

-sometitle-

Classifying N-body simulations with and without relativistic corrections
using machine learning techniques

Johan Mylius Kroken

Computational Science: Astrophysics
60 ECTS study points

Institute of Theoretical Astrophysics
Faculty of Mathematics and Natural Sciences

Johan Mylius Kroken

-sometitle-

Classifying N-body simulations with and without
relativistic corrections using machine learning
techniques

Supervisors:

A David Fonseca Mota

B Julian Adamek

C Francisco Antonio Villaescusa Navarro

Abstract

On large scales, comparable to the horizon, relativistic effects will affect the cosmological observables. In order to solve for these effects, one need to consistently solve for the metric, velocities and densities in a particular gauge. When simulating large-scale structures we use N-body simulations, which are usually performed in the Newtonian limit. However, it is not obvious that Newtonian gravity yield a good global description of an inhomogeneous cosmology across all scales (Jeong, Schmidt and Hirata 2012). However, literature suggest that Newtonian simulations are still solving the dynamics correctly, even on large scales close to the horizon where relativistic effects are important but may be corrected for (Chisari and Zaldarriaga 2011) (Green and Wald 2012).

Recently, Adamek et al. 2016 developed a relativistic N-body code, **gevolution**, which evolves large scales structures based on the weak field expansion in GR. I investigate the differences in the gravitational dynamics between structures evolved with and without relativistic effects, with focus on the gravitational potential Φ . This is a good choice for comparison as Φ is gauge invariant and the Newtonian and relativistic simulations are performed in different gauges.

The investigation is done by running 2000 simulations using identical Λ CDM cosmologies for the two gravity theories. The simulations are run using 64^4 particles on a 256^3 grid each with dimension 5120 Mpc/h which is a compromise in order to include both large and nonlinear scales. The data analysis consist of a preliminary analysis using conventional summary statistics, with focus on the bispectrum of Φ . There is a difference in the two cases for low redshifts in the equilateral and squeezed configurations. However, the main idea is to train a Convolutional Neural Network (CNN) to classify the two cases, given snapshots of Φ . The main analysis then involves interpretability of the CNN, which may be done by considering for instance saliency maps (Alqaraawi et al. 2020) or Grad-CAM (Selvaraju et al. 2020). In either case, revealing the features separating the two cases may help us understand the differences in the gravitational dynamics between the two theories. I expect that such a network is able to find relativistic corrections to the Newtonian snapshots that are of higher order than those obtained from power spectra and bispectra analysis. Further, it may also reveal which configurations of Fourier modes \mathbf{k} yield the highest bispectral power, which for now is mainly trial and error.

Contents

1	Introduction	1
1.1	Motivation	1
1.2	Outline	1
1.3	Aim	1
1.4	Nomenclature	1
I	Cosmological Structure Formation	3
2	Preliminaries	5
2.1	General Relativity	5
2.1.1	Einstein's Field Equations	5
2.1.2	Riemann Connection and Covariant Derivatives	5
2.1.3	Geodesic Equation	5
2.1.4	The Stress-Energy Tensor	5
2.2	Useful Relations	5
3	Background Cosmology	7
3.1	The homogeneous Universe	7
3.1.1	The Cosmological Principle	7
3.1.2	The Robertson-Walker Metric	7
3.1.3	The Friedmann Equations	7
3.2	My Universe is loaded with...	7
3.3	Thermal History of the Universe	7
4	Perturbation Theory	9
4.1	Initial Conditions	9
4.2	Transfer Functions	9
4.3	Power Spectra	9
4.4	Linear Evolution	9
4.5	Non-linear Evolution	9
4.6	Bispectra	9
4.6.1	Analytical Bispectrum	9
5	Simulation theory	13
5.1	N-body simulations	13
5.1.1	Describing a box of particles	13
5.1.2	Forces and Fields	13
5.1.3	Mass Assignment Schemes	13
5.1.4	Validity of Box	13
5.2	Newtonian Approach	13
5.3	General Relativistic Approach	13

II	Machine Learning	15
6	Fundamental Elements of Machine Learning	17
6.1	Introduction	17
6.2	Linear Algebra	17
6.3	Probability and Information Theory	17
6.4	Basic Machine Learning	17
6.4.1	Optimisation and Generalisation	17
6.4.2	Data and Fitting	18
6.4.3	Estimators, Bias, Variance and Error	18
6.4.4	Maximum Likelihood Estimation	19
6.4.5	Bayesian Statistics	19
6.4.6	Supervised Learning	19
6.4.7	Unsupervised Learning	19
7	Neural Networks.	21
7.1	Forward pass - Prediction	21
7.1.1	Computational Graphs	21
7.1.2	Architecture	21
7.1.3	Activation	22
7.1.4	Full graph and its output	22
7.1.5	Activation functions.	23
7.1.6	Loss functions.	24
7.1.7	Regularization	24
7.2	Backpropagation - Training	24
7.2.1	Chain rule	24
7.2.2	Backpropagation algorithm	25
7.2.3	Optimization of parameters.	25
7.3	Challenges in optimization	27
7.4	Regularization	28
7.4.1	Parameter norm penalties	28
7.4.2	Data augmentation	28
7.4.3	Early stopping.	28
7.4.4	Parameter sharing	28
7.4.5	Ensemble methods.	28
7.4.6	Dropout	28
8	Convolutional Neural Networks	29
8.1	Convolution	29
8.1.1	Basic definitions	29
8.1.2	Properties	30
8.2	New Layers	30
8.2.1	Convolutional layers	30
8.2.2	Pooling layers	32

