## NPFL114, Lecture 2



# **Training Neural Networks**

Milan Straka

**i** February 20, 2023





EUROPEAN UNION European Structural and Investment Fund Operational Programme Research, Development and Education Charles University in Prague Faculty of Mathematics and Physics Institute of Formal and Applied Linguistics



unless otherwise stated

#### **Refresh – Neural Networks**

- Neural network describes a computation, which gets an input tensor and produces an output.
  - $\circ~$  For the time being, the input tensor has a fixed size.
  - $^{\circ}~$  The input tensor is usually a vector, but it can be 2D/3D/4D tensor.
    - images, video, time sequences like speech, ...
  - $\circ~$  The output usually describes a distribution.
    - normal distribution for regression
    - Bernoulli for binary classification
    - categorical for multiclass classification
- The basic units are **nodes**, composed in an acyclic graph.
- The edges have weights, nodes have activation functions.
- Nodes of neural networks are usually composed in layers.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation SGDs

LR Schedules DoubleD

Ú F<sub>A</sub>L

We usually have a **training set**, which is assumed to consist of examples generated independently from a **data-generating distribution**.

The goal of *optimization* is to match the training set as well as possible.

However, the goal of *machine learning* is to perform well on *previously unseen* data, to achieve lowest **generalization error** or **test error**. We typically estimate it using a **test set** of examples independent of the training set, but generated by the same data-generating distribution.

The **No free lunch theorem** (Wolpert, 1996) states that averaging over *all possible* data distributions, every classification algorithm achieves the same *overall* error when processing unseen examples. In a sense, no machine learning algorithm is *universally* better than others.

Ú F<sub>Á</sub>L

Challenges in machine learning:

- *underfitting* (the model is "too weak", bad performance even on training set)
- *overfitting* (the model is "too strong", learned rules are too specific and do not generalize)



Ú F<sub>A</sub>L

We can control whether a model underfits or overfits by modifying its *capacity*.

- representational capacity (what the model could represent, depends on the model size)
- *effective capacity* (what the model actually learns, depends on training, regularization, ...)



Any change in a machine learning algorithm that is designed to *reduce generalization error* (but not necessarily its training error) is called **regularization**.

 $L^2$  regularization (also called weight decay) penalizes models with large weights (using a penalty of  $\frac{1}{2} \|\boldsymbol{\theta}\|^2$ ).



NPFL114, Lecture 2

**Hyperparameters** are not adapted by a learning algorithm itself, while the model **parameters** (weights, biases) are adapted by it.

Usually a **development set**, also called a **validation set**, is used to estimate the generalization error, allowing to update hyperparameters accordingly.

Gradient Descent

LR Schedules

#### **Loss Function**



A model is usually trained in order to minimize the **loss** on the training data.

Assuming that a model computes  $f(\boldsymbol{x}; \boldsymbol{\theta})$  using parameters  $\boldsymbol{\theta}$ , the mean square error of given N examples  $(\boldsymbol{x}^{(1)}, y^{(1)}), (\boldsymbol{x}^{(2)}, y^{(2)}), \dots, (\boldsymbol{x}^{(N)}, y^{(N)})$  is computed as

$$rac{1}{N}\sum_{i=1}^N \Big(f(oldsymbol{x}^{(i)};oldsymbol{ heta})-y^{(i)}\Big)^2.$$

A common principle used to design loss functions is the maximum likelihood principle.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

LR Schedules E

DoubleD

## **Maximum Likelihood Estimation**

Let  $X = {x^{(1)}, x^{(2)}, ..., x^{(N)}}$  be training data drawn independently from the data-generating distribution  $p_{\text{data}}$ .

We denote the **empirical data distribution** as  $\hat{p}_{\mathrm{data}}$ , where

$$\hat{p}_{ ext{data}}(oldsymbol{x}) \stackrel{ ext{\tiny def}}{=} rac{\left|\{i:oldsymbol{x}^{(i)}=oldsymbol{x}\}
ight|}{N}$$

Let  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  be a family of distributions.

- If the weights are fixed,  $p_{ ext{model}}(\mathbf{x}; \theta)$  is a probability distribution.
- If we instead consider the fixed training data  $\mathbb{X}$ , then

$$L(oldsymbol{ heta}) = p_{ ext{model}}(\mathbb{X};oldsymbol{ heta}) = \prod_{i=1}^N p_{ ext{model}}(oldsymbol{x}^{(i)};oldsymbol{ heta})$$

is called the **likelihood**. Note that even if the value of the likelihood is in range [0, 1], it is not a probability, because the likelihood is not a probability distribution.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

SGDs Adaptive LR

LR Schedules DoubleD

#### **Maximum Likelihood Estimation**

Let  $\mathbb{X} = \{ \boldsymbol{x}^{(1)}, \boldsymbol{x}^{(2)}, \dots, \boldsymbol{x}^{(N)} \}$  be training data drawn independently from the data-generating distribution  $p_{\text{data}}$ . We denote the empirical data distribution as  $\hat{p}_{\text{data}}$  and let  $p_{\text{model}}(\mathbf{x}; \boldsymbol{\theta})$  be a family of distributions.

