} \sum_{n=1}^N \nabla_{\mathbf{w}} L(y^n, \mathbf{g}(\mathbf{x}^n; \mathbf{w})).$$

- **Stochastic gradient descent:**

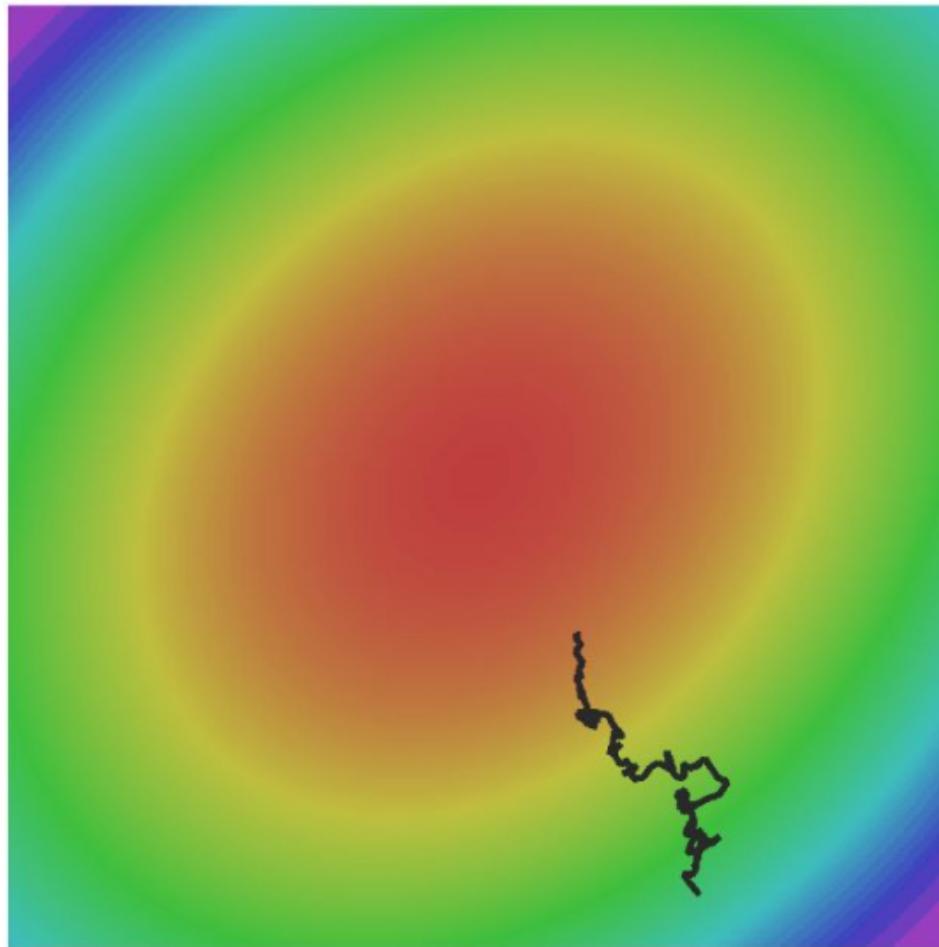
Minibatch: $\mathcal{B} = \{j_1, \dots, j_B\} \subset \{1, \dots, N\}$

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{y}) \approx \frac{1}{B} \sum_{b=1}^B \nabla_{\mathbf{w}} L(y^{j_b}, \mathbf{g}(\mathbf{x}^{j_b}; \mathbf{w})).$$

- **On-line learning:** SGD when $B = 1$.

7.3 Stochastic Gradient Descent

Noisy estimates of gradient



7.3 Stochastic Gradient Descent (SGD)

Efficiency of stochastic gradient descent:

- The expression for the empirical error is an expectation of the loss of the training set:

$$\nabla_{\mathbf{w}} R_{\text{emp}}(\mathbf{w}, \mathbf{X}, \mathbf{y}) = \mathbb{E}_{\mathbf{x}, \mathbf{y} \sim \tilde{p}_{\text{data}}} [\nabla_{\mathbf{w}} L(y, \mathbf{g}(\mathbf{x}; \mathbf{w}))],$$

- The improvement of the estimate does not scale linearly with the batch size B , since the standard error of the mean of B samples is given by

$$\frac{\sigma}{\sqrt{B}},$$

where σ is the true standard deviation in the samples.