III	Acquiring Data	31
9	Simulations.	33
9.1	Parameters	33
9.1.1	Cosmological parameters	33
9.1.2	Primordial power spectrum	33
9.1.3	Box parameters	34
9.1.4	Seeds and gravity theories	34
9.2	Output.	34
9.2.1	Datacubes	34
9.2.2	Redshift problem.	35
9.2.3	Statistics of dataset.	35
9.2.4	Normalisation	35
10	Data Verification.	37
10.1	Slices of Datacubes.	37
10.2	Powerspectra from Simulations	37
10.3	Powerspectra from Datacubes	37
10.4	Analytical Bispectra	37
10.5	Bispectra from Cube	37
10.5.1	Binning	37
10.5.2	Bispectra	37
11	Trainable Dataset	45
11.1	Gathering the cubes	45
11.2	Find image from index	45

Contents

List of Figures

4.1	Angles in arbitrary bispectrum triangle configuration where $\sum_i \mathbf{k}_i = 0$	10
7.1	Computational graph of a function $\mathbf{z} = f(g(\mathbf{x}))$	21
7.2	Computational graph of a hidden layer of a neural network.	22
7.3	Computational graph of a fully connected feed forward neural network. . . .	23
8.1	TODO: fix this figure to be correct w.r.t. cross-correlation/convolution Example of a cross-correlation operation, which we will also call a convolution. The input image is given by the matrix \mathbf{I} and the kernel is given by the matrix \mathbf{K} . The output is given by the matrix $\mathbf{K} \star \mathbf{I}$. The dashed lines indicate which elements of the input and kernel are multiplied together to produce the output. The figure is taken from TikZ.net TODO: fix this citation	30
10.1	Slice 0	38
10.2	Average matter power spectra at different redshifts.	39
10.3	Average potential power spectra at different redshifts.	40
11.1	Outline of data structure.	46

List of Figures

List of Tables

9.1	Cosmological parameters	33
9.2	Primordial power spectra parameters	33
9.3	Box parameters	34

List of Tables

Preface

Here comes your preface, including acknowledgments and thanks.

Part II

Machine Learning

Chapter 6

Fundamental Elements of Machine Learning

6.1 Introduction

In this chapter I will give a brief introduction into machine learning. This includes a mathematical description of some fundamental concepts common across numerous machine learning models. The more advanced models will be dealt with at a later stage. If not otherwise stated, the following chapter is based on Goodfellow, Bengio and Courville 2016 and Hastie, Tibshirani and Friedman 2009.

6.2 Linear Algebra

maybe

6.3 Probability and Information Theory

maybe

6.4 Basic Machine Learning

TODO: [Fill more here](#)

6.4.1 Optimisation and Generalisation

Optimisation Optimisation problems are problems in which we want to minimise some error, given some data. In other words, we want to optimise an algorithm or model given a specific dataset. We care about the error of our model for that specific dataset only.

Generalisation The concept of *generalisation* is what makes a machine learning model different from an optimisation model. We still *train* the machine learning model on some specific data, but we measure how good the model is based on how it performs on a different set of data, which it was not trained on. I.e. we need a generalised model, which is not restricted to the data it was trained on, and has the ability to perform well on unobserved data.

6.4.2 Data and Fitting

Data The key ingredient to any machine learning algorithm is the data fed into it. In an optimisation problem we deal with one single set of data on which we train the model. In a machine learning scenario we want to quantify how well the model is generalised. This is done by training the model on part of the data only, called the *training data*. The model is then assessed on the remaining data, on which it was not trained, in order to determine how general it is. I will come back to how the models are trained, but in essence they are learning features of the data they are trained on. Thus, in order to have a general model we need the data to have some inherent properties. The most important property is that the data is *independent and identically distribution* (i.i.d.). In other words, both the training and testing data are drawn from the same probability distribution.

Training and testing When training the model, we want to minimise some error with respect to the training data, which is essentially an optimisation problem. The model is then tested by measuring the same error with respect to the testing data. It is the latter error we want to be low in order to call the model general. The discrepancy between the two are summaries by the concept of fitting.

Overfitting Overfitting is when we train the model in such a way that the training error becomes too low. This optimises the model too much with respect to the training data, effectively reconstructing the training data point by point and thereby loosing trends in the data features. The result of this is poor performance on the testing data, and as a result a non-general model.

Underfitting Underfitting is the opposite, when the model is unable to achieve a low error on the training data, not capturing the main features of the dataset. This is also bad, as the training error would also be large. This also results in a non-general model.

[TODO: Figure explaining over- and underfitting](#)

6.4.3 Estimators, Bias, Variance and Error

Estimators Based on the assumption that there exists some true parameter(s) θ which remain unknown,¹ we are able to make predictions and estimations of such parameter(s). Let's say we have m independent and identically distributed (i.i.d.) random variables $\{\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m\}$ drawn from the same probability distribution $p(\mathbf{x})$. An *estimator* of the true values θ is any function of the data such that $\hat{\theta}_m = g(\mathbf{x}_1, \dots, \mathbf{x}_m)$, where $\hat{\theta}$ is the estimate of θ . This is known as point estimation, as we are estimating a single value. This definition does not pose any restrictions on the function g . However, a good estimator would yield an estimate $\hat{\theta}_m$ that is close to the true value θ .

Function estimators Say we want to predict a variable \mathbf{y} given some vector \mathbf{x} . We assume the true variable \mathbf{y} is given by some function approximation $f(\mathbf{x})$ plus some error ϵ : $\mathbf{y} = f(\mathbf{x}) + \epsilon$. The aim is then to estimate the function f with the estimator \hat{f} . If we then realise that \hat{f} is really just a point estimator in function space, the two above concepts are equivalent.

¹This is the frequentist perspective of statistics

Bias The bias of the estimator $\hat{\theta}_m$ is defined as the difference between the expected value of the estimator and the true value of the parameter: $\text{Bias}[\hat{\theta}_m] = \mathbb{E}[\hat{\theta}_m] - \theta$. An unbiased estimator has zero bias, i.e. $\mathbb{E}[\hat{\theta}_m] = \theta$. An estimator is asymptotically unbiased if its bias approaches zero as the number of data points m approaches infinity, i.e. $\lim_{m \rightarrow \infty} \mathbb{E}[\hat{\theta}_m] = \theta$.

