

CKi ein CNN-Model in C++

by Simeon Stix

22. Februar 2024

Version: 3

Betreuer: Sven Nüesch

Inhaltsverzeichnis

1	Analyse	4
1.1	Vorhaben	4
1.2	Zielgruppe	4
1.3	Anforderungen	5
1.3.1	Must-haves	5
1.3.2	Nice-to-haves	6
1.3.3	Use Cases	7
1.3.4	Use Case Diagramme	10
1.3.5	Ablaufdiagramme	11
1.4	Umsetzung	14
1.4.1	Entwicklungsumgebung	14
1.4.2	CNN	15
1.4.3	Testdaten	16
1.5	Abwägung des Nutzens und der Alternativen	17
1.5.1	Alternativen	17
1.5.2	Kosten	17
1.5.3	Nutzen	18
1.5.4	Effizienz	18
1.5.5	Vergleich	18
1.5.6	Nutzwertanalyse	19
2	Planung	21
2.1	Rollenverteilung	21

2.2	Aufgabenliste	21
2.3	Meilensteine	23
2.4	Gantt	23
3	Design	28
3.1	Konsole	28
3.1.1	Training	28
3.1.2	Verify	28
3.1.3	Prediction	29
3.1.4	Help	29
3.2	Datenbank	30
3.2.1	UByte	30
3.3	Code	31
3.3.1	Klassendiagramm	31
3.3.2	Trainingsdaten	32
3.3.3	Tests	33
4	Realisation	34
4.1	Allgemein	34
4.2	Abhängigkeiten	34
4.3	Architektur	36
4.3.1	Klassen	36
4.3.2	Ordnerstruktur	38
4.4	Code	39
4.4.1	Konzept	39
4.4.2	Algorithmen	40
4.4.3	Schlüsselpassagen & Snippets	49
5	Testing	58
5.1	Unittests	58
5.1.1	Util (Utility-Klasse)	58
5.1.2	Network	59

5.1.3	Layer	59
5.1.4	Neuron	60
5.2	Integrationstests	60
5.3	Deployment-Tests	61
6	Deployment	62
6.1	Building	62
6.1.1	Anforderungen	62
6.1.2	Ausführung	64
6.2	Installation	64
6.2.1	Installation mit CMake	65
6.2.2	Installation ohne CMake	65
6.2.3	Hinzufügung von vortrainierten Gewichtungen und Biases	65
6.3	Verwendung	66
6.3.1	Interaktion	66
	Literaturverzeichnis	68
	Abbildungsverzeichnis	70
	Tabellenverzeichnis	71
	Codeverzeichnis	72
	Anhang	73
A	Fortschrittsberichte	74

1 Analyse

1.1 Vorhaben

Die Aufgabe, die sich das Projekt CKi stellt, ist die Wissenserweiterung des Entwicklers. Dabei wird ein Programm geschrieben, welches einzelnen handgeschriebenen Ziffern erkennen kann. Dies geschieht mittels KI. Dabei wird die KI komplett vom Entwickler geschrieben. Da keine modernen externen Grundlagen verwendet werden, wird das Produkt, das Programm, nicht die Geschwindigkeit einer modernen KI erreichen. Dies spielt jedoch keine Rolle, da dieses Projekt nicht wegen des Endprodukts durchgeführt wird.

1.2 Zielgruppe

Das Projekt CKi ist als solches nicht ausgelegt einem realen Anwendungszweck zu entsprechen oder eine Lösung oder einen Lösungsansatz für einen solchen zu bieten. Diesbezüglich liegt der einzige Nutzen von CKi nicht in dessen Produkt, sondern nur im Wissensgewinn und Verständnisgewinn für den Entwickler in den Bereichen der künstlichen Intelligents oder genauer im Bereich des maschinellen Lernens mit einem *Convolutional Neural Network*, der Realisation von Anwendungen mit C++ und dessen Möglichkeiten Hardware direkt in die programminternen Abläufe einzubinden. Somit richtet sich CKi nicht nach dem Grundsatz ein bestmöglich nutzbares Produkt zu sein, sondern lediglich nach dem grössten Wissensgewinn für den Entwickler. Nach diesem Grundsatz ist die resultierende Zielgruppe der Entwickler und vereint so mehrere Rollen des Projektes CKi in einer Person.

1.3 Anforderungen

1.3.1 Must-haves

Die Must-haves wurden aus dem Themenblatt, welches am 15.09.2023 bei Walter Schnyder eingereicht wurde, übernommen und mit weiterführenden Elaborationen versehen.

- **Rückgabe in Prozentwerten, die die Wahrscheinlichkeit der Übereinstimmung mit dem digitalen Gegenstück der handgeschriebenen Zahl abbilden.** „Welche Zahl wurde (vermutlich) aufgeschrieben?“

Erläuterung: Da ein simples neuronales Netzwerk für maschinelles Lernen durch ein „Netz“ aus Knotenpunkten gebaut wird und jeder dieser Knotenpunkte, auch die Knotenpunkte, welche bei einem solchen neuralen Netzwerk als Endschnittstellen fungieren, einzeln berechnet werden, erhält man, bei Anwendungsfall von CKi, eine multiple Anzahl von Prozentzahlen, welche zur Interpretation des gelieferten Endergebnisses verwendet werden können. Diese Rückgabe der einzelnen Prozentwerte erfolgt zum Beginn über eine Konsolenausgabe. Diese wird später, wie in 1.3.2 („Nice-to-haves“) unter GUI erläutert, in ein grafisches Nutzerinterface integriert und zu diesem Zeitpunkt evtl. auch interpretiert (wobei die einzelnen Prozentwerte weiterhin einsichtig bleiben sollten).

- **funktionales neuronales Netzwerk (trainiert auf Zahlenwert)**

Erläuterung: Ein CNN-Algorithmus oder auch Convolutional Neural Network wird beim maschinellen Lernen oft bei der Interpretation von Bildern genutzt. Dabei wird das Bild in kleinere Abschnitte unterteilt und „einzeln“ an den gehirnähnlich aufgebauten Algorithmus weitergegeben.

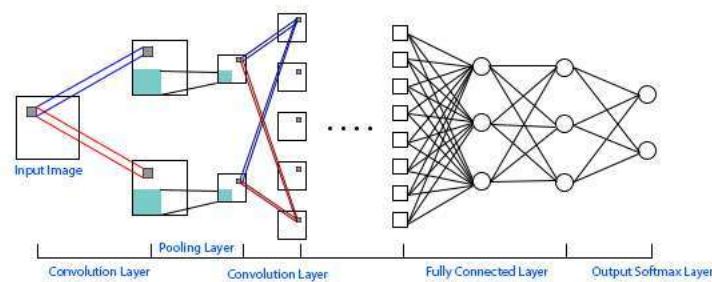


Abbildung 1.1: CNN Model Aufbau Grafik von researchgate.net

Im Projekt CKi wird dieser wie im Themenblatt beschrieben mit Bildern von einzelnen handgeschriebenen Ziffern trainiert, dem MNIST-Datensatz.

- **Nutzer-Input | Nutzer darf eine Zahl zeichnen**

Erläuterung: Ein solches CNN-Model zu erbauen und zu trainieren ist zwar die Grundlage für dieses Must-have, jedoch sollte das Produkt auch erprobbar sein. Diesbezüglich muss der Nutzer in der Lage sein, eine handschriftliche Ziffer an das neurale Netzwerk zu liefern. Hierbei ist die minimale Anforderung, dass der Nutzer in einer anderweitigen Applikation ein solches Bild erstellt hat und es nun interpretieren lassen kann. Wie in 1.3.2 („Nice-to-haves“) unter GUI beschrieben, wird diese Eingabemöglichkeit (eine Ziffer zu zeichnen oder eine schon vorhandene grafische Abbildung zu verwenden) in einem weiteren Entwicklungsschritt direkt in die zukünftige grafische Nutzeroberfläche der Applikation integriert.

1.3.2 Nice-to-haves

Die Nice-to-haves wurden aus dem Themenblatt, welches am 15.09.2023 bei Walter Schnyder eingereicht wurde, übernommen und mit weiterführenden Elaborationen versehen. Zudem behalte ich mir als Verfasser dieses Dokumentes, als Entwickler des Projektes CKi und als Zielgruppe des Projektes CKi vor, diese Liste in gegebenen Fall zu erweitern.

- **GUI**

Erläuterung: Ein GUI (oder auch grafisches User-Interface) ist die grafische Nutzeroberfläche der Applikation. Da die Applikation primär dem Wissensgewinn gewidmet ist

und dementsprechend nicht für den Nutzer optimiert wird, hat eine solche Erweiterung nur eine geringe Priorität. Diese Nutzeroberfläche wird selbst bei der Umsetzung in einer möglichst simplen Form gehalten. Dabei sollte es folgende Bestandteile beinhalten:

- Eine Möglichkeit für den Nutzer eine anderweitig gezeichnete Ziffer interpretieren zu lassen
- Eine Möglichkeit für den Nutzer eine Ziffer zu zeichnen.
- Eine (weiter-) interpretierte Ausgabe der Interpretation des CNN.
- Eine Ausgabe der nicht interpretierten Prozentwerte der Interpretation des CNN.

Bis zu dem Punkt, wo ein GUI realisiert wurde und in den Einsatz gestellt wird, sind nicht alle dieser Funktionen über die Konsole (das Interface zum Programm, welches vor dem GUI zum Einsatz kommt) verfügbar.

• GPU als Berechnungsplattform nutzen

Erläuterung: Die CPU in einem Computer ist eine sehr „fokussierte“ Hardware. So kommt es dazu, dass diese immer nur einen einzelnen Prozess berechnen kann. Dies ist für ein neuronales Netzwerk, welches Hunderte oder Tausende Knotenpunkte hat und alle einzeln berechnet werden müssen, äusserst hinderlich. Eine dezidierte Grafikkarte, wenn vorhanden, kann diesem Geschwindigkeitsverlust nachhelfen, da eine solche GPU in der Lage ist, tausende Berechnungen gleichzeitig zu tätigen. Da im Projekt mit C++, einer Hardware-nahen Programmiersprache, gearbeitet wird, kann eine Anbindung an das Rechenpotential der GPU implementiert werden. Da dies jedoch ein äusserst schwieriger Prozess ist, sich diese Anbindung bei den unterschiedlichen Herstellern von GPUs unterscheiden kann und nicht notwendig ist, um ein funktionierendes CNN zu erstellen, wurde diese Optimierung als nicht notwendig eingestuft.

1.3.3 Use Cases

Die hier aufgelisteten Use Cases entsprechen den Use Cases nach den Must-haves und sind diesbezüglich ohne die grafische Nutzeroberfläche. Dies kann dazu führen, dass ein endgültiges Produkt nicht mehr kohärent zu den Use Cases steht. Im Allgemeinen sollte aber selbst

eine solche Inkohärenz nicht in extremer Weise auftreten, da die Bedienung der Applikation als Konsole als auch als grafische Oberfläche in ähnlicher Weise auftreten sollte.

1. Prediction

- **Akteur:** Nutzer, CNN
- **Ablauf:** Siehe Ablaufdiagramm 1.3.5.1 („Interpretation“).
- **Nachbedingungen:** Der Nutzer sollte in der Konsole eine Auflistung aller möglichen Ziffern (0–9) und deren entsprechende Wahrscheinlichkeit sehen. Zudem kann der Nutzer auch die kongruierende Zahl mit der höchsten Wahrscheinlichkeit auf einer speziellen Zeile in der Konsole ablesen.
- **Ausnahmen:** Bei der Interpretation von benutzereigenen Abbildungen kann es zu multiplexen Fehlern kommen. Diese reichen von falscher Dateikodierung zu falscher Auflösung. Aufgrund dieser mannigfaltigen Möglichkeiten zu Fehlern können diese nicht alle hier erläutert werden.
- **Anmerkungen:** Offiziell ist zwar „Training“ keine Vorbedingung von „Testen“, jedoch macht die Anwendung der Applikation keinen Sinn, wenn man nicht erwarten kann, dass man ein realistisches oder sinnvolles Ergebnis erhält.