- Example: 100,000 samples in training data; gradient calculated on 1,000. 100 times less computations; gradient standard error 10 times larger

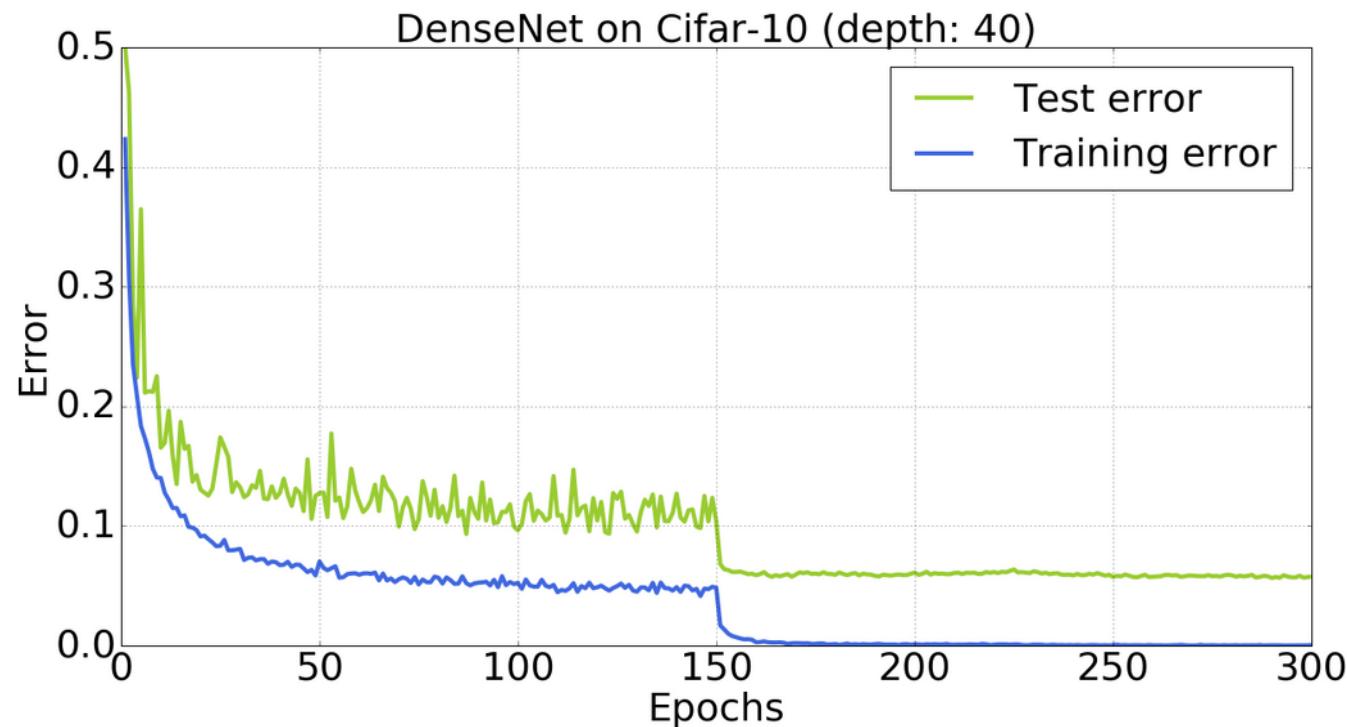
7.3 Stochastic Gradient Descent (SGD)

Implementation of learning with stochastic gradient descent:

- Sampling without replacement: A mini-batch is drawn and removed from training set
- Update step:
 - The gradient on this mini-batch is calculated
 - The parameters are updated
- When training set is depleted, an **epoch** has ended.
Number of updates per epoch: $\lceil N/B \rceil$.
- Practitioners monitor the loss and performance metric on the training set (and a validation set).