Variance Variance serves as a crucial metric in assessing the variability of an estimator concerning changes in the input data. Denoted as $\text{Var}[\hat{\theta}]$, where $\hat{\theta}$ represents the training set, low variance is desirable in an estimator. In essence, variance quantifies the extent to which we anticipate the estimator to fluctuate based on different datasets.

Standard Error The standard error, denoted as $\text{SE}[\hat{\theta}]$, is a fundamental concept closely related to variance, being the square root of the latter. Specifically, for the estimator of the mean, $\hat{\mu}$, the standard error is defined by the formula:

$$\text{SE}[\hat{\mu}] = \sqrt{\text{Var}\left[\frac{1}{m} \sum_{i=1}^m x^{(i)}\right]} = \frac{\sigma}{\sqrt{m}}, \quad (6.1)$$

Here, σ^2 represents the true variance. The standard error of $\hat{\mu}$ is pivotal in gauging the precision of the estimator for the mean, especially when applied to testing data, aiding in the evaluation of a model's generalization performance.

Mean Squared Error The Mean Squared Error (MSE) provides a comprehensive measure for evaluating the performance of an estimator. Expressed as:

$$\text{MSE} = \mathbb{E}[(\hat{\theta} - \theta)^2] = \text{Bias}[\hat{\theta}]^2 + \text{Var}[\hat{\theta}], \quad (6.2)$$

MSE encompasses both the bias and variance components. The MSE offers a unified perspective, capturing the overall accuracy and variability of the estimator in estimating the true parameter θ .

6.4.4 Maximum Likelihood Estimation

6.4.5 Bayesian Statistics

6.4.6 Supervised Learning

6.4.7 Unsupervised Learning

Chapter 7

Neural Networks

7.1 Forward pass - Prediction

7.1.1 Computational Graphs

In order to understand how information flows through a neural network we must understand computational graphs. These are collection of *nodes* and *edges*, where each node represents an *operation* and each edge represents the numerical values as they flow through the network. Nodes may be simple functions such as addition or multiplication but also take more complex. In the latter case, they can be further divided into their respective component, so that the most fundamental computational graph only consists fundamental operations.

For example, we consider the functions $\mathbf{y} = g(\mathbf{x})$ and $\mathbf{z} = f(\mathbf{y})$ such that $\mathbf{z} = f(g(\mathbf{x}))$. This can be represented as a computational graph as shown in Fig. 7.1. The computational graph is read from left to right, where the input \mathbf{x} is fed into the function g which outputs \mathbf{y} . This \mathbf{y} is then fed into the function f which outputs \mathbf{z} . The computational graph is a visual representation of the function $f(g(\mathbf{x}))$.

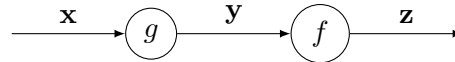


Figure 7.1: Computational graph of a function $\mathbf{z} = f(g(\mathbf{x}))$.

7.1.2 Architecture

The simplest architecture of a neural network is the fully connected feed forward (FCNN), which consist of L layers in total, the first being an input layer and the remaining $L - 1$ layers are called *hidden layers*, \mathbf{h} . Each hidden layer has an *activation* \mathbf{a} , used as an input to an *activation function*, g . Mathematically we may write this architecture as:

$$\begin{aligned} \mathbf{h}^0 &= \mathbf{x}^{(i)} \\ \mathbf{h}^1 &= g_1(\mathbf{a}^1) \\ \mathbf{h}^2 &= g_2(\mathbf{a}^2) \\ &\vdots \\ \mathbf{h}^L &= g_L(\mathbf{a}^L), \end{aligned} \tag{7.1}$$

where \mathbf{h}^L and g_L are the output layer and output function respectively. The parameter L governs the depth of the neural network. The result of the output layer is simply called the *output* or *predictor*, and typically denoted as $\hat{\mathbf{y}} = \mathbf{h}^L$. When training, we will have a true value, or label, denoted \mathbf{y} . The difference between the true value and the predicted value is called the *residual* or *error*, denoted $\tilde{\mathbf{y}} = \mathbf{y} - \hat{\mathbf{y}}$. The goal of training is to minimise the error, or equivalently, maximise the likelihood of the true value.

7.1.3 Activation

The activation \mathbf{a}^l of a layer l is an affine transformation of the output of the previous layer, \mathbf{h}^{l-1} . The intercept of this affine transformation is known as the bias \mathbf{b}^l ,¹ typically used to ensure that no activation becomes zero. The activation takes the form:

$$\mathbf{a}^l = (\mathbf{W}^{l-1 \rightarrow l})^T \mathbf{h}^{l-1} + \mathbf{b}^l, \quad (7.2)$$

where $\mathbf{W}^{l-1 \rightarrow l} \in \mathbb{R}^{\dim_{l-1} \times \dim_l}$ is the matrix of weights describing the mapping from layer $l-1$ to layer l . Each layer l has dimension (or neurons) \dim_l which governs the width of each layer.

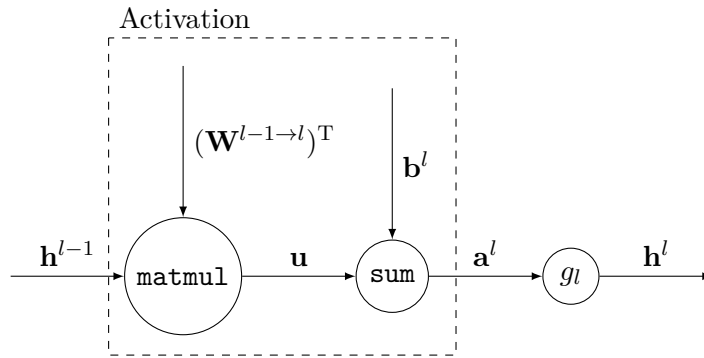


Figure 7.2: Computational graph of a hidden layer of a neural network.