The maximum likelihood estimation of  $oldsymbol{ heta}$  is:

$$egin{aligned} m{ heta}_{ ext{MLE}} &= rg\max_{m{ heta}} p_{ ext{model}}(\mathbb{X};m{ heta}) = rg\max_{m{ heta}} \prod_{i=1}^N p_{ ext{model}}(m{x}^{(i)};m{ heta}) \ &= rg\min_{m{ heta}} \sum_{i=1}^N -\log p_{ ext{model}}(m{x}^{(i)};m{ heta}) \ &= rg\min_{m{ heta}} \mathbb{E}_{m{ extbf{x}}\sim\hat{p}_{ ext{data}}}[-\log p_{ ext{model}}(m{x};m{ heta})] \ &= rg\min_{m{ heta}} H(\hat{p}_{ ext{data}}(m{ extbf{x}}), p_{ ext{model}}(m{ extbf{x}};m{ heta})) \ &= rg\min_{m{ heta}} D_{ ext{KL}}(\hat{p}_{ ext{data}}(m{ extbf{x}})) \| p_{ ext{model}}(m{ extbf{x}};m{ heta})) + H(\hat{p}_{ ext{data}}(m{ extbf{x}})) \end{aligned}$$

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

SGDs Adaptive LR

LR Schedules



#### **Maximum Likelihood Estimation**



MLE can be easily generalized to the conditional case, where our goal is to predict y given  $m{x}$ :

$$egin{aligned} eta_{ ext{MLE}} &= rg\max_{oldsymbol{ heta}} p_{ ext{model}}(\mathbb{Y}|\mathbb{X};oldsymbol{ heta}) = rg\min_{oldsymbol{ heta}} \sum_{i=1}^N -\log p_{ ext{model}}(y^{(i)}|oldsymbol{x}^{(i)};oldsymbol{ heta}) \ &= rg\min_{oldsymbol{ heta}} \mathbb{E}_{(\mathbf{x},\mathrm{y})\sim\hat{p}_{ ext{data}}}[-\log p_{ ext{model}}(y|oldsymbol{x};oldsymbol{ heta})] \ &= rg\min_{oldsymbol{ heta}} H(\hat{p}_{ ext{data}}(\mathrm{y}|\mathbf{x}), p_{ ext{model}}(\mathrm{y}|\mathbf{x};oldsymbol{ heta})) \ &= rg\min_{oldsymbol{ heta}} D_{ ext{KL}}(\hat{p}_{ ext{data}}(\mathrm{y}|\mathbf{x})) \| p_{ ext{model}}(\mathrm{y}|\mathbf{x};oldsymbol{ heta})) + H(\hat{p}_{ ext{data}}(\mathrm{y}|\mathbf{x})) \end{aligned}$$

where the conditional entropy is defined as  $H(\hat{p}_{\text{data}}) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{\text{data}}}[-\log(\hat{p}_{\text{data}}(y|\boldsymbol{x}; \boldsymbol{\theta}))]$  and the conditional cross-entropy as  $H(\hat{p}_{\text{data}}, p_{\text{model}}) = \mathbb{E}_{(\mathbf{x}, y) \sim \hat{p}_{\text{data}}}[-\log(p_{\text{model}}(y|\boldsymbol{x}; \boldsymbol{\theta}))]$ .

The resulting *loss function* is called **negative log-likelihood** (NLL), or **cross-entropy**, or **Kullback-Leibler divergence**.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation S

LR Schedules DoubleD

#### **Estimators and Bias**

Ú F<sub>A</sub>L

An **estimator** is a rule for computing an estimate of a given value, often an expectation of some random value(s). For example, we might estimate *mean* of a random variable by sampling a value according to its probability distribution.

The **bias** of an estimator is the difference of the expected value of the estimator and the true value being estimated. If the bias is zero, we call the estimator **unbiased**, otherwise **biased**.