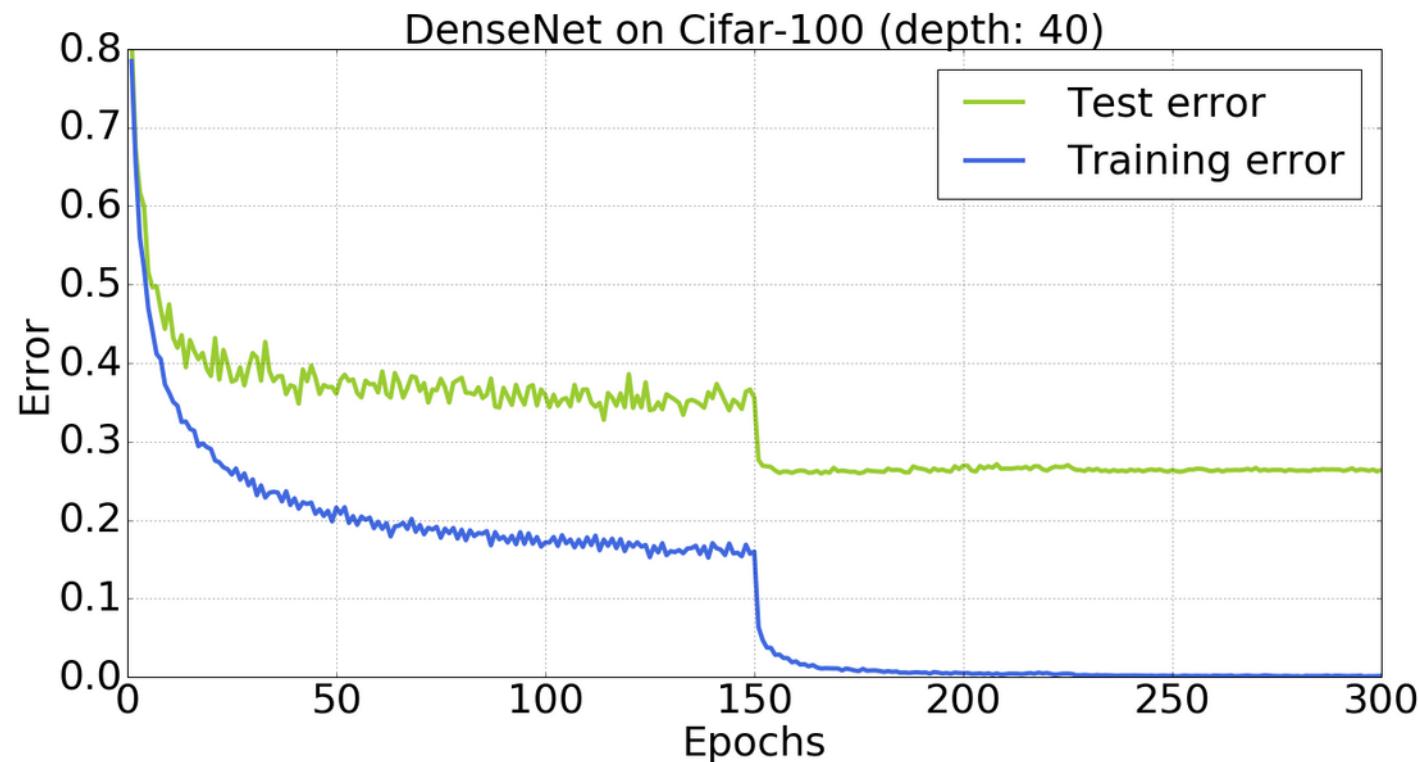
7.3 Stochastic Gradient Descent (SGD)

Learning curves for training of a neural network (DenseNet)



7.3 Stochastic Gradient Descent (SGD)

Learning curves for training of a neural network (DenseNet)