7.1.4 Full graph and its output

The full computational graph of a fully connected neural network is shown in Fig. 7.3. The output of the network is the output of the output layer, \mathbf{h}^L . The output function is typically a probability distribution, such as the softmax function, or a scalar value, such as the sigmoid function. Either way, the output is the prediction of the network, also noted $\hat{\mathbf{y}} = \mathbf{h}^L$. In order to determine the goodness of this output we must compare it with the target value \mathbf{y} . This is done using a *loss function*, $\mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$, which is a function of the prediction and the target. The loss function is a measure of how good the prediction is. The goal of training is to minimise the loss function, or equivalently, maximise the likelihood of the target value. We define the loss $\delta = \mathcal{L}(\mathbf{y}, \hat{\mathbf{y}})$ which is the error of the prediction. The loss is then used to update the weights and biases of the network, which is done using the *backpropagation* algorithm explained in section [TODO: ref this](#)

¹Must not be confused with the statistical bias of an estimator.

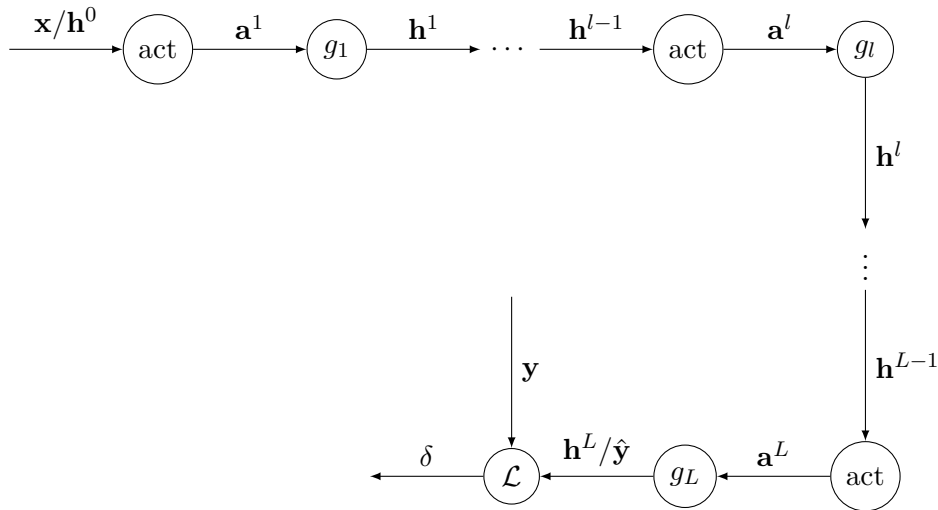


Figure 7.3: Computational graph of a fully connected feed forward neural network.

7.1.5 Activation functions

The activation functions, g_l , are what determines the output of each layer l . They can be broadly classified into two types: *saturating activation functions* and *non-saturating activation function*. For saturating activations, $\lim_{|\mathbf{x}| \rightarrow \infty} |\nabla g(\mathbf{x})| = 0$.

Linear The linear activation is as simple as it gets:

$$g(\mathbf{x}) = \mathbf{a}\mathbf{x} + b \quad \text{and} \quad g'(\mathbf{x}) = a. \quad (7.3)$$

This activation is non-saturating. There is however a problem with using linear activations, [TODO: Write about UNIVERSAL APPROXIMATION THEOREM, here or elsewhere](#)

Sigmoid The sigmoid function, or the logistic function is defined as:

$$g(\mathbf{x}) = \sigma(\mathbf{x}) = \frac{1}{1 + e^{-\mathbf{x}}} \quad \text{and} \quad g'(\mathbf{x}) = g(\mathbf{x})(1 - g(\mathbf{x})) \quad (7.4)$$

Hyperbolic tangent The hyperbolic tangent is defined as:

$$g(\mathbf{x}) = \frac{e^{\mathbf{x}} - e^{-\mathbf{x}}}{e^{\mathbf{x}} + e^{-\mathbf{x}}} \quad \text{and} \quad g'(\mathbf{x}) = 1 - g(\mathbf{x})^2 \quad (7.5)$$

Rectified linear unit The rectified linear unit, ReLU, is defined as:

$$g(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} < 0 \\ \mathbf{x}, & \mathbf{x} \geq 0 \end{cases} \quad \text{and} \quad g'(\mathbf{x}) = \begin{cases} 0, & \mathbf{x} < 0 \\ 1, & \mathbf{x} \geq 0. \end{cases} \quad (7.6)$$

(Scaled) Exponential linear unit The exponential linear unit, (S)ELU, is defined as:

$$g(\mathbf{x}) = \lambda \begin{cases} \alpha(e^{\mathbf{x}} - 1), & \mathbf{x} < 0 \\ \mathbf{x}, & \mathbf{x} \geq 0 \end{cases} \quad \text{and} \quad g'(\mathbf{x}) = \lambda \begin{cases} \alpha e^{\mathbf{x}}, & \mathbf{x} < 0 \\ 1, & \mathbf{x} \geq 0. \end{cases} \quad (7.7)$$

This function requires two hyperparameters α and λ .

Parametric rectified linear unit The parametric rectified linear unit, PReLU, is defined as:

$$g(\mathbf{x}) = \begin{cases} \alpha \mathbf{x}, & \mathbf{x} < 0 \\ \mathbf{x}, & \mathbf{x} \geq 0 \end{cases} \quad \text{and} \quad g'(\mathbf{x}) = \begin{cases} \alpha, & \mathbf{x} < 0 \\ 1, & \mathbf{x} \geq 0. \end{cases} \quad (7.8)$$

If $\alpha = 0.01$ this function is also known as the leaky rectified unit, or Leaky ReLU.

Softmax The softmax activation function is defined as:

$$g(\mathbf{x}) = \frac{e^{\mathbf{x}}}{\sum_{i=1}^n e^{x_i}} \quad \text{and} \quad g'(\mathbf{x}) = g(\mathbf{x})(1 - g(\mathbf{x})). \quad (7.9)$$

This activation function is typically used for the output layer of a neural network, where the output is a probability distribution over n classes, i.e. multiclass classification problems.