If we have a sequence of estimates, it might also happen that the bias converges to zero. Consider the well-known sample estimate of variance. Given independent and identically distributed random variables  $x_1, \ldots, x_N$ , we might estimate the mean and the variance as

$$\hat{\mu} = rac{1}{N}\sum_i x_i, ~~~ \hat{\sigma}^2 = rac{1}{N}\sum_i (x_i - \hat{\mu})^2.$$

Such a mean estimate is unbiased, but the estimate of the variance is biased, because  $\mathbb{E}[\hat{\sigma}^2] = (1 - \frac{1}{N})\sigma^2$ ; however, the bias of this estimate converges to zero for increasing N.

Also, an unbiased estimator does not necessarily have a small variance – in some cases, it can have a large variance, so a biased estimator with a smaller variance might be preferred.

NPFL114, Lecture 2

## **Properties of Maximum Likelihood Estimation**



Assume that the true data-generating distribution  $p_{\text{data}}$  lies within the model family  $p_{\text{model}}(\bullet; \boldsymbol{\theta})$ , and assume there exists a unique  $\boldsymbol{\theta}_{p_{\text{data}}}$  such that  $p_{\text{data}} = p_{\text{model}}(\bullet; \boldsymbol{\theta}_{p_{\text{data}}})$ .

• MLE is a *consistent* estimator. If we denote  $\theta_m$  to be the parameters found by MLE for a training set with m examples generated by the data-generating distribution, then  $\theta_m$  converges in probability to  $\theta_{p_{\text{data}}}$ .

Formally, for any arepsilon>0,  $P(\|m{ heta}_m-m{ heta}_{p_{ ext{data}}}\|>arepsilon) o 0$  as  $m o\infty$ .

MLE is in a sense the most statistically efficient. For any consistent estimator, let us consider the average distance of \$\mathcal{ heta}\_m\$ and \$\mathcal{ heta}\_{p\_{data}}\$: \$\mathbb{E}\_{x\_1,...,x\_m \sim p\_{data}}[\|\mathcal{ heta}\_m - \mathcal{ heta}\_{p\_{data}}\|^2]\$. It can be shown (Rao 1945, Cram\u00e9r 1946) that no consistent estimator has lower mean squared error than the maximum likelihood estimator.

Therefore, for reasons of consistency and efficiency, maximum likelihood is often considered the preferred estimator for machine learning.

## Mean Square Error as MLE

Ú F<sub>A</sub>L

During regression, we predict a number, not a real probability distribution. In order to generate a distribution, we might consider a distribution with the mean of the predicted value and a fixed variance  $\sigma^2$  – the most general such a distribution is the normal distribution.



https://upload.wikimedia.org/wikipedia/commons/3/3a/Linear\_regression.svg

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

SGDs Adaptive LR

LR LR Schedules

DoubleD

#### Mean Square Error as MLE



NPFL114, Lecture 2

Gradient Descent

Backpropagation

Adaptive LR

SGDs

LR Schedules DoubleD

#### **Gradient Descent**



Let a model compute  $f(\boldsymbol{x}; \boldsymbol{\theta})$  using parameters  $\boldsymbol{\theta}$ , and for a given loss function L denote

$$E(oldsymbol{ heta}) = \mathbb{E}_{(oldsymbol{x}, \mathrm{y}) \sim \hat{p}_{ ext{data}}} Lig(f(oldsymbol{x}; oldsymbol{ heta}), yig).$$

Backpropagation

Assuming we are minimizing an error function

$$\underset{\boldsymbol{\theta}}{\arg\min} E(\boldsymbol{\theta}),$$

we may use *gradient descent*:

$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}).$$

The constant  $\alpha$  is called a **learning rate** and specifies the "length" of a step we perform in every iteration of the gradient descent.



NPFL114, Lecture 2

ML Basics MLE

#### **Gradient Descent Variants**



The gradient of the error function  $E(\boldsymbol{\theta})$  can be computed as

$$abla_{oldsymbol{ heta}} E(oldsymbol{ heta}) = \mathbb{E}_{(\mathbf{x}, \mathrm{y}) \sim \hat{p}_{ ext{data}}} 
abla_{oldsymbol{ heta}} Lig(f(oldsymbol{x};oldsymbol{ heta}), yig).$$

- (Standard/Batch) Gradient Descent: We use all training data to compute  $\nabla_{\theta} E(\theta)$ .
- Stochastic (or Online) Gradient Descent: We estimate  $\nabla_{\theta} E(\theta)$  using a single random example from the training data. Such an estimate is unbiased, but very noisy.

 $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta}) \approx \nabla_{\boldsymbol{\theta}} L(f(\boldsymbol{x}; \boldsymbol{\theta}), y)$  for a randomly chosen  $(\boldsymbol{x}, y)$  from  $\hat{p}_{\text{data}}$ .