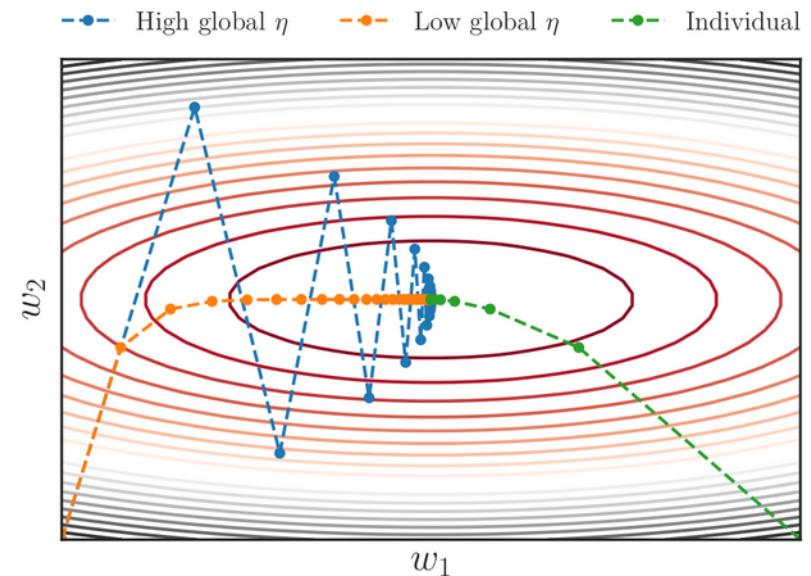
7.4 Adaptive Learning Rate Optimizers

Choosing the learning rate is critical for success

- Learning rate too high: update may “overshoot” its goal
- Learning rate too slow: progress may be slow
- Each parameter might benefit from different learning rates
- A specific learning η_i rate for each parameter w_i can be beneficial:

$$\Delta w_i = -\eta_i g_i$$

or $\Delta \mathbf{w} = -\boldsymbol{\eta} \odot \mathbf{g}$



7.4 Adaptive Learning Rate Optimizers

$\Delta\mathbf{w}$ leads to a point with lower error: $\Delta\mathbf{w} = -\boldsymbol{\eta} \odot \mathbf{g}$

$$R(\mathbf{w} + \Delta\mathbf{w}) = R(\mathbf{w}) + \mathbf{g}^T \Delta\mathbf{w} + \mathcal{O}(\Delta\mathbf{w}^2)$$

$$\begin{aligned} \text{With } -\mathbf{g}^T (\boldsymbol{\eta} \odot \mathbf{g}) &= -\sum_i g_i g_i \eta_i = \\ &\leq -\sum_i g_i g_i \eta_{\min} \\ &= -\eta_{\min} \sum_i g_i g_i \leq 0 \end{aligned}$$

we get $R(\mathbf{w} + \Delta\mathbf{w}) \leq R(\mathbf{w})$.

7.4 Adaptive Learning Rate Optimizers

Further issue of a fixed learning rate:

This assumes that the loss surface looks similar everywhere

Therefore useful: Adaption of the learning rate to the topology of the loss surface at the current point.

$$\boldsymbol{w}^{(t+1)} = \boldsymbol{w}^{(t)} - \boldsymbol{\eta}^{(t)} \odot \boldsymbol{g}^{(t)}.$$

the vector $\boldsymbol{\eta}$ provides an individual learning rate for each parameter. The elements of $\boldsymbol{\eta}$ must be non-negative.

7.4 Adaptive Learning Rate Optimizers

AdaGrad

- AdaGrad: adapts step size by remembering past gradients.
With the sum of squared past gradients for each parameter

$$v_i^{(t)} = \sum_{s=1}^t (g_i^{(s)})^2, g_i^{(s)} : i\text{-th component of the gradient at step } s$$

we get:

$$w_i^{(t+1)} = w_i^{(t)} - \eta \frac{1}{\sqrt{v_i^{(t)} + \epsilon}} g_i^{(t)},$$
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \frac{1}{\sqrt{\mathbf{v}^{(t)} + \epsilon}} \odot \mathbf{g}^{(t)}.$$

```
grad_squared=0
typical while True:
    dw = get_gradient(w)
    grad_squared = grad_squared + dw * dw
    w = w - eta*dw/(np.sqrt(grad_squared+1e-7))
```

7.4 Adaptive Learning Rate Optimizers: RMSprop, AdaDelta

- RMSprop:

$$\begin{aligned}\mathbf{v}^{(t)} &= \gamma \mathbf{v}^{(t-1)} + (1 - \gamma) (\mathbf{g}^{(t)})^2 \\ \mathbf{w}^{(t+1)} &= \mathbf{w}^{(t)} - \eta \frac{1}{\sqrt{\mathbf{v}^{(t)} + \epsilon}} \odot \mathbf{g}^{(t)}.\end{aligned}$$

```
grad_squared=0
while True:
    dw = get_gradient(w)
    grad_squared = gamma * grad_squared + (1-gamma) * dw * dw
    w = w - eta*dw/(np.sqrt(grad_squared+1e-7))
```

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \frac{\sqrt{\mathbf{d}^{(t-1)}} + \epsilon}{\sqrt{\mathbf{v}^{(t)} + \epsilon}} \odot \mathbf{g}^{(t)},$$

where: $\gamma \in [0, 1)$.