2. Training

- **Akteur:** Nutzer, CNN
- **Ablauf:** Siehe Ablaufdiagramm 1.3.5.2 („Training“).
- **Nachbedingungen:** Der Nutzer erhält nach der Beendigung der Schulung des CNN-Models eine kurze Benachrichtigung, dass das Training abgeschlossen ist. Zudem erhält der Nutzer nach jedem einzelnen Datensatz den Output, dass nun x von y Datensätzen bearbeitet wurden (kann auch als Lade-Balken oder Ähnliches (Ladetext) implementiert werden).
- **Ausnahmen:** Hierbei kann es zu zwei Fehlern kommen. Dabei handelt es sich um Fehler bei den Datensätzen. Der eine Fehler ergibt sich aus der Möglichkeit, dass CKi die entsprechende Datei aus befindlichen Gründen nicht öffnen kann.

Der andere Fehler ergibt sich, wenn der Datensatz nicht dem Standard der MNIST-Datensätze im UByte-Format entspricht.

Natürlich können auch andere Fehler eintreten, diese können nicht zu diesem Zeitpunkt abgeschätzt werden.

3. Testing

- **Akteur:** Nutzer, CNN
- **Ablauf:** Siehe Ablaufdiagramm 1.3.5.3 („Testen“).
- **Nachbedingungen:** Der Nutzer erhält eine Benachrichtigung, dass x von y Datensätzen bearbeitet wurden (kann auch als Lade-Balken oder Ähnliches (Ladetext) implementiert werden). Nach Beendigung erhält der Nutzer die Mitteilung, dass das Testen abgeschlossen ist und einen Prozentwert mit der Genauigkeit des Netzwerkes.
- **Ausnahmen:** Hierbei kann es zu zwei Fehlern kommen. Dabei handelt es sich um Fehler bei den Datensätzen. Der eine Fehler ergibt sich aus der Möglichkeit, dass CKi die entsprechende Datei aus befindlichen Gründen nicht öffnen kann. Der andere Fehler ergibt sich, wenn der Datensatz nicht dem Standard der MNIST-Datensätze im UByte-Format entspricht. Natürlich können auch andere Fehler eintreten, diese können nicht zu diesem Zeitpunkt abgeschätzt werden.
- **Anmerkungen:** Bei den Testdaten würde es evtl. Sinn ergeben, die einzelnen Konsolen Outputs auch in einer CSV-Datei festzuhalten, wobei dies mit Pipes in der Konsole dem Nutzer bereits offensteht. Offiziell ist zwar „Training“ keine Vorbedingung von „Testen“, jedoch ergibt das Testen nur begrenzt Sinn, wenn es noch nichts gibt, was sich zu testen lohnt.

1.3.4 Use Case Diagramme

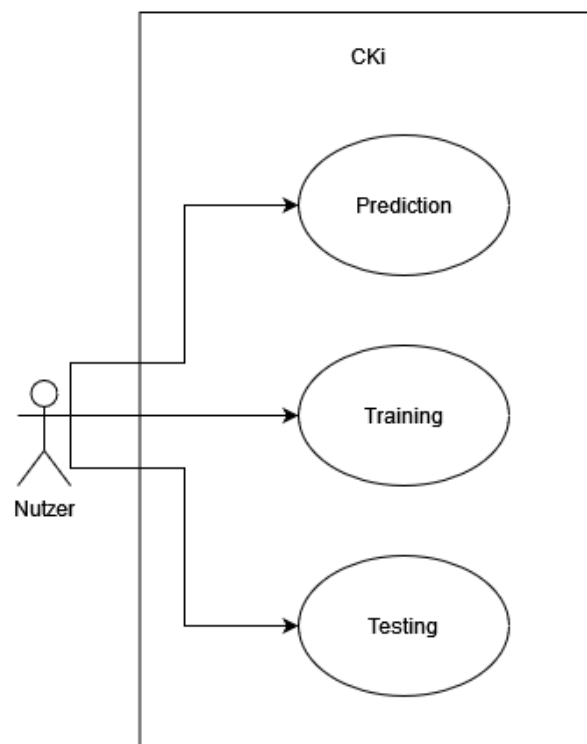


Abbildung 1.2: Use-Case-Diagramm

1.3.5 Ablaufdiagramme

1.3.5.1 Interpretation

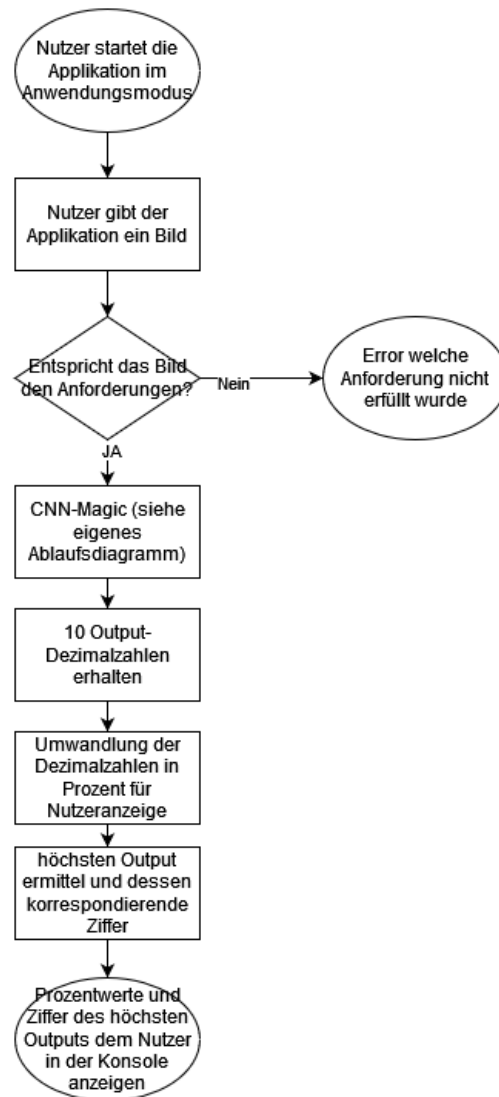


Abbildung 1.3: Ablaufdiagramm für die Interpretation von Nutzer-gelieferten Einzelbildern

