# Statistical Machine Learning (BE4M33SSU) Lecture 7: Artificial Neural Networks

Jan Drchal

Czech Technical University in Prague Faculty of Electrical Engineering Department of Computer Science

### **Overview**

р

2/49

Topics covered in the lecture:

- Neuron types
- Layers
- Loss functions
- Computing loss gradients via backpropagation
- Learning neural networks
- Regularization

#### **McCulloch-Pitts Perceptron**



$$\boldsymbol{x} = (x_1, x_2, \dots, x_n)^T \in \mathbb{R}^n$$
$$\hat{y} \in \{-1, 1\}$$
$$\boldsymbol{w} = (w_1, w_2, \dots, w_n)^T \in \mathbb{R}^n$$
$$b \in \mathbb{R}$$

$$s \in \mathbb{R}$$

$$f(s) = \begin{cases} -1 & \text{if } s < 0\\ 1 & \text{else} \end{cases}$$

input
output (activity)
weights
bias (threshold)

inner potential

activation function

$$\hat{y} = f(s) = f\left(\sum_{i=1}^{n} w_i x_i + b\right) = f\left(\boldsymbol{w} \cdot \boldsymbol{x} + b\right)$$



#### **McCulloch-Pitts Perceptron: Treating Bias**



• Treat bias as an extra fixed input  $x_0 = 1$  and weighted  $w_0 = b$ :

$$\hat{y} = f(\boldsymbol{w} \cdot \boldsymbol{x} + b) = f(\boldsymbol{w} \cdot \boldsymbol{x} + w_0 \cdot 1) = f(\boldsymbol{w}' \cdot \boldsymbol{x}')$$

• 
$$\boldsymbol{x'} = (x_0, x_1, \dots, x_n)^T \in \mathbb{R}^{n+1}$$
  
•  $\boldsymbol{w'} = (w_0, w_1, \dots, w_n)^T \in \mathbb{R}^{n+1}$ 

igstarrow Unless otherwise noted we will use x, w instead of x', w'



#### **Activation Functions**



2 m p 5/49

#### **Linear Neuron**



Single neuron with linear activation function  $\equiv$  linear regression:

$$\hat{y} = s = \boldsymbol{x} \cdot \boldsymbol{w}, \quad \hat{y}, s \in \mathbb{R}$$

• Inputs: 
$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} \boldsymbol{x}_1^T \\ \vdots \\ \boldsymbol{x}_m^T \end{pmatrix}$$

$$lacksim$$
 Targets:  $oldsymbol{y} = ig(y_1,\,\ldots,\,y_mig)^T$ ,  $y_i \in \mathbb{R}^m$ 

• Outputs: 
$$oldsymbol{\hat{y}} = ig(\hat{y}_1,\,\ldots,\,\hat{y}_mig)^T$$
,  $\hat{y}_i \in \mathbb{R}^m$ 

• For the whole dataset we get:

$$oldsymbol{\hat{y}} = \mathbf{X}oldsymbol{w}, oldsymbol{\hat{y}} \in \mathbb{R}^m$$

#### Linear Neuron: Maximum Likelihood Estimation

• Assumption: data are Gaussian distributed with mean  $x_i \cdot w$  and variance  $\sigma^2$ :

$$y_i \sim \mathcal{N}\left(\boldsymbol{x}_i \cdot \boldsymbol{w}, \sigma^2\right) = \boldsymbol{x}_i \cdot \boldsymbol{w} + \mathcal{N}\left(0, \sigma^2\right)$$