• Minibatch SGD: Trade-off between gradient descent and SGD – the expectation in  $\nabla_{\theta} E(\theta)$  is estimated using m random independent examples from the training data.

$$abla_{m{ heta}} E(m{ heta}) pprox rac{1}{m} \sum_{i=1}^m 
abla_{m{ heta}} Lig(f(m{x}^{(i)};m{ heta}),y^{(i)}ig) \,\,\, ext{for randomly chosen}\,\,\,(m{x}^{(i)},y^{(i)})\,\,\, ext{from}\,\,\,\hat{p}_{ ext{data}}.$$

NPFL114, Lecture 2

LR Schedules DoubleD

#### **Stochastic Gradient Descent Convergence**



Assume that we perform a stochastic gradient descent, using a sequence of learning rates  $\alpha_i$ , and using a noisy estimate  $J(\boldsymbol{\theta})$  of the real gradient  $\nabla_{\boldsymbol{\theta}} E(\boldsymbol{\theta})$ :

$$oldsymbol{ heta}_{i+1} \leftarrow oldsymbol{ heta}_i - lpha_i J(oldsymbol{ heta}_i).$$

It can be proven (under some reasonable conditions; see Robbins and Monro algorithm, 1951) that if the loss function is convex and continuous, then SGD converges to the unique optimum almost surely if the sequence of learning rates  $\alpha_i$  fulfills the following conditions:

$$\sum_i lpha_i = \infty, ~~ \sum_i lpha_i^2 < \infty.$$

Note that the second condition implies that  $\alpha_i \rightarrow 0$ .

For nonconvex loss functions, we can get guarantees of converging to a *local* optimum only. However, note that finding the global minimum of an arbitrary function is *at least NP-hard*.

LR Schedules DoubleD

#### **Stochastic Gradient Descent Convergence**

Convex functions mentioned on the previous slide are such that for  $m{u}, m{v}$  and real  $0 \leq t \leq 1$ ,



 $f(toldsymbol{u}+(1-t)oldsymbol{v})\leq tf(oldsymbol{u})+(1-t)f(oldsymbol{v}).$ 

A twice-differentiable function of a single variable is convex iff its second derivative is always nonnegative. (For functions of multiple variables, the Hessian must be positive semi-definite.) A local minimum of a convex function is always a global minimum. Well-known examples of convex functions are  $x^2$ ,  $e^x$ ,  $-\log x$ , MSE,  $\sigma$ +NLL, softmax+NLL. NPFL114, Lecture 2 ML Basics MLE Gradient Descent Backpropagation SGDs Adaptive LR LR Schedules DoubleD

#### **Loss Function Visualization**

Visualization of loss function of ResNet-56 (0.85 million parameters) with/without skip connections:



Figure 1 of "Visualizing the Loss Landscape of Neural Nets", https://arxiv.org/abs/1712.09913

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

SGDs Adaptive LR

LR Schedules DoubleD

#### **Loss Function Visualization**

Visualization of loss function of ResNet-110 without skip connections and DenseNet-121:





Figure 4 of "Visualizing the Loss Landscape of Neural Nets", https://arxiv.org/abs/1712.09913

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

on SGDs

Adaptive LR

LR Schedules DoubleD

#### Backpropagation



Assume we want to compute partial derivatives of a given loss function L.



#### **Backpropagation Algorithm**

#### Ú F<sub>Á</sub>L

#### **Forward Propagation**

**Input**: Network with nodes  $u^{(1)}, u^{(2)}, \ldots, u^{(n)}$  numbered in topological order. Each node's value is computed as  $u^{(i)} = f^{(i)}(A^{(i)})$  for  $A^{(i)}$  being a set of values of the predecessors  $P(u^{(i)})$  of  $u^{(i)}$ . **Output**: Value of  $u^{(n)}$ .

- $\begin{array}{ll} \bullet \ \, \mathsf{For} \ i=1,\ldots,n : \\ & \circ \ A^{(i)} \leftarrow \left\{ u^{(j)} | j \in P(u^{(i)}) \right\} \\ & \circ \ u^{(i)} \leftarrow f^{(i)}(A^{(i)}) \end{array}$
- Return  $u^{(n)}$

NPFL114, Lecture 2

## **Backpropagation Algorithm**

#### Simple Variant of Backpropagation

**Input**: The network as in the Forward propagation algorithm. **Output**: Partial derivatives  $g^{(i)} = \frac{\partial u^{(n)}}{\partial u^{(i)}}$  of  $u^{(n)}$  with respect to all  $u^{(i)}$ .