1.3.5.2 Training

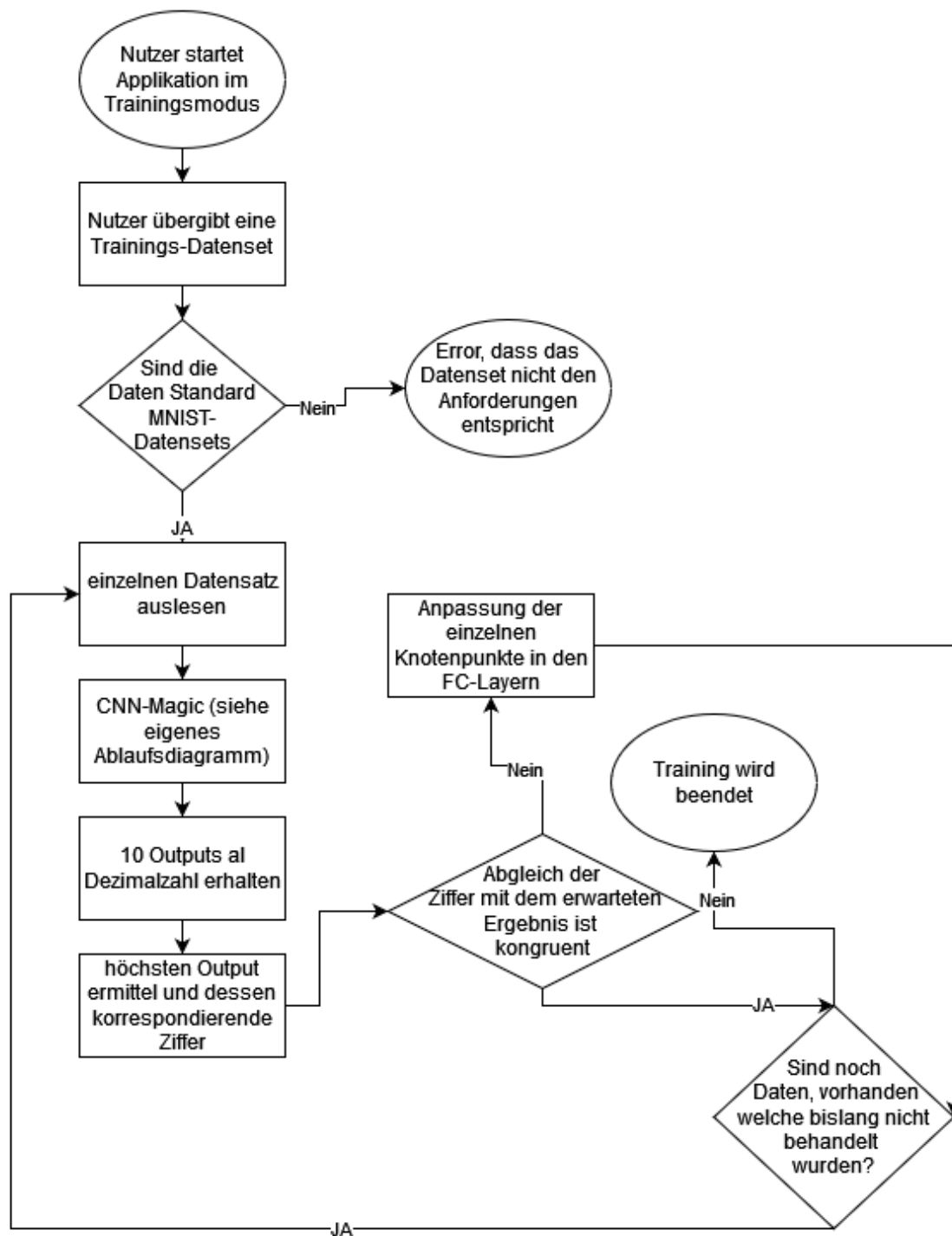


Abbildung 1.4: Ablaufdiagramm vom Training des CNN

1.3.5.3 Testen

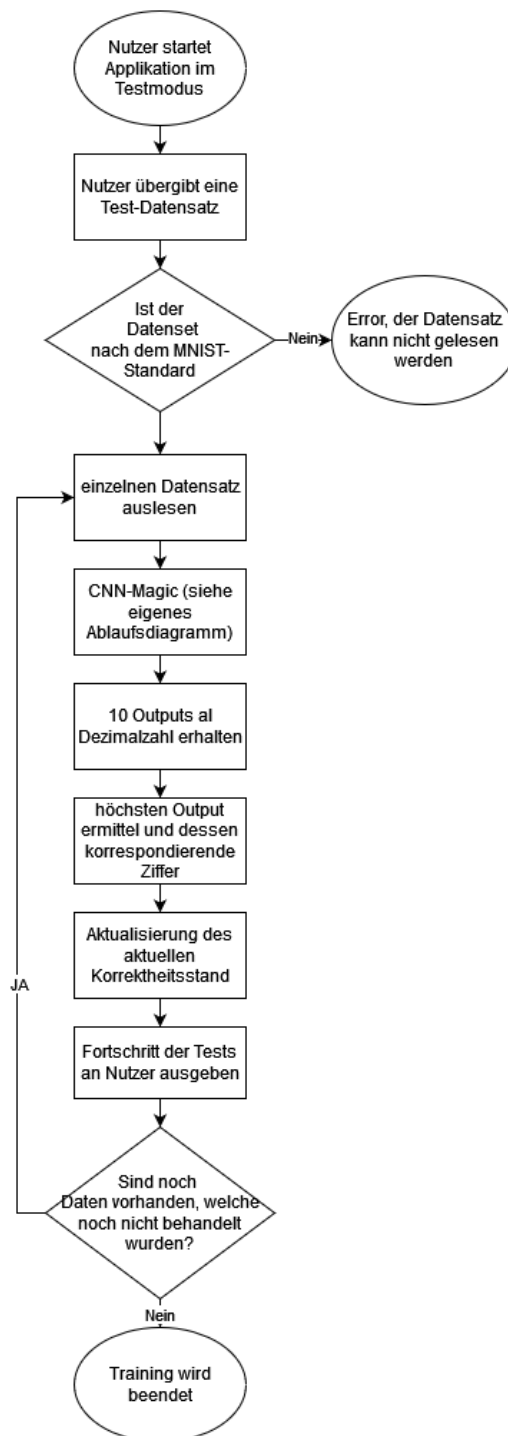


Abbildung 1.5: Ablaufdiagramm von einem Test des CNN

1.3.5.4 CNN-Magie

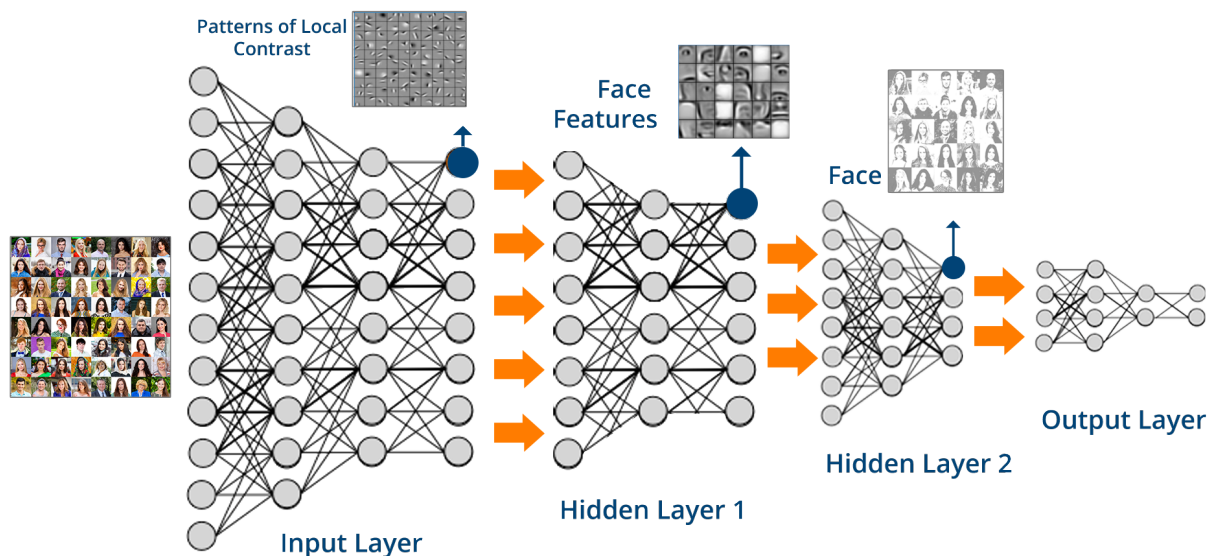


Abbildung 1.6: Visuelle Darstellung der Funktionsweise eines CNN von blog.goodaudience.com

1.4 Umsetzung

Wie im Themenblatt vom 15.09.2023 wird für die Umsetzung reines C++ verwendet.

1.4.1 Entwicklungsumgebung

Bei der Entwicklungsumgebung wird das JetBrains Produkt **CLion** zum Einsatz kommen. Dieses beinhaltet alle benötigten Funktionalitäten, die bei der Entwicklung des Produktes nötig sind. Gegebenen Falles könnte noch **SQLite Browser** zur Verwendung kommen, da dieses ein besseres (und für den Entwickler ein gewohntes) Interface zur Handhabung von lokalen Datenbanken bietet. (Dieses Produkt kam nicht zum Einsatz, da keine Datenbank verwendet wurde.) Für die Version Control des Projektes wird lokal **Git** verwendet und zur Sicherung in der Cloud wird sowohl **GitHub** als auch **GitLab** verwendet. Die Entwicklungsumgebung mit

IDE und Version Control wird primär auf dem Betriebssystem Windows verwendet. Es kann jedoch nicht garantiert sein, dass die Entwicklung nicht teilweise unter einer Distribution von Linux ablaufen wird. Hierbei sollten (evtl. abgesehen von der GUI) keine Kompatibilitätsprobleme auftreten.

1.4.1.1 Auflistung Versionen

- JetBrains Toolbox 2.1.1.18388, Windows 11, x64
- CLion (von JetBrains Toolbox) 2023.2.2
- Git 2.39.2.windows.1
- CMake 3.27.7
- GCC 13.2.0
- MiKTeX 23.10
- TeXnicCenter 2.02

1.4.2 CNN

1.4.2.1 Grundlagen

Der Kern des CNN lässt sich in drei unterschiedliche Arten von Schichten aufteilen. Zusätzlich gibt es noch zwei zusätzliche „Hilfs“-Schichten, welche für die Ein- und Ausgaben zuständig sind. Die *Eingabeschicht* ist eine dieser zwei erwähnten Schichten. Diese akzeptiert im Falle von CKi jeweils ein Bild von speziell definierten Grössen und nimmt jeden Pixel als Graustufen (Dies ist entscheidend, weil so der Farbwert des Pixels von drei zwei Byte grossen Zahlen zu nur einer zwei Byte grossen Zahl reduziert wird.) als Eingabe entgegen. Auf die Eingabeschicht folgt die *Convolutional Layers* oder auch *Faltungsschichten*. Diese Faltungsschichten sind einer der Schlüsselbestandteile eines CNN. Die Convolutional Layer dienen dazu, bestimmte Merkmale und Schlüsselemente aus dem Bild zu extrahieren. Hierbei wird bei jedem der Convolutional Layer eine mathematische Faltungsoperationen durchgeführt. Nach

den Faltungsschichten kommen die simplen *Pooling-Schichten*. Bei einer Pooling-Schicht wird die Dimension eines Bildes reduziert (Downsampling). Dies kann durch zwei Arten geschehen, entweder durch Max-Pooling, dabei wird aus einer Liste von Werten nur der höchste übermittelt oder durch Durchschnitts-Pooling, wobei der Durchschnitt dieser Liste berechnet und übermittelt wird. Durch diese Datenreduktion wird Rechenleistung gespart und so das Ergebnis schneller und stabiler geliefert. Vor der Ausgabeschicht gibt es noch die *FC Layers* oder auch *Vollständig verknüpfte Schichten (Hidden Layers)*. Bei diesen ist jeder Knotenpunkt mit jedem Knotenpunkt in der nächsten Schicht verbunden. Dies ist das Herz des gesamten Models. Dabei wird in jedem Knotenpunkt ein neuer Wert berechnet, um in der letzten Schicht, der *Ausgabeschicht* diese in Wahrscheinlichkeit zu „konvertieren“.¹

1.4.2.2 Implementation

Bei der Implementation eines neuronalen Netzwerkes wurde für CKi der Entschluss getroffen, eine C++-Klasse für jede benötigte Art von Komponente, die die Applikation benötigt, zu kreieren. Konkret bedeutet dies, dass eine Klasse für das Neuron, eine für den Layer und eine für das Netzwerk entstehen muss. Dabei werden die Klassen aufeinander aufgebaut; ein Netzwerk hat multiple Layer, ein Layer hat multiple Neuronen. Dieses Netzwerk wird danach in die Main-Funktion, den Programmeinstieg, integriert. Dies erleichtert es, in einem späteren Schritt die grosse Umstellung von einer Konsolen-Anwendung hin zu einer grafisch basierten Anwendung (auch wenn immer noch grosse Teile der Main-Klasse ausgetauscht werden müssen). Für die Speicherung der Berechnungswerte für die FC-Layer ist die Abwägung zu treffen, ob es sinnvoll ist, diese in einer lokalen Datenbank zu speichern oder in eine konkret hierfür entworfene Datei zu schreiben.

1.4.3 Testdaten

Die Testdaten für CKi können in unzähliger Ausführung auf Kaggle gefunden werden. Es wird jedoch der Standard von <http://yann.lecun.com/exdb/mnist/> verwendet.

¹vgl. Pathak, 2023; IBM, 2021; Saha, 2018; 3Blue1Brown, 2017

1.4.3.1 Speicherung

Im gegebenen Fall würde es für die Schulung, des CNN, Sinn ergeben, die einzelnen Bilder nicht über Ordner- oder Archivstrukturen zu speichern, sondern in eine oder mehrere lokale Datenbanken abzuspeichern.

Im Falle von <http://yann.lecun.com/exdb/mnist/> ist dies jedoch nicht notwendig, da diese im UByte-Format vorliegen, wie in 3.2.1 beschrieben.

1.5 Abwägung des Nutzens und der Alternativen

1.5.1 Alternativen

Die möglichen Alternativen zu einer kompletten Realisation eines CNN oder eines neuronalen Netzwerks in C++ ohne Bibliotheken oder Frameworks für neuronale Netze sind endlos. Wenn gleich man bereits bestehende Produkte begutachtet, wird man auf GitHub schnell fündig. Bei einer Suche auf GitHub mit dem Query „mnist digit recogniser“ findet man nach Stand 20.10.2023 02:13 178 Repositorien (je eines in C, C++, C#, Rust, drei in Java + Kotlin und 53 in Python). Auch unter dem Beschluss, dass man das CNN selbst bauen wolle, würde man für jede beliebige Sprache eine Bibliothek oder Framework finden, das einem diese Aufgabe erleichtern würde.

1.5.2 Kosten

Bezüglich der Kosten kann das Projekt CKi nicht mit den Alternativen mithalten. C++ ist keine Sprache, in welcher man schnell ein Produkt hervorbringt. Dies belegen etliche Studien und Artikel (Programming Languages Table, An Empirical Comparison of Seven Programming Languages, Programming Languages by Energy Efficiency, Development time in various languages). Zu dem, dass keine dezidierten Bibliotheken für maschinelles Lernen zum Einsatz kommen, erhöht den Entwicklungsaufwand, Zeitaufwand und so die Kosten. Wenn man die Stundenzahl aus der Planung mit einem Stundensatz von 50 CHF quantifizieren, käme man auf einen gesamten Kostenaufwand von 8'600.- + Lizenzkosten. Zuzüglich kämen noch Hard-

ware Anschaffungen für die Entwicklung hinzu, diese können jedoch im Falle vom Projekt CKi vernachlässigt werden.

Senior Developer	2h	100.00 Fr.
Junior Developer	89h	4'450.00 Fr.
Projektleitung	39h	1'950.00 Fr.
Hardware	-	0.00 Fr.
Software	Lizenz	827.00 Fr.
Möglicher Zusatz-Entwicklungskosten	42h	2'100.00 Fr.
Total		9'427.00 Fr.

1.5.3 Nutzen

Wie bereits im Abschnitt Zielgruppe 1.2 erwähnt, ist dieses Produkt auf den Wissensgewinn ausgerichtet und nicht auf das Produkt selbst. Dementsprechend ist in der geplanten Umsetzung des Projektes CKi der maximale Nutzen gewährleistet. Leider lässt sich dieser Nutzen nicht/schlecht quantifizieren.

1.5.4 Effizienz

Der Vergleich in Effizienz der Anwendung, der aus dem Projekt CKi resultiert und den professionell entwickelten Bibliotheken in Python, etc., wird CKi selbst mit GPU-Berechnung und dem Geschwindigkeitsvorteil von C++ gegen diese Bibliotheken in Bezug auf Run-Time-Length verlieren.

1.5.5 Vergleich

Bestehendes Produkt		Eigene Kreation	
Sektion	Kosten	Sektion	Kosten
Anschaffungskosten	0.- Fr.	Kreation	8'600.- Fr.

Wiederholende Kosten	0.- Fr.	wiederholende Kosten	0.- Fr.
Absolut	0.- Fr.	absolut	8'600.- Fr.

Eine kurze Zusammenfassung der Kosten-Nutzen-Analyse zeigt den Zustand von CKi.