Likelihood for i.i.d. data:

$$p(\boldsymbol{y}|\mathbf{X}, \boldsymbol{w}, \sigma) = \prod_{i=1}^{m} p(y_i | \boldsymbol{x}_i, \boldsymbol{w}, \sigma) = \prod_{i=1}^{m} (2\pi\sigma^2)^{-\frac{1}{2}} e^{-\frac{1}{2\sigma^2}(y_i - \boldsymbol{w} \cdot \boldsymbol{x}_i)^2} =$$
$$= (2\pi\sigma^2)^{-\frac{m}{2}} e^{-\frac{1}{2\sigma^2}\sum_{i=1}^{m}(y_i - \boldsymbol{w} \cdot \boldsymbol{x}_i)^2} =$$
$$= (2\pi\sigma^2)^{-\frac{m}{2}} e^{-\frac{1}{2\sigma^2}(\boldsymbol{y} - \mathbf{X}\boldsymbol{w})^T(\boldsymbol{y} - \mathbf{X}\boldsymbol{w})}$$

• Negative Log Likelihood:

$$\mathcal{L}(\boldsymbol{w}) = \frac{m}{2} \log \left(2\pi\sigma^2\right) + \frac{1}{2\sigma^2} \left(\boldsymbol{y} - \mathbf{X}\boldsymbol{w}\right)^T \left(\boldsymbol{y} - \mathbf{X}\boldsymbol{w}\right)$$



#### Linear Neuron: Maximum Likelihood Estimation (contd.)



- Note that  $\sum_{i=1}^{m} (y_i \boldsymbol{w} \cdot \boldsymbol{x}_i)^2 = (\boldsymbol{y} \mathbf{X}\boldsymbol{w})^T (\boldsymbol{y} \mathbf{X}\boldsymbol{w})$  is the sum-of-squares or squared error (SE)
- Minimization of  $\mathcal{L}\left(oldsymbol{w}
  ight)\equiv$  least squares estimaton
- Solving  $\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}} = 0$  we get  $\boldsymbol{w}^* = \left( \mathbf{X}^T \mathbf{X} \right)^{-1} \mathbf{X}^T \boldsymbol{y}$  (see seminar)
- Note  $\hat{\mathbf{X}}^{=} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$  is called the *Moore-Penrose pseudo-inverse*: if  $\mathbf{X}$  is square and invertible then  $\hat{\mathbf{X}}^{=} \mathbf{X}^{-1}$  (use  $(\mathbf{AB})^{-1} = \mathbf{B}^{-1} \mathbf{A}^{-1}$ )



### **Sigmoid and Probability**

• Denote: 
$$\hat{y} = \sigma(s), \ \hat{y} \in (0,1)$$

Sigmoid output can represent a parameter of the Bernoulli distribution:

$$p(y|\hat{y}) = \text{Ber}(y|\hat{y}) = \hat{y}^y (1-\hat{y})^{1-y} = \begin{cases} \hat{y} & \text{for } y = 1\\ 1-\hat{y} & \text{for } y = 0 \end{cases}$$

Motivation: log-odds linear model (see AE4B33RPZ)

• Binary classifier: 
$$h(\hat{y}) = \begin{cases} 1 & \text{if } \hat{y} > \frac{1}{2} \\ 0 & \text{else} \end{cases}$$
  
 $\sigma(s) = \frac{1}{1 + e^{-s}}$   
 $\frac{1}{2}$ 



### **Logistic Regression**



• MCP neuron using sigmoid activation function  $\equiv$  **logistic regression**:

$$\hat{y} = \sigma(\boldsymbol{w} \cdot \boldsymbol{x}), \; \hat{y} \in (0, 1)$$

• Inputs: 
$$\mathbf{X} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{mn} \end{pmatrix} = \begin{pmatrix} \boldsymbol{x}_1^T \\ \vdots \\ \boldsymbol{x}_m^T \end{pmatrix}$$

$$lacksim$$
 Target class:  $oldsymbol{y} = ig(y_1,\,\ldots,\,y_mig)^T$ ,  $y_i \in \{0,1\}$ 

• Output class: 
$$\hat{\boldsymbol{y}} = \left(\hat{y}_1, \ldots, \hat{y}_m\right)^T$$
,  $\hat{y}_i \in (0, 1)$ 