- Run forward propagation to compute all  $u^{(i)}$
- $g^{(n)} = 1$
- For  $i = n 1, \dots, 1$ :  $\circ g^{(i)} \leftarrow \sum_{j:i \in P(u^{(j)})} g^{(j)} \frac{\partial u^{(j)}}{\partial u^{(i)}}$
- Return  $\left(g^{(1)},g^{(2)},\ldots,g^{(n)}
  ight)$

In practice, we do not usually represent networks as collections of scalar nodes; instead we represent them as collections of tensor functions – most usually functions  $f : \mathbb{R}^n \to \mathbb{R}^m$ . Then  $\frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}}$  is a Jacobian matrix. However, the backpropagation algorithm is analogous.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

LR Schedules DoubleD



#### **Neural Network Activation Functions**



## **Hidden Layers Derivatives**

• *σ*:

$$rac{\partial \sigma(x)}{\partial x} = \sigma(x) \cdot ig(1 - \sigma(x)ig)$$

• tanh:

$$rac{\partial anh(x)}{\partial x} = 1 - anh(x)^2$$

• ReLU:

$$rac{\partial \operatorname{ReLU}(x)}{\partial x} = egin{cases} 1 & ext{if } x > 0 \ \operatorname{NaN} & ext{if } x = 0 \ 0 & ext{if } x < 0 \ \end{pmatrix} = rac{\operatorname{assuming } rac{\partial \operatorname{ReLU}(x)}{\partial x}(0) = 0}{=} \left[ x > 0 
ight] = \left[ \operatorname{ReLU}(x) > 0 
ight]$$

NPFL114, Lecture 2

LR Schedules

DoubleD

#### **Stochastic Gradient Descent**



#### Stochastic Gradient Descent (SGD) Algorithm

**Input**: NN computing function  $f(x; \theta)$  with initial value of parameters  $\theta$ . **Input**: Learning rate  $\alpha$ .

**Output**: Updated parameters  $\boldsymbol{\theta}$ .

Repeat until stopping criterion is met:
 Sample a minibatch of *m* training examples (*x*<sup>(i)</sup>, *y*<sup>(i)</sup>)

$$\circ ~~ oldsymbol{g} \leftarrow rac{1}{m} \sum_i 
abla_{oldsymbol{ heta}} Lig(f(oldsymbol{x}^{(i)};oldsymbol{ heta}),y^{(i)}ig)$$

$$\circ \boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \alpha \boldsymbol{g}$$

NPFL114, Lecture 2

SGDs

#### **SGD With Momentum**

#### **SGD** With Momentum

**Input**: NN computing function  $f(\boldsymbol{x}; \boldsymbol{\theta})$  with initial value of parameters  $\boldsymbol{\theta}$ .

**Input**: Learning rate  $\alpha$ , momentum  $\beta$ . **Output**: Updated parameters  $\boldsymbol{\theta}$ .

- $oldsymbol{v} \leftarrow oldsymbol{0}$
- Repeat until stopping criterion is met:
  - $\begin{array}{l} \circ \ \, \text{Sample a minibatch of } m \text{ training examples} \\ (\boldsymbol{x}^{(i)}, y^{(i)}) \\ \circ \ \, \boldsymbol{g} \leftarrow \frac{1}{m} \sum_{i} \nabla_{\boldsymbol{\theta}} L\big(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)}\big) \\ \circ \ \, \boldsymbol{v} \leftarrow \beta \boldsymbol{v} \alpha \boldsymbol{g} \end{array}$
  - $\circ \ oldsymbol{ heta} \leftarrow oldsymbol{ heta} + oldsymbol{v}$



Figure 8.5 of "Deep Learning" book, https://www.deeplearningbook.org

A nice writeup about momentum can be found on <u>https://distill.pub/2017/momentum/</u>.

Backpropagation

LR Schedules DoubleD

#### **SGD With Nesterov Momentum**

#### **SGD With Nesterov Momentum**

**Input**: NN computing function  $f(x; \theta)$ with initial value of parameters  $\theta$ . **Input**: Learning rate  $\alpha$ , momentum  $\beta$ . **Output**: Updated parameters  $\theta$ .

