

Faculty of Mathematics, Physics and Informatics
Comenius University Bratislava



Neural Networks

Lecture 5

Gradient-based learning and optimization

Optimization vs NN learning

- Although optimization provides a way to minimize the loss function for NN learning, the goals are fundamentally different:
- goal of optimization = **to reduce the error on given dataset**
- goal of NN learning (statistical inference) = **to reduce expected generalization error** (risk) => ML acts indirectly
- We have only access to a **finite** training data sample (not the whole data distribution)
- Empirical risk minimization: $E_{(\mathbf{x},d) \sim p(\text{data})}[\text{Loss}(f(\mathbf{x};\mathbf{w}),d)]$
- ... is based on a finite training sample $\{\mathbf{x},d\}$, rather than known data distribution, hence is prone to overfitting.
- Typically we apply **early stopping** criterion using the validation set.

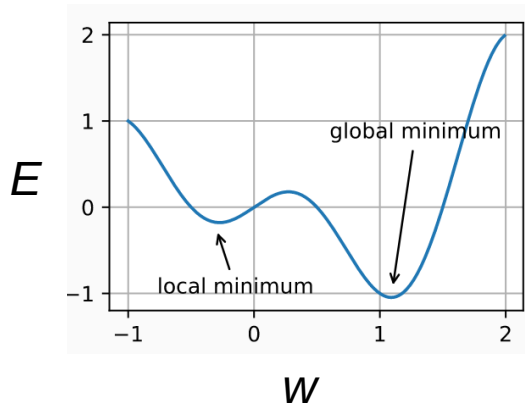
(Goodfellow et al, 2015)

Surrogate loss function

- Sometimes, the loss function we actually care about (e.g. classification error) is not the one that can be optimized efficiently.
- Exactly minimizing expected 0-1 loss is typically intractable (exponential in the input dimension)
- In such situations, one typically optimizes a surrogate loss function instead, which acts as a proxy, but has advantages:
- e.g. the negative log-likelihood of the correct class is used ($-\log P(y_i)$)
- Decrease of error, after the training set 0-1 loss has reached zero, improves the robustness of the classifier by further pushing the classes apart from each other.
- This leads to extracting more information from the training data (than would have been possible by simply minimizing the average 0-1 loss on training set).

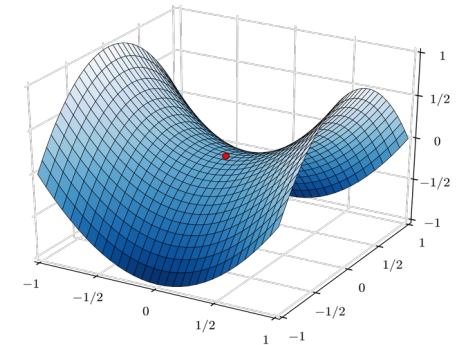
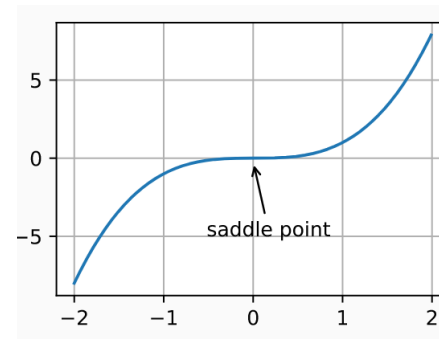
Problems in gradient-based NN learning

Local minima

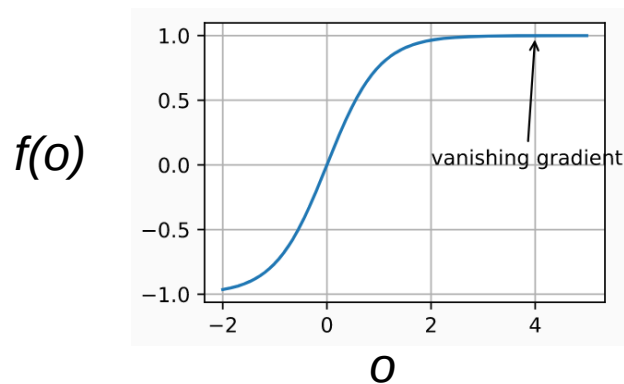


w = weight,
 E = loss f.

Saddle points



Vanishing gradients



Ill-conditioning of Hessian matrix \mathbf{H} ,
i.e. rate of its change for small $\Delta \mathbf{w}$:

- given by condition number (CN) = ratio of its max/min eigenvalues
- for large CN, \mathbf{H}^{-1} is particularly sensitive to error in the input

These problems slow down or hinder convergence.