### **Cross-Entropy**



р

11/49

• Likelihood, for the logistic regression:

$$p(\mathbf{y}|\mathbf{w}, \mathbf{X}) = \prod_{i=1}^{m} \text{Ber}(y_i | \hat{y}_i) = \prod_{i=1}^{m} \hat{y}_i^{y_i} (1 - \hat{y}_i)^{1 - y_i}$$

$$\mathcal{L}(\boldsymbol{w}) = -\sum_{i=1}^{m} \left[ y_i \log \hat{y}_i + (1 - y_i) \log \left(1 - \hat{y}_i\right) \right]$$



This loss function is called the cross-entropy

### Maximum Likelihood Estimation

- Maximum Likelihood Estimation:  $w^* = \underset{w}{\operatorname{argmin}} \mathcal{L}(\boldsymbol{w})$
- Derivative of the loss w.r.t. to the sigmoid argument:  $\frac{\partial \mathcal{L}}{\partial s_i} = \hat{y}_i y_i$
- Gradient w.r.t. logistic regression parameters:

$$\frac{\partial \mathcal{L}}{\partial \boldsymbol{w}} = \sum_{i=1}^{m} \frac{\partial \mathcal{L}}{\partial s_i} \cdot \frac{\partial s_i}{\partial \boldsymbol{w}} = \sum_{i=1}^{m} \boldsymbol{x}_i (\hat{y}_i - y_i) = \mathbf{X}^T (\boldsymbol{\hat{y}} - \boldsymbol{y})$$

•  $\frac{\partial \mathcal{L}}{\partial w} = 0$  has no analytical solution  $\implies$  use numerical methods

• Hessian: 
$$\frac{\partial^2 \mathcal{L}}{\partial w^2} = \sum_{i=1}^m \hat{y}_i (1 - \hat{y}_i) \boldsymbol{x}_i \cdot \boldsymbol{x}_i = \mathbf{X}^T \mathbf{S} \mathbf{X},$$
  
where  $\mathbf{S} \triangleq \operatorname{diag}(\hat{y}_i (1 - \hat{y}_i))$ 

 Hessian is positive definite, hence the loss is convex and has unique global minimum (see AE4B33RPZ)



# **Rectified Linear Unit (ReLU)**



• Definition  $f(s) = \max(0, s)$ 

### Fast to compute

- Helps with vanishing gradients problem: the gradient is constant for s > 0, while for sigmoid-like activations it becomes increasingly small
- igstarrow Leads to sparse representations: s < 0 turns the neuron completely off
- Unbounded: use regularization to prevent numerical problems
- ullet Might block gradient propagation ightarrow dead units ightarrow Leaky ReLU
- Satisfies the universal approximation property  $\operatorname{ReLU}_{f(s) = \max(0, s)}$

## Multilayer Perceptron (MLP)



m p

14/49

Feed-forward ANN

Fully-connected layers

MLP for regression would typically use linear output layer

### **Recurrent Neural Network (RNN)**





Fully-Connected Recurrent Neural Network (FRNN)

- Both inputs and outputs are sequences
- Feedback connections  $\rightarrow$  memory

#### **Modular and Hierarchical Architectures**



m p

16/49

- Layers can be organized in modules
- Hierarchies of modules
- Module reuse

#### **Linear Layer**

$$igle$$
 Output  $k$ :  $\hat{y}_k = oldsymbol{x} \cdot oldsymbol{w}_k$ ,  $k = 0, 1, \dots, K$ 

• All outputs using weight matrix  $\mathbf{W}$ :  $\hat{y} = x^T \mathbf{W}$ 

• Multiple samples:  $\hat{\mathbf{Y}} = \mathbf{X}\mathbf{W}$ 

$$\mathbf{W} = \begin{pmatrix} \boldsymbol{w}_1^T \\ \vdots \\ \boldsymbol{w}_K^T \end{pmatrix}^T = \begin{pmatrix} w_{01} & \dots & w_{0K} \\ \vdots & \ddots & \vdots \\ w_{n1} & \dots & w_{nK} \end{pmatrix} \qquad \begin{array}{c} x_1 \\ x_2 \\ \vdots \\ x_n \\ w_{nK} \end{array}$$