- $oldsymbol{v} \leftarrow oldsymbol{0}$
- Repeat until stopping criterion is met:
  - $^{\circ}\,$  Sample a minibatch of m training examples  $(oldsymbol{x}^{(i)},y^{(i)})$

$$egin{aligned} &\circ oldsymbol{ heta} \leftarrow oldsymbol{ heta} + eta oldsymbol{v} \ &\circ oldsymbol{g} \leftarrow rac{1}{m} \sum_i 
abla_{oldsymbol{ heta}} L(f(oldsymbol{x}^{(i)};oldsymbol{ heta}),y^{(i)}) \ &\circ oldsymbol{v} \leftarrow eta oldsymbol{v} - lpha oldsymbol{g} \ &\circ oldsymbol{ heta} \leftarrow oldsymbol{ heta} - lpha oldsymbol{g} \ \end{aligned}$$



https://github.com/cs231n/cs231n.github.io/blob/master/assets/nn3/nesterov.jpeg

Gradient Descent

Backpropagation SGDs

LR Schedules DoubleD

#### AdaGrad (2011)

**Input**: NN computing function  $f(\boldsymbol{x}; \boldsymbol{\theta})$  with initial value of parameters  $\boldsymbol{\theta}$ . **Input**: Learning rate  $\alpha$ , constant  $\varepsilon$  (usually  $10^{-7}$ ). **Output**: Updated parameters  $\boldsymbol{\theta}$ .

•  $r \leftarrow 0$ 

• Repeat until stopping criterion is met: • Sample a minibatch of m training examples  $(\boldsymbol{x}^{(i)}, y^{(i)})$ •  $\boldsymbol{g} \leftarrow \frac{1}{m} \sum_{i} \nabla_{\boldsymbol{\theta}} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$ •  $\boldsymbol{r} \leftarrow \boldsymbol{r} + \boldsymbol{g}^{2}$ •  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\alpha}{\sqrt{r} + \varepsilon} \boldsymbol{g}$ 

• The  $g^2$  and  $\frac{\alpha}{\sqrt{r}+\varepsilon}g$  are computed element-wise, i.e.,  $g^2 = g \odot g$ . It might be better to write  $\frac{\alpha}{\sqrt{r}+\varepsilon} \odot g$ , but it is not done in the papers, so we are keeping the usual notation.

DoubleD



AdaGrad has favourable convergence properties (being faster than regular SGD) for convex loss landscapes. In this settings, gradients converge to zero reasonably fast.

However, for nonconvex losses, gradients can stay quite large for a long time. In that case, the algorithm behaves as if decreasing learning rate by a factor of  $1/\sqrt{t}$ , because if each

 $oldsymbol{g} pprox oldsymbol{g}_0,$ 

then after t steps

$$oldsymbol{r} pprox t \cdot oldsymbol{g}_0^2,$$

and therefore

$$rac{lpha}{\sqrt{m{r}}+arepsilon}pprox rac{lpha/\sqrt{t}}{\sqrt{m{g}_0^2}+arepsilon/\sqrt{t}}.$$

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation SGDs

LR Schedules DoubleD

#### RMSProp (2012)

**Input**: NN computing function  $f(\boldsymbol{x}; \boldsymbol{\theta})$  with initial value of parameters  $\boldsymbol{\theta}$ . **Input**: Learning rate  $\alpha$ , momentum  $\beta$  (usually 0.9), constant  $\varepsilon$  (usually  $10^{-7}$ ). **Output**: Updated parameters  $\boldsymbol{\theta}$ .

•  $m{r} \leftarrow m{0}$ 

• Repeat until stopping criterion is met: • Sample a minibatch of m training examples  $(\boldsymbol{x}^{(i)}, y^{(i)})$ •  $\boldsymbol{g} \leftarrow \frac{1}{m} \sum_{i} \nabla_{\boldsymbol{\theta}} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$ •  $\boldsymbol{r} \leftarrow \beta \boldsymbol{r} + (1 - \beta) \boldsymbol{g}^{2}$ •  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\alpha}{\sqrt{r+\varepsilon}} \boldsymbol{g}$ 

However, after first step,  $m{r}=(1-eta)m{g}^2$ , which for default eta=0.9 is

$$oldsymbol{r}=0.1oldsymbol{g}^2,$$

so  $m{r}$  is a biased estimate of  $\mathbb{E}[m{g}^2]$  (but the bias converges to zero exponentially fast).



#### Adam (2014)

**Input**: NN computing function  $f(\boldsymbol{x}; \boldsymbol{\theta})$  with initial value of parameters  $\boldsymbol{\theta}$ . **Input**: Learning rate  $\alpha$  (default 0.001), constant  $\varepsilon$  (usually  $10^{-7}$ ). **Input**: Momentum  $\beta_1$  (default 0.9), momentum  $\beta_2$  (default 0.999). **Output**: Updated parameters  $\boldsymbol{\theta}$ .