Es werden Kosten für die Arbeitszeit von 8'600 CHF anfallen. Die Effektivität ist niedriger, als wenn dasselbe Produkt mit einer vorhandenen Bibliothek geschrieben werden würde (und die Entwicklungsdauer wäre ebenfalls geringer). Des Weiteren gibt es unzählige kostenlose Alternativen, die es nur zu verwenden gilt.

Der Nutzen des Endproduktes ist nichtig, da dieses nie zur wirklichen Anwendung kommen wird.

1.5.6 Nutzwertanalyse

Bei der Nutzwertanalyse werden unterschiedliche Alternativen systematisch und statistisch nach gewissen Kategorien verglichen. Konkret werden hier unterschiedliche Möglichkeiten der Implementierung miteinander verglichen, dabei wird eine Punkteskala von eins bis zehn verwendet. Wichtig bei der Auflistung ist es, dass eine höhere Punktzahl nicht bedeutet, wie hoch diese Kategorie ist. Eine Entwicklungsdauer von 10 bedeutet nicht, dass diese besonders lang ist, sondern dass diese Option im Bereich Entwicklungsdauer zehn Punkte erhalten hat (also besonders schnell zum implementieren ist).

Kriterien		Python mit TensorFlow		C++	
Entwicklungsdauer	10 %	10	1	2	0.2
Komplexität	15 %	10	1.5	4	0.6
Sprachkenntnisse	10 %	3	0.3	4	0.4
Frust	10 %	5	0.5	7	0.7
Typisierung	5 %	1	0.05	9	0.45
Run-Time-Length	10 %	10	1	9	0.9
Testing	5 %	6	0.3	7	0.35
Objektorientiert	5 %	5	0.25	9	0.45

Wissensgewinn	30 %	2	0.6	10	3
Total	100 %	52	5.5	61	7.05

2 Planung

2.1 Rollenverteilung

Rolle	Person
Leiter	Simeon Stix
Entwickler	Simeon Stix
Betreuer	Sven Nüesch

2.2 Aufgabenliste

1. *Planung:*

- Aufgabenliste erstellen
- Zeiteinteilung
- Meilensteine definieren
- Gantt erstellen

2. *Analyse:*

- Recherche CNN
- Recherche Testdaten
- Recherche C++ Details
- Anforderungen definieren

- Use Cases
- weitere Analysen
- Analyse schreiben

3. **Design:**

- GUI konzipieren
- Wireframes zeichnen
- Konsolenbefehle definieren
- Konsolen-Output definieren

4. **Realisation:**

- Input Layer realisieren
- Convolutional Layer realisieren
- Pooling Layer realisieren
- Dense Layer Klasse schreiben
- Output Layer erstellen
- CNN-Model erstellen
- CNN-Model trainieren
- GUI auf CNN-Model setzen
- GPU-Anbindung implementieren
- Kommentare schreiben
- weitere Verbesserungen vornehmen
- Puffer-Zeit fürs Programmieren

5. **Dokumentieren:**

- Schreiben

6. **Testen:**

- Testliste kreieren
- Tests implementieren
- Tests durchführen

7. **Deployment:**

- Projekt als EXE bauen
- Build-Anleitung kreieren

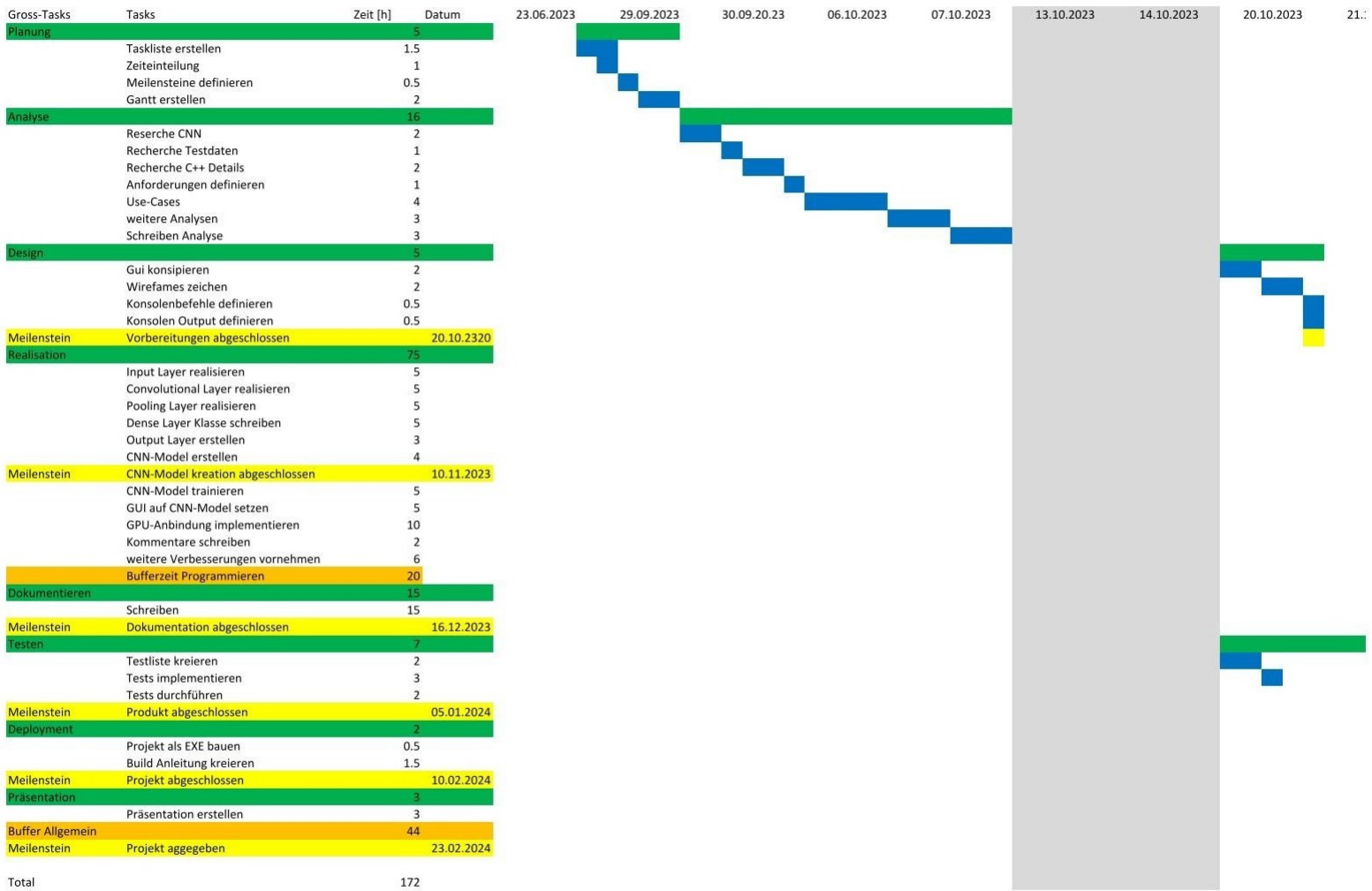
8. **Präsentation:**

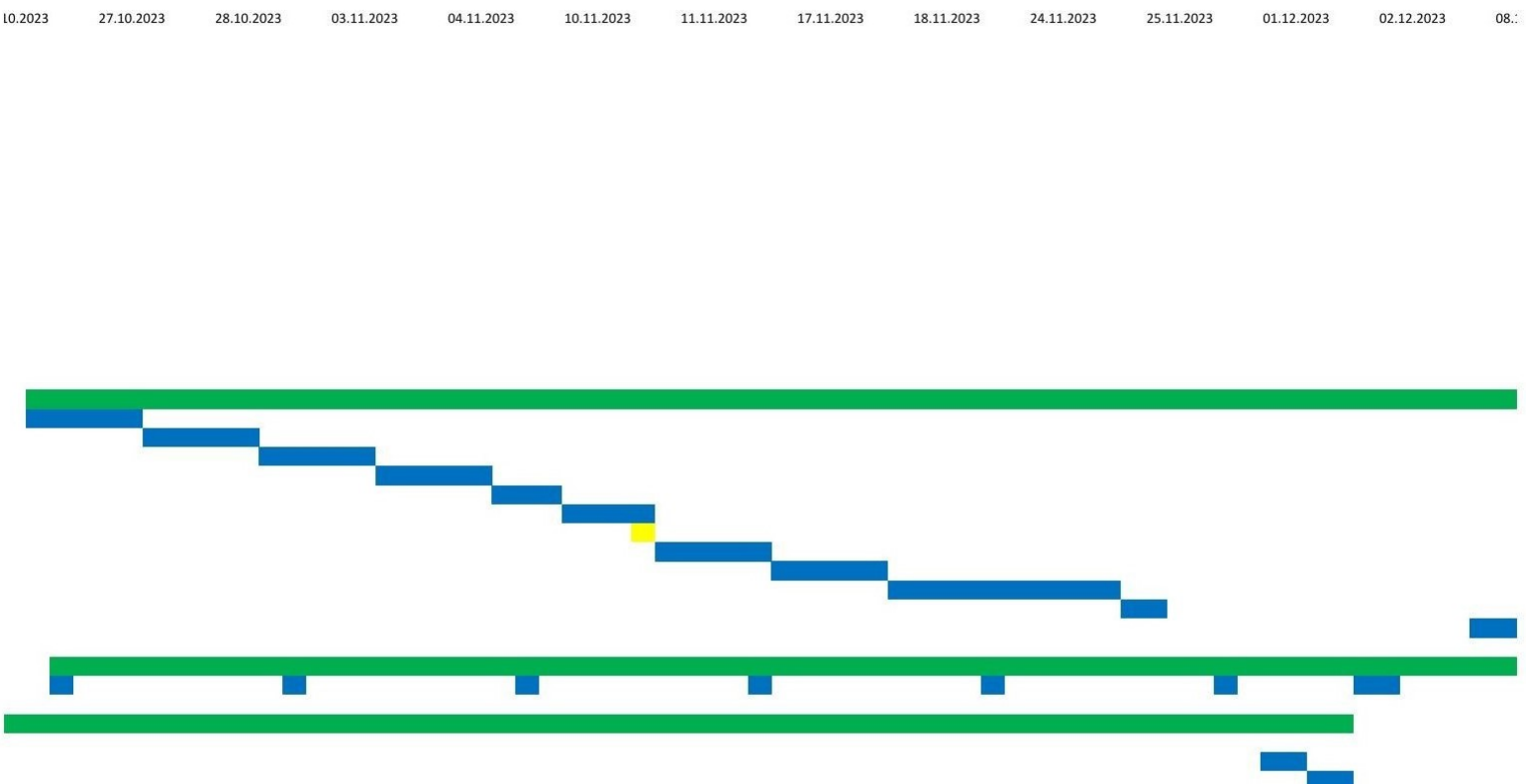
- Präsentation erstellen

2.3 Meilensteine

Meilenstein-Name	Datum
Vorbereitungen abgeschlossen	20.10.2023
CNN-Model Kreation abgeschlossen	10.11.2023
Dokumentation abgeschlossen	01.12.2023
Produkt abgeschlossen	05.01.2024
Projekt abgeschlossen	10.02.2024
Projekt abgegeben	23.02.2024

2.4 Gantt







31.2024 26.01.2024 27.01.2024 02.02.2024 03.02.2024 09.02.2024 10.02.2024 16.02.2024 17.02.2024 23.02.2024



3 Design

3.1 Konsole

Das Projekt CKi ist in seinem Grundaufbau für die Konsole konzipiert. Der grundlegende Befehl soll wie folgt aussehen:

```
1 $ ki [options] [file]
```

Dabei gibt es drei Formen von diesem Befehl.

3.1.1 Training

3.1.1.1 Input

```
1 $ ki --train -l [labels] -i [images] -e [epochs] -lr [learningrate]
2 $ #Dieser Befehl nimmt ubyte-Dateien.
3 $ #Die Optionen -e und -lr sind optional.
```

3.1.1.2 Output

Nach jedem Datensatz wird angezeigt, wie viele der zum Training gegebenen Daten bereits abgearbeitet wurden (Dies wird evtl. nicht implementiert.). Am Ende wird ein progressiv berechneter Fehler angezeigt.