$$\mathbf{X} = \begin{pmatrix} \boldsymbol{x}_1^T \\ \vdots \\ \boldsymbol{x}_m^T \end{pmatrix} = \begin{pmatrix} 1 & x_{11} & \dots & x_{1n} \\ 1 & \vdots & \ddots & \vdots \\ 1 & x_{m1} & \dots & x_{mn} \end{pmatrix} \qquad \hat{\mathbf{Y}} = \begin{pmatrix} \hat{\boldsymbol{y}}_1^T \\ \vdots \\ \boldsymbol{y}_m^T \end{pmatrix} = \begin{pmatrix} \hat{y}_{11} & \dots & \hat{y}_{1K} \\ \vdots & \ddots & \vdots \\ \hat{y}_{m1} & \dots & \hat{y}_{mK} \end{pmatrix}$$



### Softmax Layer





Multinominal classification

• Definition:  $\sigma_k(s) \triangleq \frac{e^{s_k}}{\sum_{c=1}^{K} e^{s_c}}$ , where K is the number of classes

- Softmax represents a probability distribution:  $0 \le \sigma_k \le 1$  for  $k \in \{1 \dots K\}$  and  $\sum_{c=1}^{K} \sigma_c = 1$
- Describes class membership probabilities:  $p(y = k | s) = \sigma_k(s)$

### Softmax Layer (contd.)



• One-hot encoding for sample i and class k: let  $y_{ik} = \mathbb{I}\{y_i = k\}$  and  $\hat{y}_{ik} = \mathbb{I}\{\hat{y}_i = k\}$ 

• Likelihood:

$$p(\boldsymbol{y}|\boldsymbol{w}, \mathbf{X}) = \prod_{i=1}^{m} \prod_{c=1}^{K} \hat{y}_{ic}^{y_{ic}}$$

• Negative Log Likelihood:

$$\mathcal{L}(\boldsymbol{w}) = -\sum_{i=1}^{m} \sum_{c=1}^{K} \mathbb{I}\{y_i = c\} \log(\hat{y}_{ic})$$

Again the cross-entropy

See seminar for the gradient



#### **Multinominal Logistic Regression**



• linear layer + softmax layer = multinominal logistic regression:

$$\hat{y}_k = \sigma_k(\boldsymbol{w}_k \cdot \boldsymbol{x})$$



• Classifier: 
$$h(\boldsymbol{x}, \mathbf{W}) = \operatorname*{argmax}_{k} \hat{y}_{k}$$

### **Loss Functions: Summary**



problem	suggested loss function
binary classification	cross-entropy
	$-\sum_{i=1}^{m} \left[ y_i \log \hat{y}_i + (1 - y_i) \log \left( 1 - \hat{y}_i \right) \right]$
multinominal classification	multinominal cross-entropy
	$-\sum_{i=1}^{m} \sum_{c=1}^{K} \mathbb{I}\{y_i = c\} \log(\hat{y}_{ic})$
regression	squared error
	$\sum_{i=1}^{m} \left( y_i - \hat{y}_i \right)^2$
multi-output regression	squared error
	$\sum_{i=1}^{m} \sum_{c=1}^{K} (y_{ic} - \hat{y}_{ic})^2$

• Mean w.r.t. to m is often used

### **Backpropagation Overview**

A method to compute gradient of the *loss function* with respect to its parameters

p

22/49

- Here, we present the "modular" backpropagation (see Nando de Freitas' Machine Learning course: https://www.cs.ox.ac.uk/people/ nando.defreitas/machinelearning/)
- Let us use multinominal logistic regression as an example



#### **Backpropagation: the Loss Function**

• The loss function is the multinominal cross-entropy in this case:

$$\mathcal{L}(\boldsymbol{w}) = -\sum_{i=1}^{m} \sum_{c=1}^{K} \mathbb{I}\{y_i = c\} \log \left(\frac{\exp\left(\boldsymbol{x}_i \cdot \boldsymbol{w}_c\right)}{\sum_{k=1}^{K} \exp\left(\boldsymbol{x}_i \cdot \boldsymbol{w}_k\right)}\right)$$