7.1.6 Loss functions

[TODO: write about loss functions](#)

MSE

$$\mathcal{L}_{\text{MSE}}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{1}{n} \tilde{\mathbf{y}}^T \tilde{\mathbf{y}}, \quad (7.10)$$

where $\tilde{\mathbf{y}} \equiv \mathbf{y} - \hat{\mathbf{y}}$

Cross entropy

$$\mathcal{L}_{\text{CE}}(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^n y_i \log(\hat{y}_i) \quad (7.11)$$

where \hat{y}_i is the predicted probability of class i and y_i is the true probability of class i , i.e. the target probability for that class. n is the number of classes. In the binary case, $n = 2$, this is equivalent to the binary cross entropy:

$$\mathcal{L}_{\text{BCE}}(\mathbf{y}, \hat{\mathbf{y}}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}). \quad (7.12)$$

7.1.7 Regularization

7.2 Backpropagation - Training

7.2.1 Chain rule

Basics Lets say we have a function $f(x)$ that is composed of two functions $g(x)$ and $h(x)$, such that $f(x) = g(h(x))$. The chain rule states that the derivative of $f(x)$ with respect to x is the product of the derivative of $g(x)$ with respect to $h(x)$ and the derivative of $h(x)$ with respect to x :

$$\frac{df}{dx} = \frac{dg}{dh} \frac{dh}{dx}. \quad (7.13)$$

This can be generalised into the vector case where $f(\mathbf{x}) = g(h(\mathbf{x}))$ and $\mathbf{x} \in \mathbb{R}^n$ and [TODO: write more ahaha](#)

7.2.2 Backpropagation algorithm

The learnable parameters of a neural network are the weights and biases of each layer. In order to update these in an educated² way, we need to know how the loss, δ , changes with respect to each parameter, $\nabla_{\mathbf{b}^l} \delta, \nabla_{\mathbf{W}^{l-1 \rightarrow l}} \delta \quad \forall l \in [1, L]$. This is done using the chain rule, starting from the output layers and working our way backwards through the network. Recalling that the activation is written as Eq. (7.2) it makes sense to find the derivative of the loss with respect to the output activation:

$$\frac{d\delta}{d\mathbf{a}_i^L} = \sum_j \frac{\partial \delta}{\partial h_j^L} \frac{dh_j^L}{da_i^L} = \sum_j \frac{\partial \delta}{\partial \hat{y}_j} \frac{dg_L}{da_i^L} \Big|_{\mathbf{a}_j^L} \leftrightarrow \nabla_{\mathbf{a}^L} \delta = \nabla_{\hat{\mathbf{y}}} \delta \odot \frac{dg_L}{d\mathbf{a}^L} \Big|_{\mathbf{a}^L}, \quad (7.14)$$

where \odot is the *element-wise Hadamard product*. The quantity $\nabla_{\mathbf{a}^L} \delta$ is the gradient of the loss with respect to the output activation. This is propagated backwards through the network as follows:

$$\nabla_{\mathbf{h}^{l-1}} \delta = \mathbf{W}^{l-1 \rightarrow l} \nabla_{\mathbf{a}^l} \delta \quad \text{and} \quad \nabla_{\mathbf{a}^{l-1}} \delta = \nabla_{\mathbf{h}^{l-1}} \delta \odot \frac{dg_{l-1}}{d\mathbf{a}^{l-1}} \Big|_{\mathbf{a}^{l-1}}. \quad (7.15)$$

When we know these loss gradients we can compute the gradient of the weight and biases:

$$\nabla_{\mathbf{W}^{l-1 \rightarrow l}} \delta = \mathbf{h}^{l-1} (\nabla_{\mathbf{a}^l} \delta)^T \quad \text{and} \quad \nabla_{\mathbf{b}^l} \delta = \nabla_{\mathbf{a}^l} \delta. \quad (7.16)$$

We further define a parameter tuple $\Theta^l \equiv (\mathbf{W}^{l-1 \rightarrow l}, \mathbf{b}^l)$ and a gradient tuple $\nabla_{\Theta^l} \delta \equiv (\nabla_{\mathbf{W}^{l-1 \rightarrow l}} \delta, \nabla_{\mathbf{b}^l} \delta)$. The weight and biases of each layer is then updated as $\Theta^l \leftarrow \mathcal{U}(\Theta^l, \nabla_{\Theta^l} \delta)$, where \mathcal{U} is an update function. This is done for each layer l in the network. The backpropagation algorithm is summarised in Algorithm 1. [TODO: Revise this](#)

Algorithm 1 Backpropagation algorithm

Require: Loss δ , output activation \mathbf{a}^L , output layer \mathbf{h}^L , weight and bias tuples $\Theta^l \equiv (\mathbf{W}^{l-1 \rightarrow l}, \mathbf{b}^l)$, activation functions g_l and update function \mathcal{U} .

Ensure: Updated weight and bias tuples $\Theta^l \equiv (\mathbf{W}^{l-1 \rightarrow l}, \mathbf{b}^l)$.

```

1:  $\nabla_{\mathbf{a}^L} \delta \leftarrow \nabla_{\hat{\mathbf{y}}} \delta \odot \frac{dg_L}{d\mathbf{a}^L} \Big|_{\mathbf{a}^L}$ 
2: for  $l = L, L-1, \dots, 1$  do
3:    $\nabla_{\mathbf{h}^{l-1}} \delta \leftarrow \mathbf{W}^{l-1 \rightarrow l} \nabla_{\mathbf{a}^l} \delta$ 
4:    $\nabla_{\mathbf{a}^{l-1}} \delta \leftarrow \nabla_{\mathbf{h}^{l-1}} \delta \odot \frac{dg_{l-1}}{d\mathbf{a}^{l-1}} \Big|_{\mathbf{a}^{l-1}}$ 
5:    $\nabla_{\mathbf{W}^{l-1 \rightarrow l}} \delta \leftarrow \mathbf{h}^{l-1} (\nabla_{\mathbf{a}^l} \delta)^T$ 
6:    $\nabla_{\mathbf{b}^l} \delta \leftarrow \nabla_{\mathbf{a}^l} \delta$ 
7:    $\Theta^l \leftarrow \mathcal{U}(\Theta^l, \nabla_{\Theta^l} \delta)$ 
8: end for
```