7.4 Adaptive Learning Rate Optimizers

Adam

- Combines the idea of momentum term and Adelta by memorizing both exponentially decaying averages of *past gradients* $\mathbf{g}^{(t)}$ and *past squared gradients* $(\mathbf{g}^{(t-1)})^2$.

$$\mathbf{m}^{(t)} = \beta_1 \mathbf{m}^{(t-1)} + (1 - \beta_1) \mathbf{g}^{(t)}$$

$$\mathbf{v}^{(t)} = \beta_2 \mathbf{v}^{(t-1)} + (1 - \beta_2) (\mathbf{g}^{(t)})^2$$

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \frac{1}{\sqrt{\mathbf{v}^{(t)}} + \epsilon} \odot \mathbf{m}^{(t)}$$

```
first_moment = 0
second_moment = 0
while True:
    dw = get_gradient(w)
    first_moment = beta1 * first_moment + (1-beta1) * dw
    second_moment = beta2 * second_moment + (1-beta2) * dw * dw
    w = w - eta*first_moment / (np.sqrt(second_moment)+1e-7)
```

J:

7.4 Adaptive Learning Rate Optimizers: Delta-delta rule; delta-bar-delta rule

- Delta-delta rule:
Learning rate at step t

$$\eta_i^{(t)} = \eta_i^{(t-1)} + \gamma g_i^{(t)} g_i^{(t-1)} \quad , \gamma : \text{step size}$$

- Delta-bar-delta rule:

$$\eta_i^{(t)} = \eta_i^{(t-1)} + \Delta\eta_i^{(t)}$$
$$\Delta\eta_i^{(t)} = \begin{cases} \kappa & \bar{g}_i^{(t)} g_i^{(t)} > 0 \\ -\phi\eta_i^{(t-1)} & \bar{g}_i^{(t)} g_i^{(t)} < 0 \\ 0 & \text{otherwise} \end{cases}$$

where

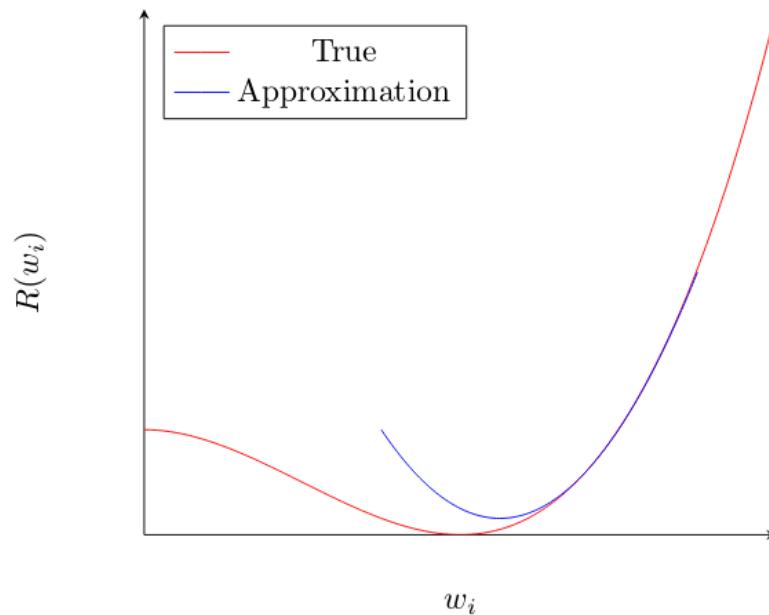
$$\bar{g}_i^{(t)} = \theta \bar{g}_i^{(t-1)} + (1 - \theta) g_i^{(t)}$$

Demo on Optimizers

7.5 First Order Gradient Alternatives

Quickprop

- Quickprop:
Main Idea: Assume that loss surface is locally quadratic.



The resulting update is: $\Delta^{\text{new}} w_i = \frac{g_i^{\text{new}}}{g_i^{\text{old}} - g_i^{\text{new}}} \Delta^{\text{old}} w_i.$

7.5 First Order Gradient Alternatives

Quickprop

- Assume $R(\mathbf{w})$ is a piece-wise function so that: $R(\mathbf{w})_i = R(w_i)$.
- Taylor expansion:

$$R(w_i + \Delta w_i) = R(w_i) + g_i \Delta w_i + \frac{1}{2} g'_i (\Delta w_i)^2 + \mathcal{O}((\Delta w_i)^3)$$

- To maximize the change in the objective, one can minimize

$$R(w_i + \Delta w_i) - R(w_i) = g_i \Delta w_i + \frac{1}{2} g'_i (\Delta w_i)^2$$