3.1.2 Verify

3.1.2.1 Input

```
1 $ ki --verify -l [labels] -i [images]
2 $ #Dieser Befehl nimmt ubyte-Dateien.
```

3.1.2.2 Output

Nach jedem Datensatz wird angezeigt, wie viele der zum Testen bereitgestellten Daten bereits abgearbeitet wurden. Zudem wird angezeigt, wie viele der Datensätze richtig, genauer gesagt falsch erkannt wurden (Dies wird evtl. nicht implementiert.). Am Ende wird aber angezeigt, wie akkurat das Netzwerk die Test-Datensätze bewertet hat.

3.1.3 Prediction

3.1.3.1 Input

```
1 $ ki [file]
2 $ #Dieser Befehl nimmt jpg-Dateien & png-Dateien.
```

3.1.3.2 Output

Als Output werden alle Ziffern von 0 bis 9 mit den entsprechenden Prozentwerten zurückgegeben (Dies kommt daher, dass die KI die Wahrscheinlichkeit der Übereinstimmung für jede Ziffer berechnet.). Zudem wird auch angezeigt, welche Ziffer die höchste Übereinstimmung hat, da diese die erkannte Ziffer darstellt.

Eine interessante Erweiterung wäre es, bei diesem Befehl (oder in einer leicht abgewandelten Form) ein kleines Fenster zu öffnen, um dort direkt die Zahl zu zeichnen. Dies hätte den Vorteil, dass Parameter wie Grösse direkt bekannt und kontrolliert werden können.

3.1.4 Help

3.1.4.1 Input

```
1 $ ki --help
2 $ #Dieser Befehl liefert eine Liste an Befehlen.
```

3.1.4.2 Output

Dieser Befehl liefert eine Liste an Befehlen, die in der Konsole verwendet werden können.

3.2 Datenbank

Die Applikation benötigt keine Datenbank. Zur Speicherung der Test- und Trainings-Datensätze werden sogenannte ubyte-Dateien eingesetzt. Diese enden auf der Dateierweiterung „.ubyte“.

3.2.1 UByte

Die Bilddateien im MNIST-Datensatz werden im sogenannten „ubyte“ (unsigned byte) Format gespeichert. Dies bedeutet, dass die Pixelwerte der Bilder als Bytes gespeichert werden. Jedes Bild im Datensatz wird als Reihe von Bytes dargestellt, wobei jedes Byte den Graustufenwert eines Pixels repräsentiert. Hinzu kommt, dass auch die dazugehörigen Labels in einer „ubyte“-Datei gespeichert werden.

Aufgrund dieser Speicherung müssen die Datensets vor deren Verwendung eingelesen und interpretiert werden.

Für diese Interpretation muss die standardisierte Höhe und Weite für jedes auszulesende Bild multipliziert werden. Danach müssen so viele Bytes aus der entsprechenden „ubyte“-Datei ausgelesen und als vorzeichenlose ganze Zahlen gespeichert werden, denn jede dieser Zahlen stellt einen Pixel dar.

Für eine genaue Anleitung für das Einlesen und Interpretieren in Python ist die Webseite AndroidKT zu kontaktieren.¹

¹vgl. LeCun & Cortes & Burges, 2002; AndroidKT, 2023

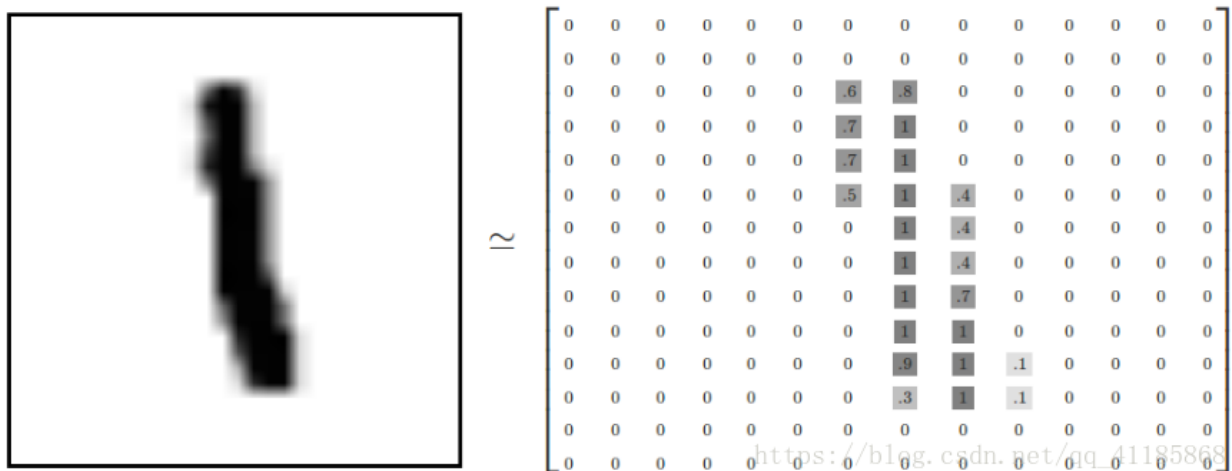


Abbildung 3.1: Darstellung einer Ziffer und deren Code von blog.csdn.net

3.3 Code

3.3.1 Klassendiagramm

Es gibt vier Klassen im Projekt CKi. Dabei handelt es sich um auf sich selbst aufbauende Strukturen. Als Code-Einstieg fungiert die Main-Klasse, welche die Main-Funktion bereitstellt. In dieser wird die Klasse „Network“ aufgesetzt. Es werden mehrere Hidden-Layer erstellt, die Datensätze der MNIST-Datenkollektion eingelesen und die Funktionen „Train“, „Verify“ und „Predict“ sowie „Load Network“ ausgeführt.

In der Network-Klasse werden mehrere Layer initialisiert und verwendet. Diese Layer sind in der Layer-Klasse definiert und besitzen wiederum multiple Neuronen.

Die Neuronen sind die Knotenpunkte, welche mit der mathematischen Sigmoid-Funktion einen Output-Wert berechnen. Zusätzlich ermitteln diese auch die Fehler-Abweichung bei der Back-propagation, um die einzelnen Werte anzupassen.

Um Sigmoid und die Ableitung von Sigmoid zu berechnen, gibt es noch eine Utility-Klasse „Util“. Diese besteht nur aus statischen Funktionen, welche an anderen Orten benötigt werden.

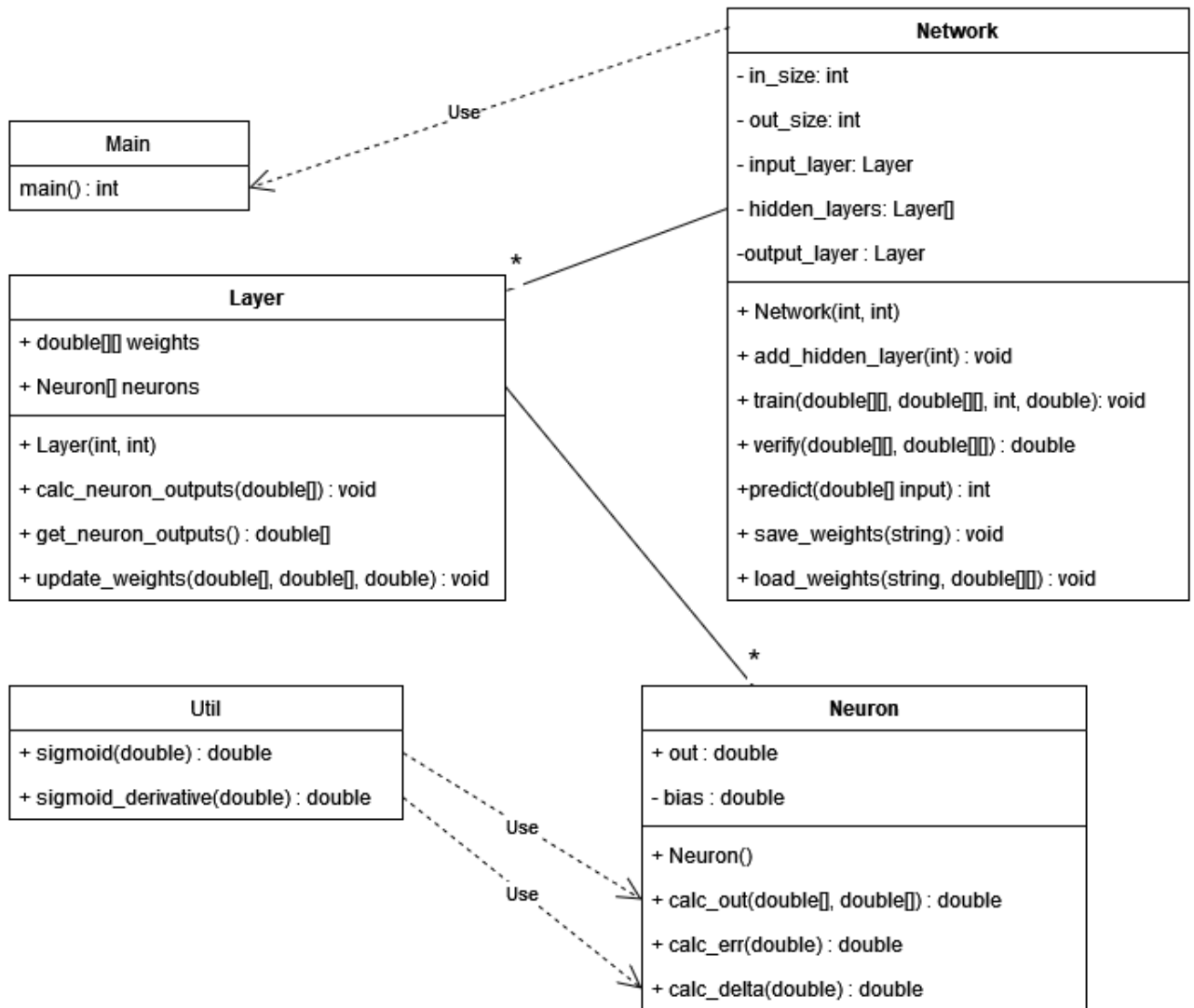


Abbildung 3.2: Das Klassendiagramm, Grundaufbau der Applikation (überholt)

3.3.2 Trainingsdaten

Die Trainingsdaten sind das MNIST-Datenset mit den handschriftlichen Zahlen. (<http://yann.lecun.com/exdb/mnist/>)

3.3.3 Tests

Im Projekt CKi gibt es drei Arten von Tests. Es gibt das simple Ausprobieren. Da es nicht allzu viele Nutzerschnittstellen gibt, kann man diese ausführen und begutachten, ob diese mit der Beschreibung übereinstimmen. Bei einer dieser Schnittstellen wird die KI getestet. Dies geschieht, indem die KI ihr unbekannten Datensätze zu sehen bekommt und das Endergebnis mit einem vordefinierten Ergebnis abgeglichen wird. Egal, ob das Ergebnis korrekt oder inkorrekt erkannt worden ist, wird es statistisch aufgenommen. Am Ende wird dem Nutzer eine Prozentzahl der korrekten Erkenntnisse präsentiert. Diesbezüglich ist dies kein Test, in dem die Applikation versagen könnte, es ist eine reine Leistungsüberprüfung, ob mehr Training vonnöten ist. Zusätzlich zu diesen zwei Testmöglichkeiten gibt es noch die Unit-Tests. Diese werden dazu genützt, um einzelne Funktionen und Klassen noch vor deren Verwendung zu überprüfen.

Zusätzlich zu diesen drei Tests wäre es überaus interessant zu begutachten und zu vergleichen, wie hoch die Leistungsunterschiede zwischen CKi und einem in Python mit TensorFlow gebauten Programm mit vergleichbarer Grundstruktur sind. Diesbezüglich wird eventuell neben CKi auch ein kleines vergleichbares Modell in Python realisiert, um diese Daten zu sammeln.