[algorithm](#)

7.2.3 Optimization of parameters

Gradient descent The ultimate goal is to arrive at some parameters Θ that minimize the output loss δ . In order to do so, we must update the parameters until this goal is reached. The question is then *how* we should update these parameters. In the most

²Instead of just wildly guessing

general case, we need to find some $\Delta\Theta$ such that $\Theta \leftarrow \Theta + \Delta\Theta$, where the only condition is that $\Delta\Theta$ yield a new Θ whose loss decrease. The backbone of this thinking is the *gradient descent* algorithm, which is a first order optimization algorithm. The idea is that since we want to minimize the loss, which is a function of Θ , we should move in the direction of the negative gradient of the loss with respect to Θ . the most naïve implementation of this is to simply advance in the direction of the negative gradient by a step given by the learning rate³ γ : $\Delta\Theta = -\gamma\nabla_{\Theta}\delta$. This is known as *batch gradient descent* and is the most basic form of gradient descent. This would leave us with the following update function:

$$\mathcal{U} = \Theta - \gamma\nabla_{\Theta}\delta. \quad (7.17)$$

There are several problems with this approach, the most important of them being:

- The learning rate γ is constant, which means that the step size is constant. This can lead to the algorithm getting stuck in local minima.
- The gradient is computed using the entire dataset, which can be very computationally expensive.

Stochastic gradient descent To our rescue comes stochastic gradient descent (SGD), which solves both of the above issues, by estimating the gradient using only a subset of the data, called a *mini-batch*. This means that the gradient is computed as an average over the mini-batch, which is a good approximation of the true gradient since the gradient of a dataset is statistically an expectation value. This allows us to increase the dataset without increasing the cost of computing the gradients. We may use several such minibatches chosen at random, say m batches. This would yield the following estimate for the gradient:

$$\nabla_{\Theta}\delta \approx \frac{1}{m} \sum_{i=1}^m \nabla_{\Theta}\delta_i, \quad (7.18)$$

where δ_i is the loss from minibatch i . This is known as *mini-batch gradient descent*. The update function for SGD is then identical to Eq. (7.17) but with the gradient estimate from Eq. (7.18). Choosing minibatches at random also greatly reduce the chances of being stuck in a local minimum.

Adding momentum In order to speed up convergence we may add a momentum term, which takes into account the previously found gradients. This is done by adding a momentum term, called the velocity \mathbf{v} , to the update function. We start with some initial velocity, and it is updated as $\mathbf{v} \leftarrow \alpha\mathbf{v} - \gamma\nabla_{\Theta}\delta$, where α is the momentum hyperparameter, which controls the weight of the previous velocity (and thus the previous gradients). This yields simply $\Delta\Theta = \mathbf{v}$ and the following update function:

$$\mathcal{U} = \Theta + \mathbf{v}. \quad (7.19)$$

[TODO: Nesterov momentum?](#)

Adaptive momentum One algorithm that adapts the momentum term is the *Adam* algorithm [TODO: ref Kinga and Ba](#) which uses the first and second moments of the gradient to adapt the momentum term. The first moment is the mean of the gradient and the second moment is the variance of the gradient.

Algorithm 2 Adam

Require: Learning rate γ , decay rates $\beta_1, \beta_2 \in [0, 1)$, first moment $\mathbf{m} = \mathbf{0}$, second moment $\mathbf{v} = \mathbf{0}$, weight decay rate ϵ (for numerical stability), initial parameters Θ_0 , and timestep $t = 0$.

```

1: while  $\Theta_t$  not converged do
2:    $t \leftarrow t + 1$  (Update timestep)
3:    $\mathbf{g}_t \leftarrow \nabla_{\Theta_{t-1}} \delta$  (Compute gradient at timestep  $t$ )
4:    $\mathbf{m}_t \leftarrow \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \mathbf{g}_t$  (Update biased first moment estimate)
5:    $\mathbf{v}_t \leftarrow \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) \mathbf{g}_t \odot \mathbf{g}_t$  (Update biased second moment estimate)
6:    $\hat{\mathbf{m}}_t \leftarrow \frac{\mathbf{m}_t}{1 - \beta_1^t}$  (Bias-corrected first moment estimate)
7:    $\hat{\mathbf{v}}_t \leftarrow \frac{\mathbf{v}_t}{1 - \beta_2^t}$  (Bias-corrected second moment estimate)
8:    $\Delta \Theta_t = -\gamma \frac{\hat{\mathbf{m}}_t}{\sqrt{\hat{\mathbf{v}}_t + \epsilon}}$  (Compute update)
9:    $\Theta_t \leftarrow \Theta_{t-1} + \Delta \Theta_t$  (Update parameters)
10: end while
    return  $\Theta_t$  (Resulting parameters)

```

The adam algorithm is given as: The hyperparameters β_1 and β_2 are typically set to 0.9 and 0.999 respectively, and ϵ is typically set to 10^{-8} . The learning rate γ is typically set to 0.001.

Learning rate schedules?**Batch normalization**

7.3 Challenges in optimization

Vanishing gradient The vanishing gradient problem is a problem that arises when using saturating activation functions, such as the sigmoid function. The problem is that the gradient of the loss with respect to the weights and biases of the network becomes very small, which means that the weights and biases are not updated. This is a problem because the network is not learning anything. Another cause of vanishing gradients are flat regions of the loss in parameter space which are neither local minima nor global minima. This is a problem because the gradient is zero in these regions, which means that the network is not learning anything. Examples of such points are saddle points and plateaus.