### **Backpropagation Based on Modules**



- Computation of  $abla \mathcal{L}(m{w})$  involves repetitive use of the *chain rule*
- We can make things simpler by divide and conquer approach
- Divide to simplest possible modules (these can be later combined into complex hierarchies)
- Represent even the loss function as a module
- Passing messages



#### **Backpropagation: Backward Pass Message**

• Let  $\delta^l = \frac{\partial \mathcal{L}}{\partial z^l}$  be the sensitivity of the loss to the module output for layer l, then:

$$\delta_i^l = \frac{\partial \mathcal{L}}{\partial z_i^l} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial z_i^l} = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial z_i^l}$$

We need to know how to compute derivatives of outputs w.r.t. inputs only!





#### **Backpropagation:** Parameters



 Similarly if the module has parameters we want to know how the loss changes w.r.t. them:

$$\frac{\partial \mathcal{L}}{\partial w_i^l} = \sum_j \frac{\partial \mathcal{L}}{\partial z_j^{l+1}} \cdot \frac{\partial z_j^{l+1}}{\partial w_i^l} = \sum_j \delta_j^{l+1} \frac{\partial z_j^{l+1}}{\partial w_i^l}$$

• Derivatives of module outputs w.r.t. to the parameters are all we need



### **Backpropagation: Steps**



• So for each module we need only to specify these three messages:

forward:  $z^{l+1} = f(z^l)$ backward:  $\frac{\partial z^{l+1}}{\partial z^l}$ parameter (optional):  $\frac{\partial z^{l+1}}{\partial w^l}$ 





#### **Example: Linear Layer**







### **Example: Squared Error**

• forward: 
$$z^{l+1} = \frac{1}{2} \sum_{i=1}^{K} (y_i - z_i^l)^2, \quad i \in \{1, \dots, n\}$$

• **backward**: 
$$\frac{\partial z_i^{l+1}}{\partial z^l} = y_i - z_i^l$$
,  $i \in \{1, \dots, n\}$ 

#### **Gradient Descent**



• Task: find parameters which minimize loss over the training dataset:

$$oldsymbol{ heta}^* = \operatorname*{argmin}_{oldsymbol{ heta}} \mathcal{L}(oldsymbol{ heta})$$

where  $\theta$  is a set of all parameters defining the ANN (e.g., all weight matrices)

Gradient descent:  $\theta^{(t+1)} = \theta^{(t)} - \eta^{(t)} \nabla \mathcal{L}(\theta^{(t)})$ where  $\eta^{(t)} > 0$  is the **learning rate** or **step size** at iteration t





### Batch, Online and Mini-Batch Learning

### When to update weights?

- (Full) Batch learning: after all patterns are used (epoch)
  - inefficient for redundant datasets
- **Online learning**: after each training pattern
  - noise can help overcome local minima but can also harm the convergence in the final stages while fine-tuning
  - Stochastic Gradient Descent (SGD) does this
  - convergence almost surely to local minimum when  $\eta^{(t)}$  decreases appropriately in time
- Mini-batch learning: after a small sample of training patterns

#### Momentum



• Simulate inertia to overcome plateaus in the error landscape:

$$\boldsymbol{v}^{(t+1)} = \mu \boldsymbol{v}^{(t)} - \eta^{(t)} \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)})$$
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \boldsymbol{v}^{(t+1)}$$

where  $\mu \in [0,1]$  is the momentum parameter

- Momentum damps oscillations in directions of high curvature
- It builds velocity in directions with consistent (possibly small) gradient



#### **Nesterov's Momentum**



Different approach by Nesterov (1983, convex optimization):

$$\boldsymbol{v}^{(t+1)} = \mu \boldsymbol{v}^{(t)} - \eta^{(t)} \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)} + \mu \boldsymbol{v}^{(t)})$$
$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} + \boldsymbol{v}^{(t+1)}$$