4 Realisation

4.1 Allgemein

Die Realisation des Projekts CKI zeichnet sich durch die Implementierung eines modularen neuronalen Netzwerks aus, das für Aufgaben wie die Erkennung handschriftlicher Ziffern konzipiert wurde. Zum Einsatz kamen dabei Standard-C++-Technologien sowie eine CMake-basierte Projektstruktur für das Build-Management. Die Kernstruktur besteht aus Klassen für Neuron, Layer, und Network, die die Basis des Netzwerks bilden, ergänzt durch eine Util-Klasse für Hilfsfunktionen, wie Aktivierungsfunktionen und Datenverarbeitung. Zusätzlich gibt es im Ordner „Test“ Unit-Tests mit den entsprechenden Attrappen für die MNIST-Datensätze. Die Unit-Tests sind mit Google-Tests implementiert.

4.2 Abhängigkeiten

Das Projekt CKI kommt mit nur wenigen Abhängigkeiten aus. Die, die es dennoch gibt, sind umso wichtiger für das erfolgreiche Ausführen der Applikation.

Abhängigkeiten:

1. **googletest:**

Version (Tag): v1.14.0

Git: <https://github.com/google/googletest>

Beschreibung: GoogleTest ist ein C++-Framework für Unit-Tests, das Assertions für die Überprüfung von Code und Funktionen für die Organisation und Ausführung von Tests bietet.

2. **nlohmann/json:**

Version (Tag): v3.11.3

Git: <https://github.com/nlohmann/json>

Beschreibung: Nlohmann/json ist eine moderne, header-only C++ Bibliothek für die Verarbeitung von JSON-Daten, die einfache Integration und intuitive Nutzung bietet.

3. **wichtounet/mnist:** (nicht mehr in Verwendung)

Version (Commit): 3b65c35

Git: <https://github.com/wichtounet/mnist>

Beschreibung: Wichtounet/mnist ist ein einfacher C++-Reader für den MNIST-Datensatz, der es ermöglicht, Trainings- und Testbilder sowie Labels zu lesen und zu verwenden.

4. **nothings/stb:**

Version (Commit): 5736b15

Git: <https://github.com/nothings/stb>

Beschreibung: STB ist eine Sammlung plattformübergreifender Header-Dateien für C, die umfassende Funktionen für Grafik, Audio und Textverarbeitung ohne externe Abhängigkeiten bereitstellt.

4.3 Architektur

4.3.1 Klassen

4.3.1.1 Klassendiagramm

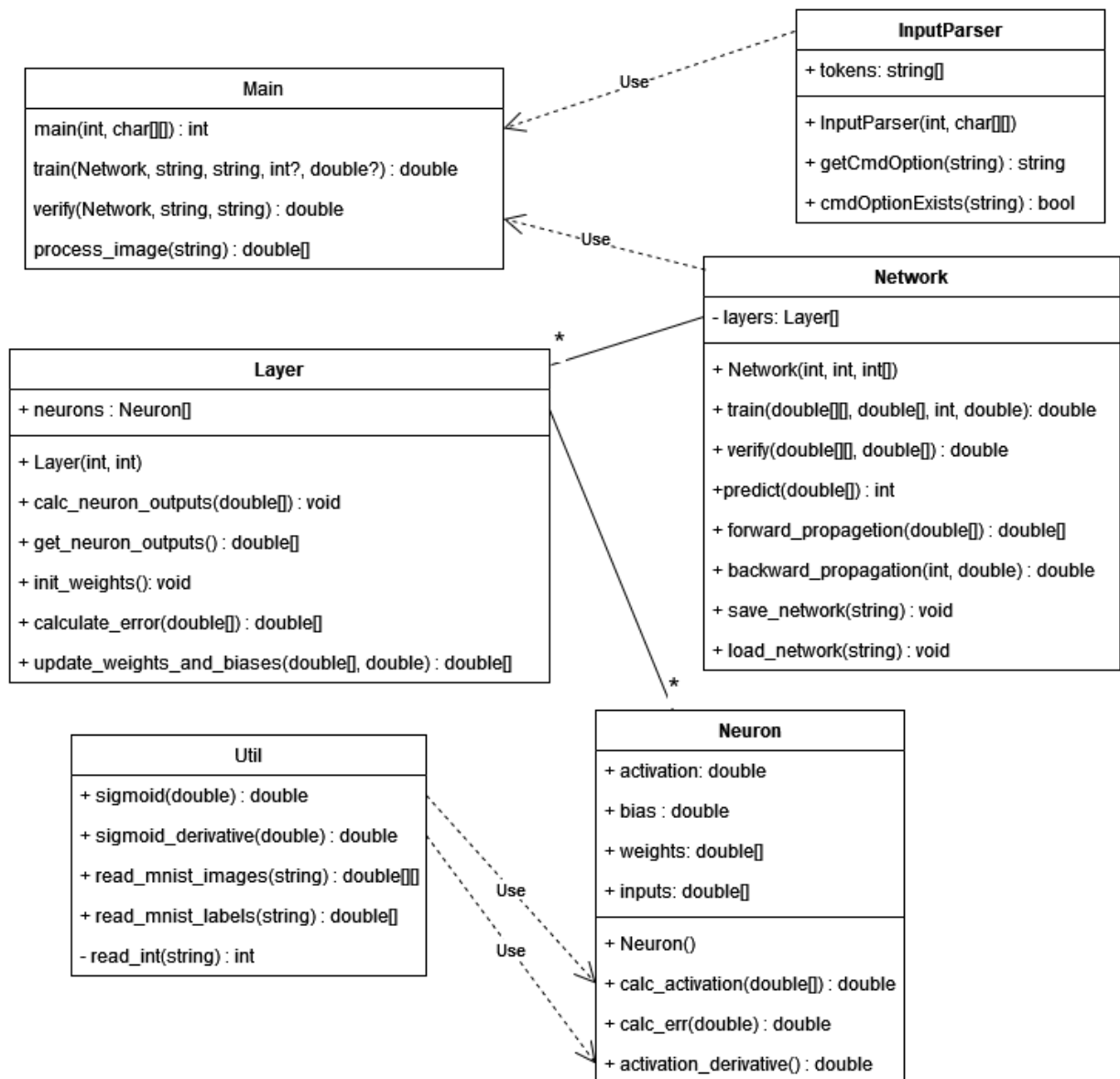


Abbildung 4.1: Klassendiagramm nach der Realisation

Das ist das Klassendiagramm, wie die Applikation des Projektes CKi nach der Implementation aussieht. Somit sind die nötigen Anpassungen, die seit der Designphase zu treffen waren, ersichtlich.

Wie zu sehen ist, gibt es etliche markante Änderungen in den Klassen Netzwerk und Layer. Dies kommt daher, dass in der Designphase aufgrund fehlendem oder unzureichendem Wissen die Struktur angepasst werden musste. Wie diese Klassen funktionieren sowie die wichtigsten Funktionen sind in den Abschnitten 4.3.1, bis 4.4.3 zu finden.

Interessant ist die Klasse InputParser, welche beim ursprünglichen Entwurf nicht anzutreffen war. Diese Klasse ist zuständig, einzelne Attribute aus dem Aufruf der CKI-Main-Funktion auszulesen und eine einfache Verwendung dieser zu gewährleisten. Dabei ist diese Klasse nicht selbst implementiert, sondern stammt von <https://stackoverflow.com/questions/865668/parsing-command-line-arguments-in-c>. Dementsprechend wird im weiteren Verlauf nicht weiter auf diese Funktion eingegangen.

4.3.1.2 Netzwerk

Die Network-Klasse repräsentiert das neuronale Netzwerk. Es unterstützt die Initialisierung des Netzwerks mit einer bestimmten Anzahl von Eingabe- und Ausgabeneuronen sowie eine variable Anzahl von versteckten Schichten und deren Grössen. Die Klasse bietet Funktionen für das Training des Netzwerks mit gegebenen Eingaben und Zielwerten, die Überprüfung der Netzwerkleistung und Vorhersagen für neue Eingaben. Dies ist möglich wegen der Durchführung von Vorwärts- und Rückwärtspropagation, die wichtigsten Algorithmen im Bereich des maschinellen Lernens (Mehr dazu unter Algorithmen 4.4.2).

Zusätzlich existieren noch zwei Funktionen in der Klasse, welche zuständig sind für die Speicherung und das Laden des Netzwerks in bzw. aus einer Datei.

4.3.1.3 Layer

Die Layer-Klasse definiert die Layer im neuronalen Netzwerk, bestehend aus mehreren Neuronen. Sie bietet Funktionen zum Initialisieren und Setzen von Gewichten, Berechnen der Ausgaben der Neuronen basierend auf Eingaben, Berechnen der Fehler im Vergleich zu den Zielwerten und der Aktualisierung von Gewichten und Biases auf Basis von Fehlern und Lern-

rate. Jedes Layer-Objekt enthält eine Liste von Neuron-Objekten, die die Neuronen in diesem Layer repräsentieren.

4.3.1.4 Neuron

Die Neuron-Klasse stellt ein einzelnes Neuron dar, inklusive seiner Gewichte, Bias, Eingaben, Aktivierungsfunktion und Summe der gewichteten Eingaben. Sie bietet Funktionen zur Berechnung der Aktivierung basierend auf den Eingaben, der Ableitung der Aktivierungsfunktion, der Fehlerberechnung im Vergleich zu einem Zielwert, sowie zur Speicherung der Gewichte und des Biases.

Zur Berechnung der Aktivierung und deren Ableitung werden die statischen Utility-Funktionen Sigmoid und die Ableitung von Sigmoid aus der Utility-Klasse genutzt (Siehe auch 4.4.3.7).

4.3.1.5 Utility

Die Util-Klasse bietet statische Hilfsfunktionen für das neuronale Netzwerk, darunter die Sigmoid-Aktivierungsfunktion und ihre Ableitung sowie Funktionen zum Lesen von MNIST-Bilddaten und Labels aus Dateien. Diese Hilfsfunktionen sind essenziell für die Vorverarbeitung von Eingabedaten und die Implementierung der Lernmechanismen im Netzwerk.

4.3.2 Ordnerstruktur

Im Kern des Projekts stehen die Hauptdateien, welche die essenziellen Klassen wie Neuron, Layer, Network und Util enthalten. Diese sind grundlegend für die Funktionalität des neuronalen Netzwerks. Die CMakeLists.txt unterstützt das Build-Management, vereinfacht die Kompilierung und Konfiguration. Der .gitignore sorgt dafür, dass unnötige Dateien und Ordner nicht in die Versionskontrolle einfließen. Ein speziell dafür vorgesehener Test-Ordner beinhaltet Tests zur Überprüfung der Funktionalität, was die Zuverlässigkeit des Systems sicherstellt. Zusätzlich runden Dummy-Files und ein eigener Ordner für die MNIST-Datasets das Projekt ab, indem sie das Vorhandensein der Datensets für Training und Tests des Netzwerks garantieren.

Das CKI-Projekt ist strukturiert in:

- **Hauptdateien:** (kein Ordner, Rootverzeichnis) Enthalten die Klassen Neuron, Layer, Network, und Util für die Kernlogik des neuronalen Netzwerks.
- **CMakeLists.txt:** (Im Rootverzeichnis) Für das Build-Management erleichtert das Kompilieren und die Konfiguration des Projekts.
- **.gitignore:** (Im Rootverzeichnis) Definiert Dateien und Ordner, die von Git-Versionierung ausgeschlossen sind, wie Build-Artefakte und IDE-spezifische Dateien.
- **Test-Ordner:** Beinhaltet Testfälle zur Überprüfung der Funktionalität einzelner Komponenten und des Gesamtsystems.
 - **Dummy-Files:** Zusätzliche Dateien für Testzwecke.
 - **Unittests:** Tests zur Sicherstellung korrekter Funktionalität, Isolation von Fehlern.
- **MNIST-Datasets Ordner:** Speichert die Datensätze für das Training und Testen des Netzwerks, insbesondere für die Erkennung handschriftlicher Ziffern.