Exploding gradient Another issue are cliffs in parameter space, or very steep regions. This is a problem because the gradient is very large in these regions, which means that the network is learning too much. This can lead to the network diverging.

Local minima Being stuck in a local minimum means that the network appear to converge and learning stops. However, the optimal parameters are not found, and the converged model is not the best one. It is hard to know if the network is stuck in a local minimum or not.

³ γ is also often referred to as the step size

7.4 Regularization

In order combat some of the challenged posed in the previous section we use regularization. Regularization is a technique used to prevent overfitting, which is when the network learns the training data too well, and thus does not generalize well to unseen data. This is done by adding a penalty term to the loss function, which penalizes large weights and biases. This is done by adding a regularization term to the loss function:

7.4.1 Parameter norm penalties

7.4.2 Data augmentation

7.4.3 Early stopping

7.4.4 Parameter sharing

7.4.5 Ensemble methods

7.4.6 Dropout

Chapter 8

Convolutional Neural Networks

8.1 Convolution

8.1.1 Basic definitions

Mathematically, convolution is given as:

$$s(t) = (x * w)(t) = \int x(a)w(t - a)da. \quad (8.1)$$

This definition is purely mathematical. For machine learning purposes we refer to the function $x(a)$ as the *input* and $w(t - a)$ as the *kernel*. The output of the convolution is often called the *feature map*. This is easily generalized to higher dimensions. For example, in two dimensions the discrete convolution between an input image I and a kernel K is given as:

$$S(i, j) = (I * K)(i, j) = \sum_m \sum_n I(m, n)K(i - m, j - n). \quad (8.2)$$

If we choose to *flip the kernel* relative to the input I , the convolution becomes commutative:

$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i - m, j - n)K(m, n). \quad (8.3)$$

The latter is preferred in a machine learning context as there is less variation in the values of m and n because the size of the kernel is often much smaller than the size of the input. However, there is yet another quantity that is often used in machine learning, namely the *cross-correlation*, which is more intuitive in some sense, because the kernel is not flipped. The cross-correlation is given as:

$$S(i, j) = (K \star I)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n). \quad (8.4)$$

Most convolutional layers in machine learning use the cross-correlation. However, it is often referred to as just a convolution without the kernel flip, so this is the convention used further. An example of such a convolution (kernel not flipped) is illustrated in figure 8.1.

Whether the kernel is flipped does not play a huge role in a convolutional neural networks, because the kernel is learned by the network. However, it is important to be aware of the difference between the two definitions, because it can lead to confusion when interpreting both the output of the network and the convolutional layers themselves.

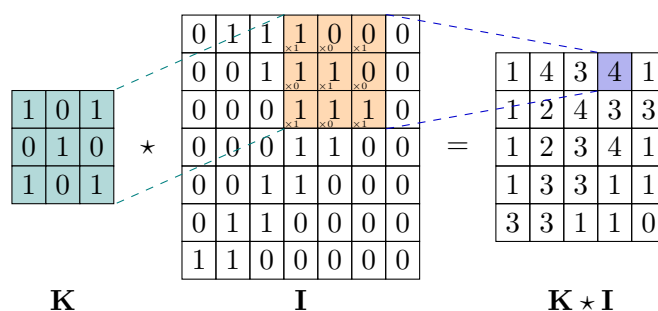


Figure 8.1: [TODO: fix this figure to be correct w.r.t. cross-correlation/convolution](#) Example of a cross-correlation operation, which we will also call a convolution. The input image is given by the matrix I and the kernel is given by the matrix K . The output is given by the matrix $K \star I$. The dashed lines indicate which elements of the input and kernel are multiplied together to produce the output. The figure is taken from TikZ.net [TODO: fix this citation](#)

8.1.2 Properties

Feature detection Different kernel has the ability to detect different features in the input. For example, different kernels can detect vertical edges, horizontal edges, corners or even more complex features. These kernels may be manually defined, however this defies the basic concepts of machine learning. These kernels are instead learned by the network, analogous to the weight matrices \mathbf{W} in a fully connected network. Hence, the output of a convolutional operation is a feature map.

Sparse interactions The kernels \mathbf{K} are usually a lot smaller than the input images \mathbf{I} . This means that each element of the kernel is used multiple times when producing the output feature map, which again allows for fewer parameters stored and faster computations.

Parameter sharing Exactly what it sounds like, the same parameters are re-used to produce more than one output. The learnable weights are the same for each element of the feature map. This is also a way to reduce the number of parameters stored and the number of computations needed.

Equivariance Convolutional layers are equivariant to translation. This means that if the input is translated, the output is translated by the same amount. This is a very useful property when dealing with images, because the location of a feature is not important. For example, if we want to detect an edge in an image, it does not matter where in the image the edge is located. All of the above properties are related to each other, but their combinations makes convolutional layers in a network very powerful.

8.2 New Layers

8.2.1 Convolutional layers

A convolutional layer is just a layer in the neural network that uses the convolution (or rather the cross-correlation) operation rather than just plain matrix multiplication. Both the input and output may consist of several planes or channels. For example, an input image may have three channels, one for each color (RGB). The output of a

convolutional layer may also have several channels. This is useful because it allows the network to learn different features in parallel. For example, one channel may learn to detect vertical edges, while another channel may learn to detect horizontal edges. The output of the convolutional layer is then a stack of feature maps, one for each channel. Let's consider convolutions in 2 dimensions, where we have C_{in} input channels and C_{out} output channels. The convolutional layer is then defined as:

$$\mathbf{S}(N_i, C_{\text{out}_j}) = \text{bias}(C_{\text{out}_j}) + \sum_{k=1}^{C_{\text{in}}} \mathbf{K}(C_{\text{out}_j}, k) \star \mathbf{I}(N_i, k), \quad (8.5)$$

where N_i is some batch of the dataset. In the case where we omit the bias for the convolutional layers and only have one input channel, this reduces to the simpler form:

$$\mathbf{S}(N_i, C_{\text{out}_j}) = \mathbf{K}(C_{\text{out}_j}) \star \mathbf{I}(N_i). \quad (8.6)$$

The above equation uses C_{out} different kernels to produce the same number of feature maps.