- While classic momentum *corrects* the velocity using gradient at  $\theta^{(t)}$ , Nesterov uses  $\theta^{(t)} + \mu v^{(t)}$  which is similar to  $\theta^{(t+1)}$
- Stronger theoretical convergence guarantees for convex functions
- Slightly better in practice
- For more details see Sutskever et al.: On the importance of initialization and momentum in deep learning, 2013



### **Annealing the Learning Rate**



- Decrease the learning rate in the course of optimization.
- **Step decay**: reduce the learning rate by a factor (e.g.,  $\frac{1}{2}$ ) every few iterations
- **Exponential decay**: set  $\eta_t = \eta_0 e^{-kt}$  for the iteration t
- Hyperbolic decay: set  $\eta_t = \frac{\eta_0}{1+kt}$

### **Resilient Propagation (Rprop)**

Motivation: a magnitude of gradient differs a lot for different Parameters

35/49

 Rprop (Riedmiller and Braun, 1992) does not use gradient value - the step size for each weight is adapted using its sign, only

1 for each  $\theta_i$ 2 if  $\frac{\partial \mathcal{L}^{(t-1)}}{\partial \theta_i} \cdot \frac{\partial \mathcal{L}^{(t)}}{\partial \theta_i} > 0$ 3  $\Delta_i^{(t)} = \min\left(\Delta_i^{(t-1)} \cdot \eta^+, \Delta_{max}\right)$ 4 elseif  $\frac{\partial \mathcal{L}^{(t-1)}}{\partial \theta_i} \cdot \frac{\partial \mathcal{L}^{(t)}}{\partial \theta_i} < 0$ 5  $\Delta_i^{(t)} = \max\left(\Delta_i^{(t-1)} \cdot \eta^-, \Delta_{min}\right)$ 6  $\theta_i^{(t+1)} = \theta_i^{(t)} - \operatorname{sign}\left(\frac{\partial \mathcal{L}^{(t)}}{\partial \theta_i}\right) \cdot \Delta_i^{(t)}$ 

where the step size  $\Delta_i^{(t)} \in [\Delta_{min}, \Delta_{max}]$  and  $0 < \eta^- < 1 < \eta^+$ 

Typically order of magnitude faster than basic Gradient Descent
 Does not work well for mini-batches

Igel and Hüsken: Improving the Rprop Learning Algorithm, 2000

### Adagrad

Adaptive Gradient method

• Idea: reduce learning rates for parameters having high values of gradient

р

36/49

$$\begin{split} g_i^{(t+1)} &= g_i^{(t)} + \left(\frac{\partial \mathcal{L}}{\partial \theta_i^{(t)}}\right)^2 \\ \theta_i^{(t+1)} &= \theta_i^{(t)} - \frac{\eta}{\sqrt{g_i^{(t+1)}} + \epsilon} \cdot \frac{\partial \mathcal{L}}{\partial \theta_i^{(t)}} \end{split}$$

g<sub>i</sub> accumulates squared partial derivatives w.r.t. to the parameter θ<sub>i</sub>
 ϵ is a small positive number to prevent division by zero
 Weakness: ever increasing g<sub>i</sub> leads to slow convergence eventually


### RMSProp

Similar to Adagrad but employs a moving average:

$$g_i^{(t+1)} = \gamma g_i^{(t)} + (1 - \gamma) \left(\frac{\partial \mathcal{L}}{\partial \theta_i^{(t)}}\right)^2$$

•  $\gamma$  is a *decay* parameter (typical value  $\gamma = 0.9$ )