4.4 Code

4.4.1 Konzept

Das Projekt CKi nutzt grundlegende Konzepte des maschinellen Lernens wie Vorwärts- und Rückwärts-Propagierung, Aktivierungsfunktionen (hier Sigmoid) und die Anpassung von Gewichten während des Trainingsprozesses. Die modulare Struktur, unterteilt in Schichten, Neuronen und Hilfsfunktionen, sowie die reine objektorientierte Implementation in C++ ohne externe Bibliotheken, abgesehen von Bibliotheken für Tests und Datenhandling, zeichnen die Architektur der Applikation aus.

Um mehr über die hier erwähnten Funktionen zu erfahren, ist der Abschnitt 4.4.2 zu konsultieren.

4.4.2 Algorithmen

4.4.2.1 Forward-Propagation

Die Forward-Propagation ist ein grundlegender Prozess in neuronalen Netzwerken, der es ermöglicht, Vorhersagen auf Basis von Eingabedaten zu treffen. Dabei werden die Eingabedaten durch das Netzwerk von der Eingabeschicht über eine oder mehrere versteckte Schichten bis zur Ausgabeschicht vorwärts geleitet. Jede Schicht besteht aus Neuronen, die über Gewichte mit den Neuronen der vorherigen Schicht verbunden sind. Die Daten werden in jedem Neuron durch eine Summationsfunktion verarbeitet, die die gewichteten Eingaben aufsummiert und einen Bias-Wert hinzufügt. Das Ergebnis dieser Summation wird dann durch eine Aktivierungsfunktion geleitet, um die Ausgabe des Neurons zu bestimmen.

Die Aktivierungsfunktion bestimmt, wie Neuronen ihre Eingaben in Ausgaben umwandeln und ist entscheidend für die Fähigkeit des Netzwerks, komplexe Muster in den Daten zu erkennen. Beliebte Aktivierungsfunktionen sind die Sigmoid-, Tanh- und ReLU-Funktion.¹

Sobald die Eingabedaten durch das Netzwerk propagiert worden sind und die Ausgabeschicht erreicht haben, wird das Ergebnis mit dem tatsächlichen Wert verglichen, um den Fehler der Vorhersage zu bestimmen. Dieser Fehler wird dann in einem separaten Prozess, der als Backpropagation bekannt ist, verwendet, um die Gewichte im Netzwerk anzupassen und die Vorhersagegenauigkeit zu verbessern.²

Mathematische Beschreibung der Forward-Propagation

Dieser Abschnitt ist eine kurze mathematische Repräsentation des oben beschriebenen Algorithmus. Die Informationen für den gesamten Abschnitt stammen von Sven Nüesch (2023) und 3Blue1Brown (2017).

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Die Sigmoid-Funktion ist eine beliebte Aktivierungsfunktion, die die Summation der gewichteten Eingaben und des Bias in einen Wert zwischen 0 und 1 umwandelt. Diese Funktion

¹vgl. Li & Johnson & Yeung, 2017; databasecamp, 2023

²vgl. Nüesch, 2023; 3Blue1Brown, 2017;

ermöglicht es dem Netzwerk, Nichtlinearitäten zu modellieren.

$$b = Bias$$

$$w = Weight$$

Gewichte sind die Parameter, die die Stärke der Verbindung zwischen den Neuronen in aufeinanderfolgenden Schichten steuern, während der Bias eine zusätzliche Anpassung ermöglicht, die zu der gewichteten Summe der Eingaben hinzugefügt wird, um die Aktivierung eines Neurons zu beeinflussen.

$$a = Activation$$

Die Aktivierung ist die Ausgabe eines Neurons nach Anwendung der Aktivierungsfunktion auf die gewichtete Summe der Eingaben plus Bias.

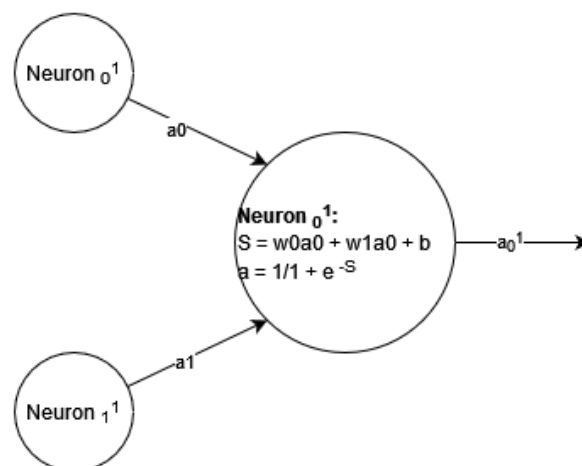


Abbildung 4.2: Abbildung von 3 isolierten Neuronen (Inspiriert von 3Blue1Brown)

$$a_0^{(1)} = \sigma(w_0 a_0^{(0)} + w_1 a_1^{(0)} + b) = \sigma\left(\sum_j^n w_j a_j^{(0)} + b\right)$$

Diese Formel repräsentiert die gewichtete Summe der Eingabesignale $a_j^{(0)}$, wobei jedes Eingabesignal mit einem Gewicht w_j multipliziert wird. Der Index j läuft über alle Eingaben von 0 bis n , wobei n die Anzahl der Eingabeneuronen ist. Zusätzlich wird ein Bias-Term b

hinzugefügt, der es dem Neuron ermöglicht seine Ausgabe unabhängig von der Eingabe anzupassen.

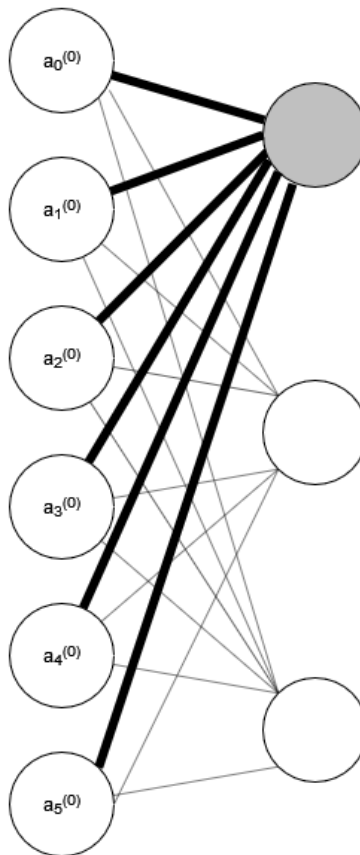


Abbildung 4.3: Abbildung von einigen Neuronen (Inspiriert von 3Blue1Brown)

$$a_0^{(1)} = \sigma(w_{0,0} a_0^{(0)} + w_{0,1} a_1^{(0)} + \dots + w_{0,n} a_n^{(0)} + b_0)$$

Diese Gleichung beschreibt die Aktivierung $a_0^{(1)}$ eines Neurons in der ersten Schicht (ausser der Eingabeschicht) basierend auf den Eingaben $a_i^{(0)}$ (wo i von 0 bis n läuft), den Gewichten $w_{0,i}$, die diese Eingaben mit dem Neuron verbinden, und dem Bias b_0 .

Um von dieser Gleichung eines einzelnen Neurons zu den Gleichungen für einen gesamten Layer mit $k+1$ Neuronen zu gelangen, erweitern wir die vorherige Gleichung für jedes Neuron j im Layer:

$$a_j^{(1)} = \sigma(w_{j,0} a_0^{(0)} + w_{j,1} a_1^{(0)} + \dots + w_{j,n} a_n^{(0)} + b_j)$$

Dabei läuft j von 0 bis k .

Die Gleichungen für alle Neuronen in einem Layer können effizient in Matrixform ausgedrückt werden. Die Gewichte zwischen den Eingaben und den Neuronen im Layer können in einer Gewichtsmatrix W zusammengefasst werden, die Eingaben in einem Vektor $a^{(0)}$, und die Bias-Werte in einem Bias-Vektor b . Die Operation wird dann als Matrix-Vektor-Multiplikation dargestellt:

$$\sigma \left(\begin{pmatrix} w_{0,0} & w_{0,1} & \cdots & w_{0,n} \\ w_{1,0} & w_{1,1} & \cdots & w_{1,n} \\ \vdots & \vdots & \ddots & \vdots \\ w_{k,0} & w_{k,1} & \cdots & w_{k,n} \end{pmatrix} \begin{pmatrix} a_0^{(0)} \\ a_1^{(0)} \\ \vdots \\ a_n^{(0)} \end{pmatrix} + \begin{pmatrix} b_0 \\ b_1 \\ \vdots \\ b_n \end{pmatrix} \right)$$

Hier repräsentiert W die Gewichtsmatrix, $a^{(0)}$ den Eingabevektor (aus der vorherigen Schicht), und b den Bias-Vektor. Die Aktivierungsfunktion σ wird elementweise auf das Ergebnis der Matrix-Vektor-Addition angewendet.

Schliesslich kann die obige Gleichung noch weiter vereinfacht werden, indem die Matrix, der Eingabevektor und der Bias-Vektor einfach durch ihre Symbole dargestellt werden, ohne ihre Elemente explizit aufzuführen. Dies führt zur kompaktesten Form:

$$a^{(1)} = \sigma(W a^{(0)} + b)$$

In dieser Gleichung steht $a^{(1)}$ für den Vektor der Aktivierungen der Neuronen im betrachteten Layer nach Anwendung der Aktivierungsfunktion auf die gewichtete Summe der Eingaben und der Bias-Werte.

4.4.2.2 Back-Propagation

Die Backpropagation, kurz für „backward propagation of errors“, ist ein Schlüsselmechanismus im Training neuronaler Netzwerke. Dieser Algorithmus ermöglicht es, die Gewichte des Netzwerks so anzupassen, dass der Gesamtfehler bei der Vorhersage minimiert wird. Backpropagation wird nach der Forward-Propagation angewendet, nachdem eine Vorhersage durch das Netzwerk gemacht und der Fehler zwischen der Vorhersage und dem tatsächlichen Wert

berechnet wurde.

Der Prozess der Backpropagation besteht aus zwei Hauptphasen: der Berechnung des Gradienten des Fehlers bezüglich aller Gewichte im Netzwerk und der anschliessenden Anpassung dieser Gewichte in die Richtung, die den Fehler minimiert. Der Gradient gibt an, in welche Richtung die Gewichte verändert werden müssen, um den Fehler zu verringern, und die Grösse der Anpassung wird durch die Lernrate bestimmt.

Backpropagation nutzt die Kettenregel der Differenzialrechnung, um die Fehlergradienten für die Gewichte jeder Schicht vom Ausgang zurück zum Eingang effizient zu berechnen. Der berechnete Fehlergradient für jede Gewichtung zeigt, wie eine kleine Änderung in diesem Gewicht den Gesamtfehler beeinflusst. Durch die systematische Anpassung der Gewichte basierend auf diesen Gradienten kann das Netzwerk schrittweise verbessert werden, um genauere Vorhersagen zu liefern.

Insgesamt ermöglicht Backpropagation das effiziente Training tiefer neuronaler Netzwerke, indem es systematisch die Netzwerkgewichte anpasst, um den Fehler zwischen den Vorhersagen des Netzwerks und den tatsächlichen Werten zu minimieren, was zu einer verbesserten Modellleistung führt.³

Mathematische Beschreibung der Back-Propagation

Dieser Abschnitt ist eine kurze mathematische Repräsentation des oben beschriebenen Algorithmus. Die Informationen für den gesamten Abschnitt stammen von Sven Nüesch (2023) und 3Blue1Brown (2017).