- We set the derivative w. r. t. Δw_i to zero and obtain: $\Delta w_i = -\frac{g_i}{g'_i}$
- Now, approximate: $g'_i \approx \frac{g_i^{\text{new}} - g_i^{\text{old}}}{w_i - (w_i - \Delta^{\text{old}} w_i)}$
- Plugging this approximation into $\Delta w_i = -\frac{g_i}{g'_i}$ gives the update

7.5 First Order Gradient Alternatives

Rprop

- Gradient is completely replaced

$$\Delta w_i = -\eta_i^{\text{new}} \text{sign}(g_i^{\text{new}}).$$

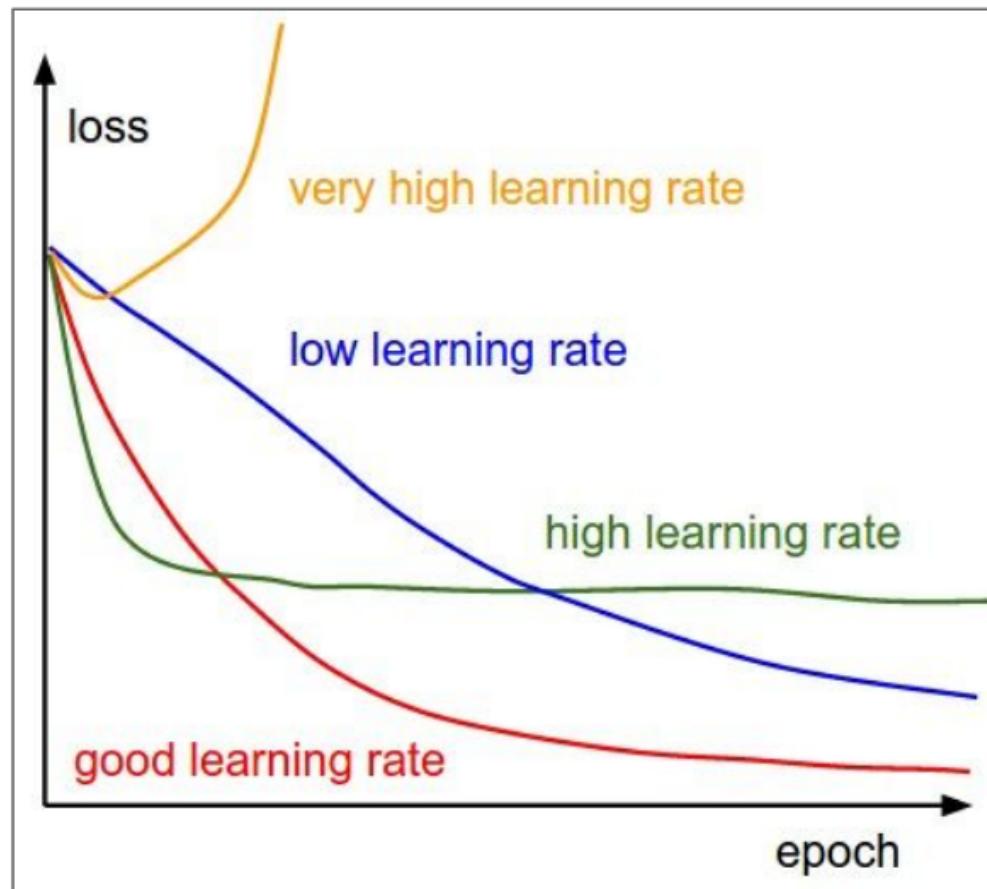
- Learning rate adjustment similar to delta-bar-delta rule:

$$\eta_i^{\text{new}} = \begin{cases} \eta^+ \eta_i^{\text{old}} & g_i^{\text{old}} g_i^{\text{new}} > 0 \\ \eta^- \eta_i^{\text{old}} & g_i^{\text{old}} g_i^{\text{new}} < 0 \\ \eta_i^{\text{old}} & \text{otherwise} \end{cases}$$