Role of the Hessian matrix

- \mathbf{H} plays an important role in supervised training of neural networks:

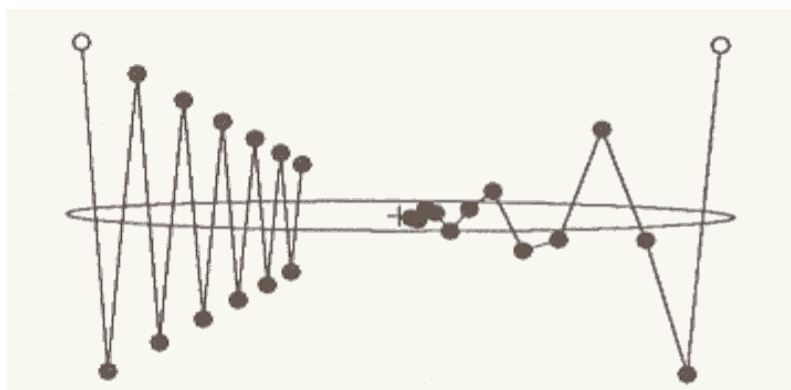
$$\mathbf{H} = \left[\frac{\partial^2 E(\mathbf{w})}{\partial w_i \partial w_j} / \mathbf{w}_0 \right]_{ij}$$

- Spread of eigenvalues of \mathbf{H} has a profound influence on the dynamics of back-propagation learning (condition number).
- The inverse of \mathbf{H} provides a basis for pruning (i.e., deleting) insignificant synaptic weights from a multilayer perceptron.
- \mathbf{H} is basic to the formulation of second-order optimization methods as an alternative to BP learning.
- Typical profile of \mathbf{H} in BP learning (LeCun et al., 1998): a few small eigenvalues, many medium-sized eigenvalues, and a few large eigenvalues => a wide spread in the eigenvalues of the Hessian.

Early modifications of gradient descent learning

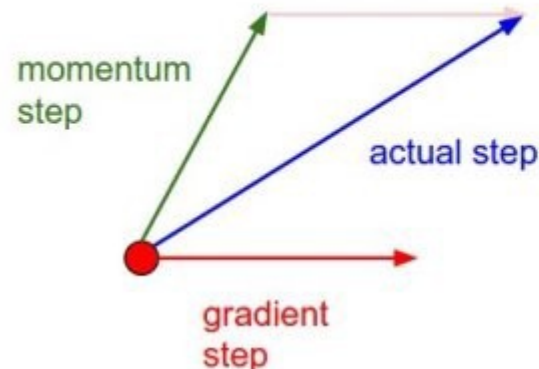
Adding a momentum: $\Delta \mathbf{w}(t) = -\alpha \nabla E(\mathbf{w}(t)) + \gamma \Delta \mathbf{w}(t-1)$ $0 \leq |\gamma| < 1$

- helps speed up SGD and dampen oscillations



w/out

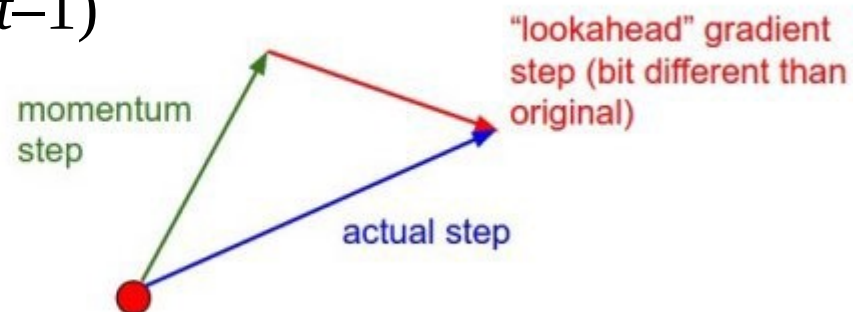
with



Nesterov accelerated gradient (Nesterov, 1983)

$$\Delta \mathbf{w}(t) = -\alpha \nabla E[\mathbf{w}(t) + \gamma \Delta \mathbf{w}(t-1)] + \gamma \Delta \mathbf{w}(t-1)$$

- helps adjust learning speed by looking into the near future



Towards second-order optimization methods

$$E(\mathbf{w}) = E(\mathbf{w}_0) + \mathbf{g}^T(\mathbf{w}_0)\Delta\mathbf{w} + 1/2 \Delta\mathbf{w}^T \mathbf{H}(\mathbf{w}_0)\Delta\mathbf{w} + O^{3+}(\Delta\mathbf{w})$$

$$\Delta\mathbf{w} = \mathbf{w} - \mathbf{w}_0$$

Taylor expansion:

Gradient vector:

$$\mathbf{g}(\mathbf{w}_0) = \nabla E(\mathbf{w}_0) = \left[\frac{\partial E}{\partial w_1} / \mathbf{w}_0, \dots, \frac{\partial E}{\partial w_{|\mathbf{w}|}} / \mathbf{w}_0 \right]^T$$

$$1D: f(x) = \sum_{i=0}^{\infty} \frac{f^{(i)}(x_0)}{i!} (x - x_0)^i$$

- Error back-propagation is a linear approximation of E : $\Delta\mathbf{w}(t) = -\alpha \mathbf{g}(t)$
- Quadratic approx. of $E(\mathbf{w}) \rightarrow$ **Newton's method**: $\Delta\mathbf{w} = -\mathbf{H}^{-1}(t) \mathbf{g}(t)$
- **Quasi-Newton method** approximates $\mathbf{H}^{-1}(t)$ with a positive definite matrix
- **Conjugate-gradients methods** are intermediate between the steepest descent and the Newton's method, by achieving faster convergence (than the former) and lower computational complexity (than the latter).

Conjugate-gradient methods

- 2nd order optimization methods
- minimize the quadratic function $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} + \mathbf{b}^T \mathbf{x} + c$ (*)
- -> set of linear equations: $\mathbf{A} \mathbf{x} = \mathbf{b}$ (\mathbf{A} = positive definite and symmetric)
- Solution: $\mathbf{x}^* = \mathbf{A}^{-1} \mathbf{b}$ ($\mathbf{A} \equiv \mathbf{H}$)
- Given the matrix \mathbf{A} , a set of nonzero vectors $\mathbf{s}(0), \mathbf{s}(1), \dots$, (up to $\dim(\mathbf{A})$) is \mathbf{A} -conjugate (i.e., non-interfering with each other in the context of \mathbf{A}) if: $\mathbf{s}(i)^T \mathbf{A} \mathbf{s}(j) = 0$. (for $\mathbf{A} = \mathbf{I}$, conjugacy = orthogonality).
- Find conjugate directions (for minibatch) without computing $\mathbf{H}^{-1}(t)$, e.g. via Polak-Ribière method:

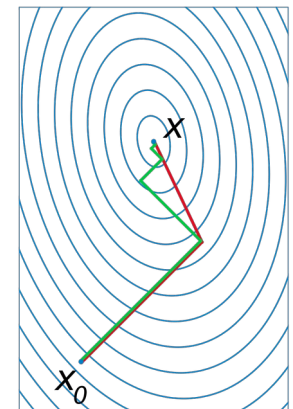
$$\beta(t) = \frac{(\mathbf{g}(t) - \mathbf{g}(t-1))^T \mathbf{g}(t)}{\mathbf{g}(t-1)^T \cdot \mathbf{g}(t-1)}$$

$$\rho(t) = -\mathbf{g}(t) + \beta(t) \cdot \rho(t-1)$$

$$\mathbf{w}(t+1) = \mathbf{w}(t) + \alpha^*(t) \cdot \rho(t)$$

(Haykin, 2009)

Computed analytically



Regularization

Goal – minimize overfitting

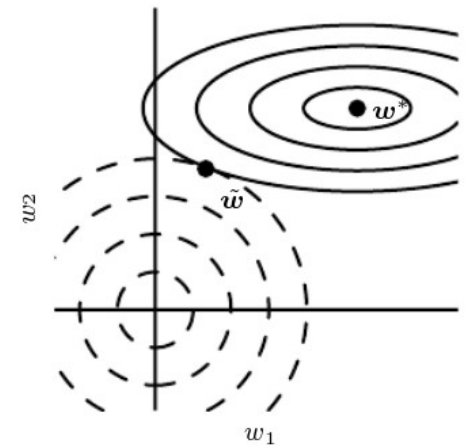
Risk function: $R(\mathbf{w}) = E(\mathbf{w}) + \lambda C(\mathbf{w})$ [performance + complexity]

- **Explicit:**

$$L_1(\mathbf{w}) = \epsilon \sum_i |w_i| \quad L_2(\mathbf{w}) = \frac{\epsilon}{2} \|\mathbf{w}\|^2$$