[TODO: Figures to show kernel size, padding, stride and dilation](#)

Kernel size The kernel size is usually chosen to be odd, because this allows for a central pixel. This is not necessary, but it is a common choice. The kernel size is usually chosen to be small, because this allows for more parameters to be stored in the network. However, the kernel size is usually chosen to be large enough to capture the features we want to detect.

Padding The padding is used to control the size of the output feature maps. If we do not use any padding, the output feature maps will be smaller than the input feature maps. This is because the kernel cannot be centered on the edges of the input feature maps. The padding is usually chosen to be zero padding, which means that the input feature maps are padded with zeros. This is the simplest choice, but there are other choices as well. For example, we may choose to pad the input feature maps with the edge values. This is called *edge padding*. Another choice is to pad the input feature maps with the reflection of the edge values. This is called *reflection padding*. The padding is usually chosen to be symmetric, which means that the same amount of padding is added to each side of the input feature maps. However, this is not necessary.

Stride The stride is used to control the size of the output feature maps. If we do not use any stride, the output feature maps will be smaller than the input feature maps. This is because the kernel cannot be centered on the edges of the input feature maps. The stride is usually chosen to be one, which means that the kernel is moved one pixel at a time. However, the stride can be chosen to be larger than one, which means that the kernel is moved several pixels at a time. The stride is usually chosen to be the same in both the horizontal and vertical directions, but this is not necessary.

Dilation The dilation is used to control the size of the output feature maps. If we do not use any dilation, the output feature maps will be smaller than the input feature maps. This is because the kernel cannot be centered on the edges of the input feature maps. The dilation is usually chosen to be one, which means that the kernel is moved one pixel at a time. However, the dilation can be chosen to be larger than one, which

means that the kernel is moved several pixels at a time. The dilation is usually chosen to be the same in both the horizontal and vertical directions, but this is not necessary.

All of the above would define how convolution should take place, and also control the dimensions of the resulting feature map. If we, along an axis, have an input dimension of D_{in} , a kernel size of K , a padding of P , a stride of S and a dilation of D , the output dimension, D_{out} , is given as:

$$D_{\text{out}} = \left\lfloor \frac{D_{\text{in}} + 2P - D(K - 1) - 1}{S} + 1 \right\rfloor \quad (8.7)$$

8.2.2 Pooling layers

Pooling is a form of dimensional reduction. It is used to reduce the size of the feature maps. This is useful because it reduces the number of parameters stored and the number of computations needed. Pooling is usually done in the spatial dimensions, but it can also be done in the channel dimension. Pooling often consists of replacing the output of a small region of the feature map with some summary statistics for that region. Examples of summary statistics used may be max pooling, average pooling or L^2 pooling. Max pooling is the most common choice. A consequence of pooling is invariance to local translation. Also note that pooling may not have to be a separate operation. Similar effects of dimension reduction and translation invariance may equally be obtained by choosing kernel size, padding, stride and dilation such that the dimension of the feature map is reduced.

Bibliography

- Adamek, Julian et al. (29th July 2016). ‘gevolution: a cosmological N-body code based on General Relativity’. In: *Journal of Cosmology and Astroparticle Physics* 2016.7, pp. 053–053. ISSN: 1475-7516. DOI: 10.1088/1475-7516/2016/07/053. arXiv: 1604.06065[astro-ph,physics:gr-qc,physics:physics]. URL: <http://arxiv.org/abs/1604.06065> (visited on 23/08/2023).
- Alqaraawi, Ahmed et al. (3rd Feb. 2020). *Evaluating Saliency Map Explanations for Convolutional Neural Networks: A User Study*. arXiv: 2002.00772[cs]. URL: <http://arxiv.org/abs/2002.00772> (visited on 23/10/2023).
- Chisari, Nora Elisa and Matias Zaldarriaga (2nd June 2011). ‘Connection between Newtonian simulations and general relativity’. In: *Physical Review D* 83.12, p. 123505. ISSN: 1550-7998, 1550-2368. DOI: 10.1103/PhysRevD.83.123505. arXiv: 1101.3555[astro-ph,physics:gr-qc]. URL: <http://arxiv.org/abs/1101.3555> (visited on 23/08/2023).
- Goodfellow, Ian, Yoshua Bengio and Aaron Courville (2016). *Deep Learning*. <http://www.deeplearningbook.org>. MIT Press.
- Green, Stephen R. and Robert M. Wald (15th Mar. 2012). ‘Newtonian and Relativistic Cosmologies’. In: *Physical Review D* 85.6, p. 063512. ISSN: 1550-7998, 1550-2368. DOI: 10.1103/PhysRevD.85.063512. arXiv: 1111.2997[astro-ph,physics:gr-qc]. URL: <http://arxiv.org/abs/1111.2997> (visited on 23/08/2023).
- Hastie, T., R. Tibshirani and J.H. Friedman (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*. Springer series in statistics. Springer. ISBN: 9780387848846. URL: <https://books.google.no/books?id=eBSgoAEACAAJ>.
- Jeong, Donghui, Fabian Schmidt and Christopher M. Hirata (4th Jan. 2012). ‘Large-scale clustering of galaxies in general relativity’. In: *Physical Review D* 85.2, p. 023504. ISSN: 1550-7998, 1550-2368. DOI: 10.1103/PhysRevD.85.023504. arXiv: 1107.5427[astro-ph]. URL: <http://arxiv.org/abs/1107.5427> (visited on 23/08/2023).
- Selvaraju, Ramprasaath R. et al. (Feb. 2020). ‘Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization’. In: *International Journal of Computer Vision* 128.2, pp. 336–359. ISSN: 0920-5691, 1573-1405. DOI: 10.1007/s11263-019-01228-7. arXiv: 1610.02391[cs]. URL: <http://arxiv.org/abs/1610.02391> (visited on 23/10/2023).