- Hyperparameters: η^+, η^-

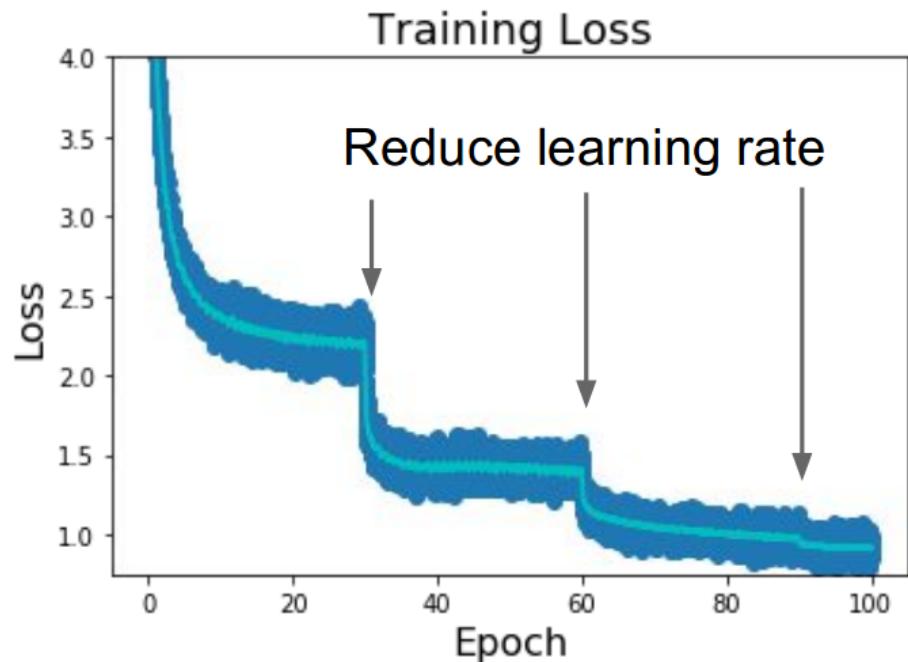
Choosing learning rates; learning rate decay

- Learning rate (LR) is a crucial hyperparameter
- LR can be decreased during learning



Choosing learning rates; learning rate decay

- Learning rate decay strategies:
 - No decay
 - Exponential decay
 - $1/t$ decay $\eta = \eta_0 / (1 + kt)$
 - Decay at plateaus
 - “Cosine annealing”
$$\eta = 1/2\eta_0(1 + \cos(t\pi/T))$$

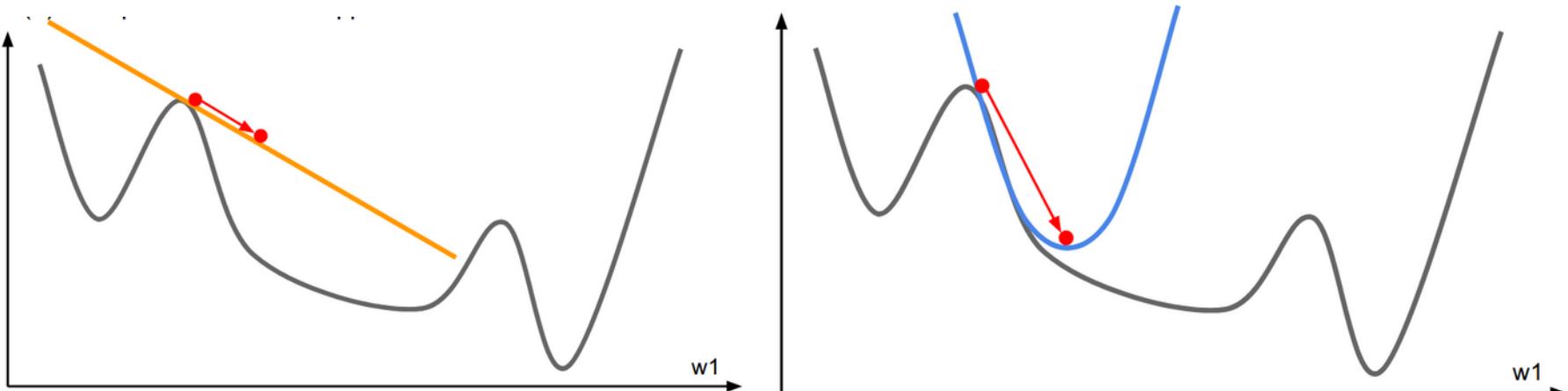


Loshchilov, I., & Hutter, F. (2016). Sgdr: Stochastic gradient descent with warm restarts. arXiv preprint arXiv:1608.03983.