- **Implicit:**

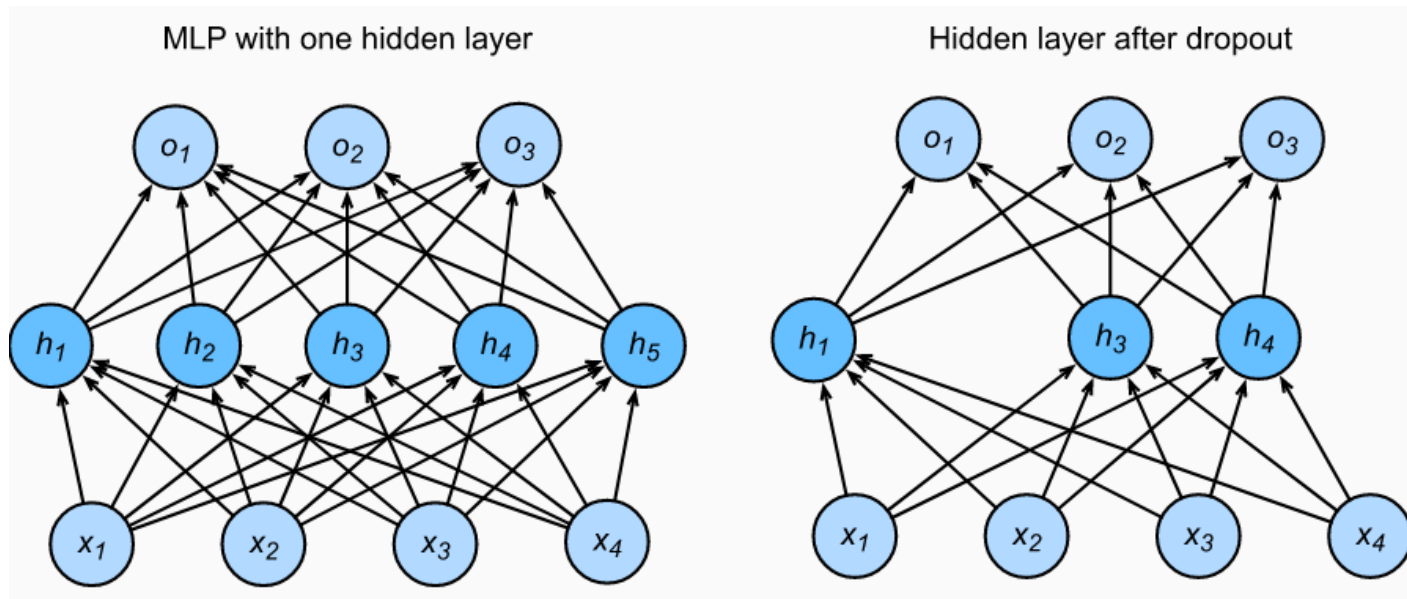
- weight decay $w_l^{\text{new}}(t) = \epsilon \cdot w_l^{\text{new}}(t)$, $0 \ll \epsilon < 1$,
 - ...leads to L2-regul.
 - dropout (Hinton, 2012): random turning off neurons during training
- data augmentation – increasing the size of the training set, e.g. by elastic distortions



(Goodfellow, 2015)

Dropout

- Applied only during training
- Helps to avoid overfitting
- Free parameter = number of (randomly) dropped units



(Zhang et al, 2019)

AdaGrad algorithm

- Back to first-order methods
- Introduces the variable that accumulates gradient variance (vector)

$$\mathbf{g}(t) = \nabla E(\mathbf{w}(t)) \quad \mathbf{s}(t) = \mathbf{s}(t-1) + \mathbf{g}^2(t) \quad \Delta \mathbf{w}(t) = -\frac{\alpha}{\sqrt{\mathbf{s}(t)} + \epsilon} \cdot \mathbf{g}(t)$$

- decreases the learning rate dynamically on per-coordinate basis
- Uses the magnitude of the gradient as a means of adjusting how quickly progress is achieved – coordinates with large gradients are compensated with a smaller learning rate.
- first-order method (the gradient can be a useful proxy)
- On deep learning problems AdaGrad can sometimes be too aggressive in reducing learning rates. Mitigating strategies exist.

RMSprop

- Decouples rate scheduling from coordinate-adaptive learning rates

$$\mathbf{s}(t) = \gamma \mathbf{s}(t-1) + (1 - \gamma) \mathbf{g}^2(t) \quad \Delta \mathbf{w}(t) = - \frac{\alpha}{\sqrt{\mathbf{s}(t) + \epsilon}} \cdot \mathbf{g}(t) \quad \epsilon = 10^{-6}$$

- Coefficient γ determines how long the history is when adjusting the per-coordinate scale.
- RMSprop shares with momentum the leaky averaging. However, RMSProp uses the technique to adjust the coefficient-wise preconditioner (for reducing the condition number).