Unlike for Adagrad updates do not get infinitesimally small

### Adam (Adaptive Moment Estimation)

• Kingma and Ba: Adam: A Method for Stochastic Optimization, 2014

38/49

$$\begin{split} m_{i}^{(t+1)} &= \beta_{1} m_{i}^{(t)} + (1 - \beta_{1}) \frac{\partial \mathcal{L}}{\partial \theta_{i}^{(t)}} & \hat{m}_{i}^{(t+1)} = \frac{m_{i}^{(t+1)}}{1 - \beta_{1} m_{i}^{(t)}} \\ v_{i}^{(t+1)} &= \beta_{2} v_{i}^{(t)} + (1 - \beta_{2}) \left( \frac{\partial \mathcal{L}}{\partial \theta_{i}^{(t)}} \right)^{2} & \hat{v}_{i}^{(t+1)} = \frac{v_{i}^{(t+1)}}{1 - \beta_{2} v_{i}^{(t)}} \\ \theta_{i}^{(t+1)} &= \theta_{i}^{(t)} - \eta \frac{\hat{m}_{i}^{(t+1)}}{\sqrt{\hat{v}_{i}^{(t+1)}} + \epsilon} \end{split}$$

- $m_i$  and  $v_i$  are first and second raw moment estimates of gradient (mean and uncentered variance)
- $\hat{m}_i$  and  $\hat{v}_i$  are corrections of  $m_i$  and  $v_i$  as these are biased to zero due to initialization  $m_i = v_i = 0$

 $\beta_1$  and  $\beta_2$  are *decay* parameters

### **Second Order Methods**



Newton's method for optimization:

$$\boldsymbol{\theta}^{(t+1)} = \boldsymbol{\theta}^{(t)} - \left(\nabla^2 \mathcal{L}(\boldsymbol{\theta}^{(t)})\right)^{-1} \nabla \mathcal{L}(\boldsymbol{\theta}^{(t)})$$

- No parameters needed
- Need to compute the Hessian matrix and invert it!
- Approximations, e.g., L-BFGS are not suitable for mini-batches



### Regularization



40/49

- How to deal with overfitting?
  - get more data
  - find simpler model, search for optimal architecture, e.g., number, type and size of layers
  - use *regularization*
- Most types of regularization are based on penalties for model complexity
- Bayesian point of view: introduce prior distribution on model parameters

### **L2** Regularization

- Recall the solution for the linear regression w<sup>\*</sup> = (X<sup>T</sup>X)<sup>-1</sup> X<sup>T</sup>y
   What if X<sup>T</sup>X has no inverse?
- We can modify the solution by adding a small element to the diagonal:

$$\boldsymbol{w}^* = \left( \mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X}^T \boldsymbol{y}, \quad \lambda > 0$$

- It turns out that this approach no only helps with inverting  $\mathbf{X}^T \mathbf{X}$  but it also improves model generalization
- It is the solution of the *regularized* loss function:

$$\mathcal{L}(\boldsymbol{w}) = \left(\boldsymbol{y} - \mathbf{X}\boldsymbol{w}\right)^{T} \left(\boldsymbol{y} - \mathbf{X}\boldsymbol{w}\right) + \lambda \boldsymbol{w}^{T}\boldsymbol{w}$$

this one is called the **L2 regularization**, see seminar for the derivation • The term  $\lambda w^T w = \lambda ||w||_2^2$  minimizes the size of the weight vector • Note that we omit bias in  $\lambda w^T w$ 



# L2 Regularization (contd.)



• L2 regularization is also called the ridge regression or the weight decay



 $w_1$ 

# L2 Regularization as Gaussian Prior



• Recall the likelihood:

$$p(\boldsymbol{y}|\boldsymbol{w}, \mathbf{X}) = (2\pi\sigma^2)^{-\frac{m}{2}} e^{-\frac{1}{2\sigma^2}(\boldsymbol{y} - \mathbf{X}\boldsymbol{w})^T(\boldsymbol{y} - \mathbf{X}\boldsymbol{w})}$$

• Define a Gaussian prior with zero mean and variance  $\sigma_0^2$  for the parameters:

$$p(\boldsymbol{w}) = \left(2\pi\sigma_0^2\right)^{-\frac{1}{2}} e^{-\frac{1}{2\sigma_0^2}\boldsymbol{w}^T\boldsymbol{w}}$$

Then the posterior is:

$$p(\boldsymbol{w}|\boldsymbol{y}, \mathbf{X}) = \frac{p(\boldsymbol{y}|\boldsymbol{w}, \mathbf{X}) \cdot p(\boldsymbol{w})}{p(\boldsymbol{y}|\mathbf{X})}$$