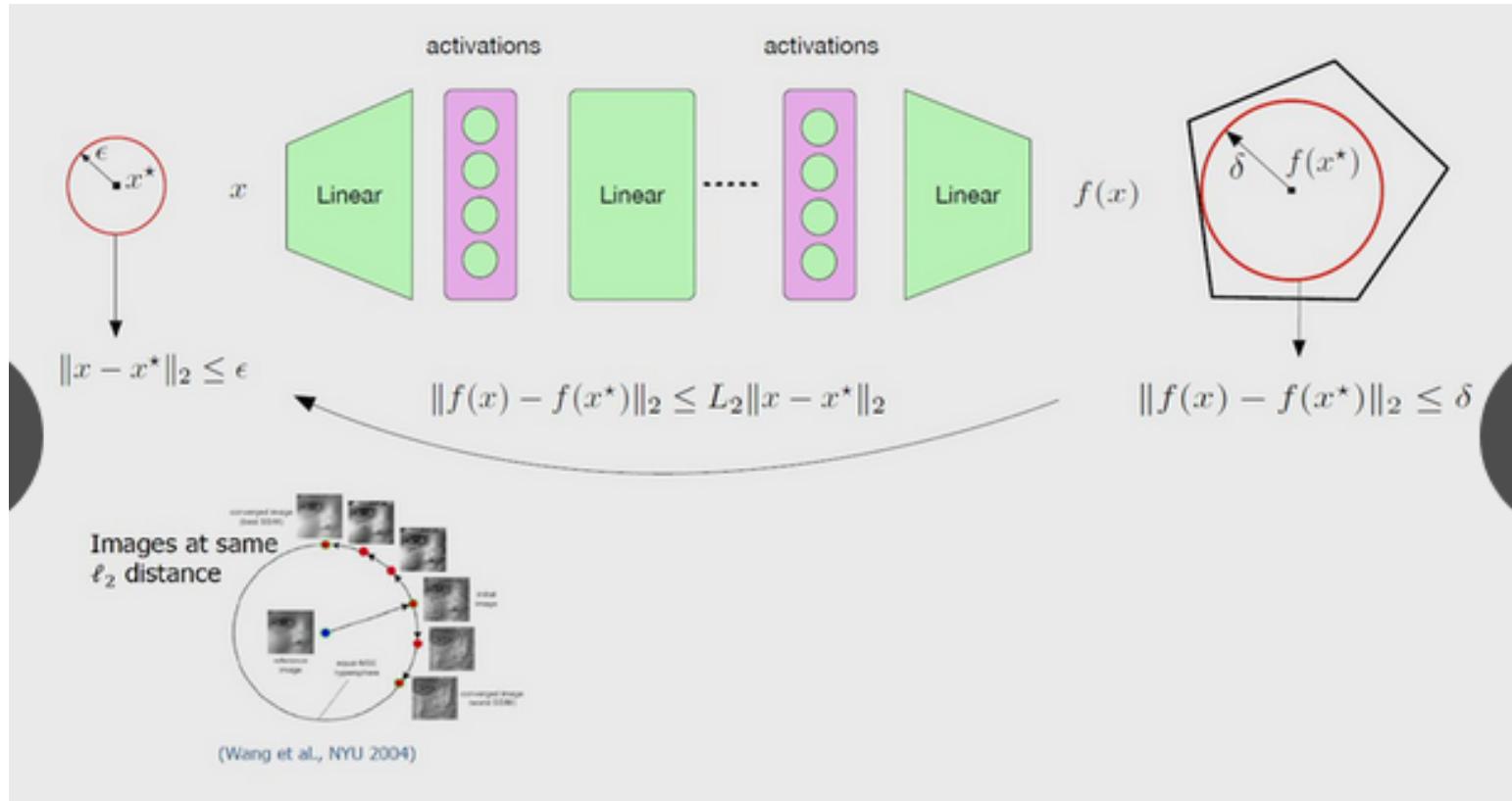
- Linear,
Inverse sqrt,..

7.5 Second order methods

- Use the information on curvature (Hessian)
- Currently not used to train neural networks
 - Why?
- See also part II of this lecture.



NeurIPS2019 and Lipschitz continuity



Fazlyab, M., Robey, A., Hassani, H., Morari, M., & Pappas, G. J. (2019). Efficient and Accurate Estimation of Lipschitz Constants for Deep Neural Networks. arXiv preprint arXiv:1906.04893.