AdaDelta

- Yet another variant of AdaGrad: it decreases the amount by which the learning rate is adaptive to coordinates
- It does not literally have a learning rate since it uses the amount of change itself as calibration for future change:

$$s(t) = \rho s(t-1) + (1-\rho) \mathbf{g}^2(t) \qquad \mathbf{g}'(t) = \sqrt{\frac{\Delta \mathbf{w}(t-1) + \epsilon}{s(t) + \epsilon}} \cdot \mathbf{g}(t)$$

$$\mathbf{w}(t) = \mathbf{w}(t-1) - \mathbf{g}'(t)$$

$$\Delta \mathbf{w}(t) = \rho \Delta \mathbf{w}(t-1) + (1-\rho) \mathbf{g}'^2(t)$$

Adam algorithm

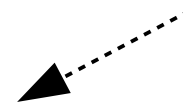
- Combines 3 preceeding techniques into one efficient algorithm
- uses leaky averaging to obtain an estimate of both the momentum and also the second moment of the gradient

$$\begin{aligned} \mathbf{v}(t) &= \beta_1 \mathbf{v}(t-1) + (1 - \beta_1) \mathbf{g}(t) & \mathbf{v}'(t) &= \mathbf{v}(t) / (1 - \beta_1^t) & \beta_1 &= 0.9 \\ s(t) &= \beta_2 s(t-1) + (1 - \beta_2) \mathbf{g}^2(t) & s'(t) &= s(t) / (1 - \beta_2^t) & \beta_2 &= 0.999 \end{aligned}$$

$$\Delta \mathbf{w}(t) = -\frac{\alpha}{\sqrt{s'(t)} + \epsilon} \cdot \mathbf{v}'(t) \quad (\text{Kingma \& Ba, 2014})$$

- Still, gradients with significant variance may hinder convergence ($s(t)$ can blow up)
- Yogi algorithm addresses this: $s(t) = s(t-1) + (1 - \beta_2) (\mathbf{g}^2(t) - s(t-1))$

(Zaheer et al, 2018)



$$s(t) = s(t-1) + (1 - \beta_2) \mathbf{g}^2(t) \cdot \text{sgn}(\mathbf{g}^2(t) - s(t-1))$$

Natural gradient learning

- uses **Fisher information**: a positive semidefinite matrix ($|\mathbf{w}| \times |\mathbf{w}|$), defines a Riemannian metric (\rightarrow information geometry) (Amari, 1998)
- look at p.d.f. via $KL(p(\mathbf{x}; \mathbf{w}) \parallel p(\mathbf{x}; \mathbf{w} + \Delta \mathbf{w})) = \dots \approx \frac{1}{2} (\Delta \mathbf{w})^T \mathbf{F} \Delta \mathbf{w}$
- matrix \mathbf{F} is the negative expected Hessian of $\log p(\mathbf{x}; \mathbf{w})$
- $\Delta \mathbf{w}^* = \arg \min_{\Delta \mathbf{w}} \{L(\mathbf{w} + \Delta \mathbf{w}) + \lambda \cdot KL(p(\mathbf{x}; \mathbf{w}) \parallel p(\mathbf{x}; \mathbf{w} + \Delta \mathbf{w})) - c\}$
- $\Delta \mathbf{w}(t) = -\alpha \mathbf{F}^{-1}(\mathbf{w}(t)) \cdot \mathbf{g}(t)$, i.e. **natural gradient** $\mathbf{g}_{nat}(t) = \mathbf{F}^{-1}(\mathbf{w}) \mathbf{g}(t)$
- can be interpreted as curvature of the log-likelihood function p
- in NG descent, we control movement in **prediction space** (rather than parameter space)
- efficient approximations of \mathbf{F}^{-1} possible (Amari et al, 2019)

Summary

- NN learning and classical optimization have different objectives
- NN goal = minimize generalization error (sometimes using surrogate loss functions)
- Various known problems hinder first-order gradient methods
- Second-order methods provide more information but are much more costly
- Earlier methods focused on approximating the Hessian
- Recent methods focus only on gradients and its adaptive versions
- Natural gradient learning (2nd order) uses Riemannian metric
- Further improvements possible (found useful in deep learning), to be mentioned later