The denominator does not depend on the parameters w:

$$p(\boldsymbol{w}|\boldsymbol{y}, \mathbf{X}) \propto p(\boldsymbol{y}|\boldsymbol{w}, \mathbf{X}) \cdot p(\boldsymbol{w})$$

#### **MAP Estimate**

Maximizing p(w|y, X) gives us the Maximum a posteriori (MAP) estimate:

$$\boldsymbol{w}_{MAP} = \operatorname*{argmax}_{\boldsymbol{w}} p(\boldsymbol{w}|\boldsymbol{y}, \mathbf{X}) = \operatorname*{argmin}_{\boldsymbol{w}} \left( -\log p(\boldsymbol{w}|\boldsymbol{y}, \mathbf{X}) \right)$$

where

$$-\log p(\boldsymbol{w}|\boldsymbol{y}, \mathbf{X}) = \frac{1}{2\sigma^2} (\boldsymbol{y} - \mathbf{X}\boldsymbol{w})^T (\boldsymbol{y} - \mathbf{X}\boldsymbol{w}) + \frac{1}{2\sigma_0^2} \boldsymbol{w}^T \boldsymbol{w} + C$$

• We can omit C, define  $\lambda = \frac{\sigma^2}{\sigma_0^2}$  and minimize the loss function we already know:

$$\mathcal{L}(\boldsymbol{w}) = (\boldsymbol{y} - \mathbf{X}\boldsymbol{w})^T (\boldsymbol{y} - \mathbf{X}\boldsymbol{w}) + \lambda \boldsymbol{w}^T \boldsymbol{w}$$



### Weight Decay Discussion



- Having zero mean Gaussian prior keeps the weights smaller
- Weight decay is widely used for most types of layers in ANNs
- Intuition: sigmoid-like neurons kept near zero potential (via small weights) behave similarly to linear neurons
- The same works for other models, e.g., polynomial regression
- ullet  $\lambda$  is usually set using cross validation



# Dropout

- Idea: average many neural networks, share weights to make this computationally feasible
- For each training example omit a unit with probability p (often p = 0.5)
- This is like sampling from  $2^U$  networks where U is the number of units
- Typically only a small amount of  $2^U$  networks is actually sampled

Prevents coadaptation of feature detectors



(a) Standard Neural Net (b) After applying dropout. Srivastava et al.: A Simple Way to Prevent Neural Networks from Overfitting, 2014



# **Dropout (contd.)**

- How to make predictions with networks using dropout?
- Scale outputs (output weights) of all affected units by p
- For a linear unit taking inputs from the dropout layer we have:

$$s = \sum_{i=1}^{n} w_i \delta_i x_i, \quad p(\delta_i | p) = \text{Ber}(\delta_i | p)$$
$$\mathbb{E}(s) = \sum_{i=1}^{n} w_i \mathbb{E}(\delta_i) x_i = p \sum_{i=1}^{n} w_i x_i$$

where the expectation is computed over all  $2^n$  configurations

- For general neural networks we still get a good approximation of the expectation when scaling by p
- See Baldi and Sadowski: *The Dropout Learning Algorithm*, 2014
- $\bullet$  When used for inputs higher values of p are suggested



### **Other Regularization Approaches**



- ullet L1 regularization: sum absolute values, i.e., use  $\lambda \left\| oldsymbol{w} 
  ight\|_1$
- Early stopping: start with small weights, stop when validation loss starts to grow
- Randomize inputs: same as the weight decay for linear neurons
- Noisy weights
- DropConnect: connection-based dropout
- Other weight sharing approaches: Convolutional Neural Networks
- Model averaging

### **Next Lecture**

- Deep Neural Networks
- Convolutional Neural Networks
- Autoencoders
- Transfer learning





































# 1. Forward Pass













classic

Nesterov



 $heta_1$






(a) Standard Neural Net



(b) After applying dropout.