•  $oldsymbol{s} \leftarrow oldsymbol{0}$ ,  $oldsymbol{r} \leftarrow oldsymbol{0}$ ,  $t \leftarrow 0$ 

• Repeat until stopping criterion is met: • Sample a minibatch of m training examples  $(\boldsymbol{x}^{(i)}, y^{(i)})$ •  $\boldsymbol{g} \leftarrow \frac{1}{m} \sum_{i} \nabla_{\boldsymbol{\theta}} L(f(\boldsymbol{x}^{(i)}; \boldsymbol{\theta}), y^{(i)})$ •  $t \leftarrow t + 1$ •  $\boldsymbol{s} \leftarrow \beta_1 \boldsymbol{s} + (1 - \beta_1) \boldsymbol{g}$  (biased first n•  $\boldsymbol{r} \leftarrow \beta_2 \boldsymbol{r} + (1 - \beta_2) \boldsymbol{g}^2$  (biased second •  $\hat{\boldsymbol{s}} \leftarrow \boldsymbol{s}/(1 - \beta_1^t), \ \hat{\boldsymbol{r}} \leftarrow \boldsymbol{r}/(1 - \beta_2^t)$  (unbiased estin •  $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \frac{\alpha}{\sqrt{\hat{\boldsymbol{x}}} + \hat{\boldsymbol{s}}} \hat{\boldsymbol{s}}$ 

(biased first moment estimate)
(biased second moment estimate)
(unbiased estimates of the moments)



#### **Adam Bias Correction**

To allow analysis, we add indices to the update

$$oldsymbol{s}_t \leftarrow eta_1 oldsymbol{s}_{t-1} + (1-eta_1)oldsymbol{g}_t,$$



Because  $\sum_{i=0}^{\infty} \beta_1^i = \frac{1}{1-\beta_1}$ ,  $s_{\infty}$  is computed as a weighted average of infinitely many elements.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

DoubleD



#### **Adam Bias Correction**

However, for  $t < \infty$ , the sum of weights in the computation of  $s_t$  does not sum to one.

To obtain an unbiased estimate, we therefore need to account for the "missing" elements; in other words, we need to scale the weights so that they sum to one.

The sum of weights after t steps is

$$(1-eta_1)\sum_{i=1}^teta_1^{t-i}=\sum_{i=1}^teta_1^{t-i}-\sum_{i=0}^{t-1}eta_1^{t-i}=1-eta_1^t,$$

so we obtain an unbiased estimate by dividing  $s_t$  with  $(1 - \beta_1^t)$ , and analogously for the correction of r.





 $(1-\beta)$ 

ML Basics MLE

Backpropagation SGDs

LR Schedules DoubleD





NPFL114, Lecture 2 ML Basics MLE Gradient Descent Backpropagation SGDs Adaptive LR LR Schedules DoubleD 35/45





http://2.bp.blogspot.com/-L98w-SBmF58/VPmICIjKEKI/AAAAAAAACCs/rrFz3VetYmM/s400/ found at http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html



http://3.bp.blogspot.com/-nrtJPrdBWuE/VPmIB46F2al/AAAAAAAACCw/vaE\_B0SVy5k/s400/ found at http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html

NPFL114, Lecture 2 ML Basics MLE Gradient Descent Backpropagation SGDs Adaptive LR LR Schedules

37/45

DoubleD





http://1.bp.blogspot.com/-K\_X-yud8nj8/VPmIBxwGlsI/AAAAAAACC0/JS-h1fa09EQ/s400/ found at http://www.denizyuret.com/2015/03/alec-radfords-animations-for.html

NPFL114, Lecture 2 ML Basics

Gradient Descent

MLE

Backpropagation

SGDs Adaptive LR

LR Schedules DoubleD

#### **Learning Rate Schedules**

Even if RMSProp and Adam are adaptive, they still usually require carefully tuned decreasing learning rate for top-notch performance.

- **Polynomial decay**: learning rate is multiplied by some polynomial of the current update number *t*.
  - Linear decay uses  $\alpha_t = \alpha_{\text{initial}} \cdot \left(1 \frac{t}{\max \text{ steps}}\right)$  and has theoretical guarantees of convergence, but is usually too fast for deep neural networks.
  - $\circ~$  Inverse square root decay uses  $lpha_t = lpha_{ ext{initial}} \cdot rac{1}{\sqrt{t}}$  and

is currently used by best machine translation models.

• **Exponential decay**: learning rate is multiplied by a constant each minibatch/epoch/several epochs.

$$\circ ~~ lpha_t = lpha_{ ext{initial}} \cdot c^t$$

• Often used for convolutional networks (image recognition etc.).





NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

Adaptive LR

SGDs

LR Schedules DoubleD

#### **Learning Rate Schedules**

• **Cosine decay**: The cosine decay has became quite popular in the past years, both for training and finetuning.

$$lpha_t = lpha_{ ext{initial}} \cdot rac{1}{2} igg( 1 + \cos igg( \pi \cdot rac{t}{ ext{max steps}} igg) igg)$$



• Cyclic restarts, warmup, ...

The tf.optimizers.schedules offers several such learning rate schedules, which can be passed to any Keras optimizer directly as a learning rate.

- tf.optimizers.schedules.PiecewiseConstantDecay
- tf.optimizers.schedules.PolynomialDecay
- tf.optimizers.schedules.ExponentialDecay
- tf.optimizers.schedules.CosineDecay

NPFL114, Lecture 2

SGDs

DoubleD

LR Schedules

#### Why do Neural Networks Generalize so Well – Double Descent <sup>Ú</sup>



Figure 1: Curves for training risk (dashed line) and test risk (solid line). (a) The classical *U-shaped risk curve* arising from the bias-variance trade-off. (b) The *double descent risk curve*, which incorporates the U-shaped risk curve (i.e., the "classical" regime) together with the observed behavior from using high capacity function classes (i.e., the "modern" interpolating regime), separated by the interpolation threshold. The predictors to the right of the interpolation threshold have zero training risk.

Figure 1 of "Reconciling modern machine learning practice and the bias-variance trade-off", https://arxiv.org/abs/1812.11118

#### **Deep Double Descent**





Figure 1: Left: Train and test error as a function of model size, for ResNet18s of varying width on CIFAR-10 with 15% label noise. **Right:** Test error, shown for varying train epochs. All models trained using Adam for 4K epochs. The largest model (width 64) corresponds to standard ResNet18.

NPFL114, Lecture 2

ML Basics MLE

Gradient Descent

Backpropagation

SGDs

Adaptive LR

LR Schedules DoubleD

#### **Deep Double Descent – Effective Model Complexity**



The authors define the **Effective Model Complexity** (EMC) of a training procedure  $\mathcal{T}$  with respect to distribution  $\mathcal{D}$  and parameter  $\varepsilon > 0$  as

$$\mathrm{EMC}_{\mathcal{D},arepsilon}(\mathcal{T}) \stackrel{\scriptscriptstyle\mathrm{def}}{=} \maxig\{n \, ig|\, \mathbb{E}_{\mathrm{S}\sim\mathcal{D}^n}[\mathrm{Error}_S(\mathcal{T}(S))] \leq arepsilonig\},$$

where  $\operatorname{Error}_{S}(M)$  is the mean error of a model M on the train samples S.

**Hypothesis:** For any natural data distribution  $\mathcal{D}$ , neural-network-based training procedure  $\mathcal{T}$ , and small  $\varepsilon > 0$ , if we consider the task of predicting labels based on n samples from  $\mathcal{D}$ , then:

- Under-parametrized regime. If  $EMC_{\mathcal{D},\varepsilon}(\mathcal{T})$  is sufficiently smaller than n, any perturbation of  $\mathcal{T}$  that increases its effective complexity will decrease the test error.
- Over-parametrized regime. If  $\text{EMC}_{\mathcal{D},\varepsilon}(\mathcal{T})$  is sufficiently larger than n, any perturbation of  $\mathcal{T}$  that increases its effective complexity will decrease the test error.
- Critically parametrized regime. If  $\mathrm{EMC}_{\mathcal{D},\varepsilon}(\mathcal{T}) \approx n$ , then a perturbation of  $\mathcal{T}$  that increases its effective complexity might decrease or increase the test error.

LR Schedules

#### Why do Neural Networks Generalize so Well



Figure 2: Left: Test error as a function of model size and train epochs. The horizontal line corresponds to model-wise double descent–varying model size while training for as long as possible. The vertical line corresponds to epoch-wise double descent, with test error undergoing double-descent as train time increases. **Right** Train error of the corresponding models. All models are Resnet18s trained on CIFAR-10 with 15% label noise, data-augmentation, and Adam for up to 4K epochs. *Figure 2 of "Deep Double Descent: Where Bigger Models and More Data Hurt", https://arxiv.org/abs/1912.0229* 

**NPFL114, Lecture 2** ML Basics MLE Gradient Descent Backpropagation SGDs Adaptive LR LR Schedules **DoubleD** 44/45

#### Why do Neural Networks Generalize so Well



(a) **CIFAR-100.** There is a peak in test error even with no label noise.

(b) **CIFAR-10.** There is a "plateau" in test error around the interpolation point with no label noise, which develops into a peak for added label noise.

Figure 4 of "Deep Double Descent: Where Bigger Models and More Data Hurt", https://arxiv.org/abs/1912.02292