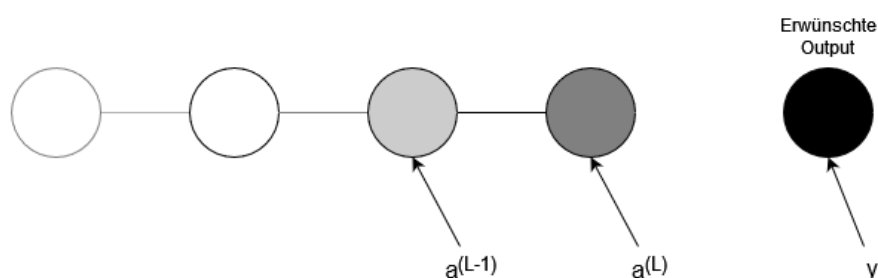


Abbildung 4.4: 4 Neuronen für die Back-Propagation (Bild von 3Blue1Brown)

³vgl. Nüesch, 2023; 3Blue1Brown, 2017; Karpathy, 2016; Roy, 2022; Li & Johnson & Yeung, 2017;

$$\sigma'(x) = \sigma(x) (1 - \sigma(x))$$

Die Funktion ist die Ableitung der Sigmoid-Funktion, die in der neuronalen Netzwerkoptimierung für die Gewichtsaktualisierung mittels Backpropagation verwendet wird.

$$z^{(L)} = w^{(L)} a^{(L-1)} + b^{(L)}$$

$$a^{(L)} = \sigma(z^{(L)})$$

$y = \text{erwartete Output}$

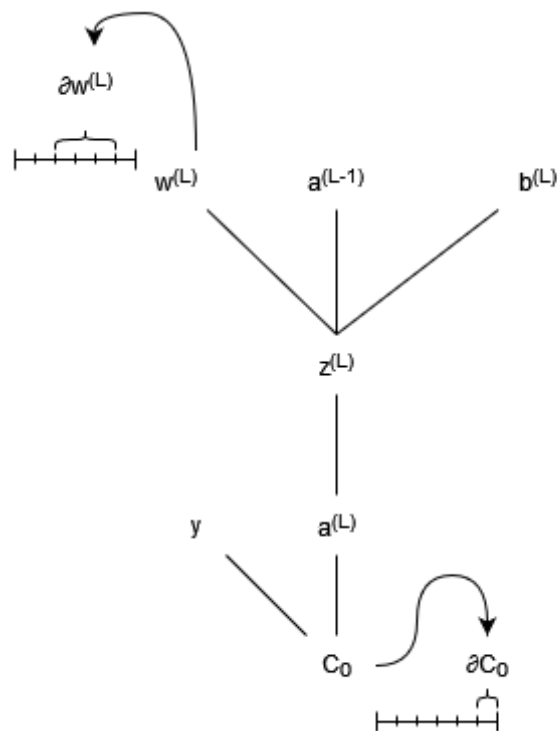


Abbildung 4.5: Kostenbaum in der Back-Propagation (Bild von 3Blue1Brown)

$$Cost_0(\dots) = C_0 = (a^{(L)} - y)^2$$

Die Formel beschreibt die Kostenfunktion für ein einzelnes Trainingsbeispiel in einem neuronalen Netzwerk. Diese Kostenfunktion, oft als quadratischer Fehler oder Mean Squared Error

(MSE) bezeichnet, misst den quadratischen Unterschied zwischen der vorhergesagten Ausgabe und dem tatsächlichen Wert, um die Leistung des Modells zu bewerten und während des Trainingsprozesses zu optimieren.

$$\frac{\partial C_0}{\partial w^{(L)}} = \text{Empfindlichkeit der Cost - Funktion in Bezug auf das Gewicht}$$

$$\frac{\partial C_0}{\partial w^{(L)}} = \frac{\partial z^{(L)}}{\partial w^{(L)}} \frac{\partial a^{(L)}}{\partial z^{(L)}} \frac{\partial C_0}{\partial a^{(L)}}$$

Diese Formel beschreibt die Berechnung des Gradienten der Kostenfunktion C_0 bezüglich des Gewichts $w^{(L)}$ in der letzten Schicht LL eines neuronalen Netzes. Sie nutzt die Kettenregel der Differentialrechnung, um den Einfluss von Änderungen im Gewicht auf die Kosten zu ermitteln, indem sie das Produkt aus der Ableitung des Netzeingangs $z^{(L)}$ nach dem Gewicht, der Ableitung der Aktivierungsfunktion $a^{(L)}$ nach dem Netzeingang und der Ableitung der Kostenfunktion nach der Aktivierung berechnet.

$$\frac{\partial C_0}{\partial a^{(L)}} = 2(a^{(L)} - y)$$

$$\frac{\partial a^{(L)}}{\partial z^{(L)}} = \sigma'(z^{(L)})$$

$$\frac{\partial z^{(L)}}{\partial w^{(L)}} = a^{(L-1)}$$

$$\frac{\partial C_0}{\partial w^{(L)}} = a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

Dies ist eine spezifischere Formulierung, die direkt die Komponenten aus der oberen Gleichung zusammenfasst und annimmt. Diese Gleichungen verbinden also die Änderung des Fehlers in Bezug auf ein Gewicht direkt mit der Aktivierung der vorherigen Schicht.

$$\frac{\partial z^{(L)}}{\partial b^{(L)}} = 1$$

$$\frac{\partial C_0}{\partial b^{(L)}} = 1 \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

Diese Funktion berechnet wie die Funktion $\frac{\partial C_0}{\partial w^{(L)}} = a^{(L-1)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$ den Gradienten

der Kostenfunktion C_0 , jedoch nicht in Bezug auf die Gewichtung $w^{(L)}$, sondern auf den Bias $b^{(L)}$.

$$\frac{\partial z^{(L)}}{\partial a^{(L-1)}} = w^{(L)}$$

$$\frac{\partial C_0}{\partial a^{(L-1)}} = w^{(L)} \sigma'(z^{(L)}) 2(a^{(L)} - y)$$

Diese Formel zeigt, im Gegensatz zu den beiden anderen, wie der Gradient der Kostenfunktion C_0 bezüglich der Aktivierungen $a^{(L-1)}$ des vorherigen Layers berechnet wird. So kann diese Funktion, gleich wie die für die Gewichtung und dem Bias, dazu verwendet werden, zu verstehen, wie Änderungen in den Aktivierungen des vorherigen Layers $L - 1$ den Gesamtfehler beeinflussen können.

$$\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$$

$$\nabla C = \begin{pmatrix} \frac{\partial C}{\partial w^{(1)}} \\ \frac{\partial C}{\partial b^{(1)}} \\ \vdots \\ \frac{\partial C}{\partial w^{(L)}} \\ \frac{\partial C}{\partial b^{(L)}} \end{pmatrix}$$

Die Gleichung $\frac{\partial C}{\partial w^{(L)}} = \frac{1}{n} \sum_{k=0}^{n-1} \frac{\partial C_k}{\partial w^{(L)}}$ beschreibt den Gradienten der Kostenfunktion C in Bezug auf die Gewichte $w^{(L)}$ in der L -ten Schicht des Netzwerks. Der Gradient wird als Durchschnitt der Gradienten aus n Trainingsbeispielen berechnet.

Die 2 Notation (∇C) stellt den Gradienten der Kostenfunktion C als Vektor dar, der alle partiellen Ableitungen der Kostenfunktion in Bezug auf alle Gewichte $w^{(L)}$ und alle Bias-Werte $b^{(l)}$ für jede Schicht L des Netzwerks enthält. Dieser Gradientenvektor ist wesentlich für die Anwendung des Gradientenabstiegsverfahrens, da er die Richtung angibt, in der die Gewichte und Bias-Werte angepasst werden müssen, um die Kostenfunktion zu minimieren.

Mit mehreren Outputs:

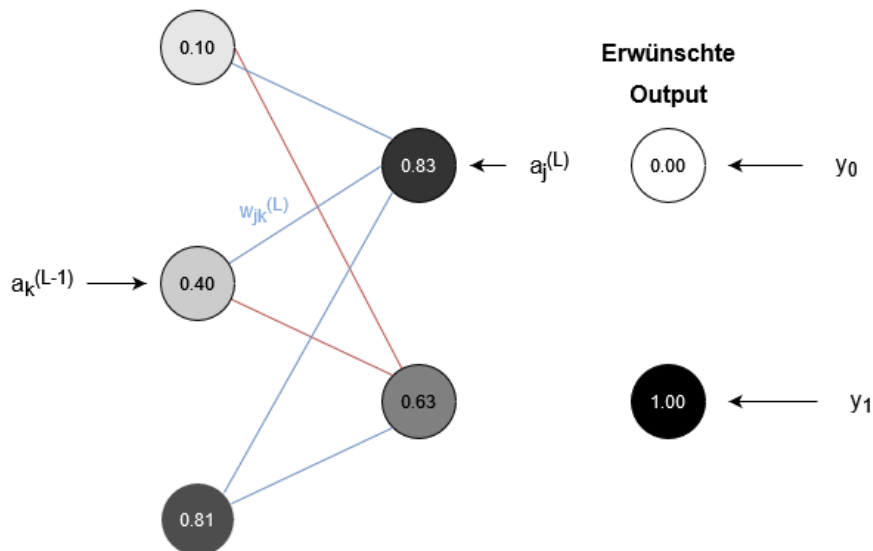


Abbildung 4.6: 5 Neuronen in einem Netzwerk (Bild von 3Blue1Brown)

Dies sind im Grundsatz die gleichen Funktionen wie oben, es wurde ihnen aber eine weitere Komplexitätsstufe hinzugefügt, indem nun das Netzwerk nicht mehr aus nur einem Neuron pro Layer besteht.

$$C_0 = \sum_{j=0}^{n_L-1} (a_j^{(L)} - y_j)^2$$

$$z_j^{(L)} = \dots + w_{jk}^{(L)} a_k^{(L-1)} + \dots + b_j^{(L)}$$

$$a_j^{(L)} = \sigma(z_j^{(L)})$$

$$\frac{\partial C_0}{\partial w_{jk}^{(L)}} = \frac{\partial z_j^{(L)}}{\partial w_{jk}^{(L)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}$$

$$\frac{\partial C_0}{\partial a_k^{(L-1)}} = \sum_{j=0}^{n_L-1} \frac{\partial z_j^{(L)}}{\partial a_k^{(L-1)}} \frac{\partial a_j^{(L)}}{\partial z_j^{(L)}} \frac{\partial C_0}{\partial a_j^{(L)}}$$

$$\frac{\partial C}{\partial w_{jk}^{(l)}} = a_k^{(l-1)} \sigma'(z_j^{(l)}) \frac{\partial C}{\partial a_j^{(l)}}$$

$$\frac{\partial C}{\partial a_j^{(l)}} = \sum_{k=0}^{n_{l+1}-1} w_{jk}^{(l+1)} \sigma'(z_j^{(l+1)}) \frac{\partial C}{\partial a_j^{(l+1)}}$$

or

$$2(a_j^{(L)} - y_j)$$

4.4.3 Schlüsselpassagen & Snippets

Die wichtigsten Funktionen in einem neuronalen Netzwerk sind die beiden Propagations, vorwärts als auch rückwärts. Wobei diese das Grundgerüst für die drei Hauptfunktionen, Train, Verify & Predict bilden.

4.4.3.1 Train

```

1 double Network::train(std::vector<std::vector<double>> &inputs, std::vector<↵
  ↵ double> &labels, int epochs, double learning_rate){
2     double total_error = 0;
3     for(int epoch = 0; epoch<epochs; epoch++){
4         for(std::size_t i = 0; i < inputs.size(); i++){
5             std::vector<double> outputs = Network::forward_propagation(inputs[i↵
              ↵ ]);
6             total_error += Network::backward_propagation(labels[i], ↵
              ↵ learning_rate);
7             total_error /= 2;
8         }
9     }
10    return total_error;
11 }
```

Die Funktion, welche für das Lernen des neuronalen Netzwerks verantwortlich ist, ist simpel aufgebaut. Sie erhält primär zwei Listen: eine Liste aus Bildern und eine Liste aus Beschriftungen für die Bilder. Danach nutzt Train die Funktionen der Forward- und Back-Propagation, um zuerst das Bild durch das Netzwerk bewerten zu lassen und dann diese Bewertung mithilfe der Beschriftungen zu korrigieren. Dies tut sie für den gesamten Datensatz oder durch eine zukünftige Erweiterung in einzelnen Blöcken (für höhere Effektivität) und der Anzahl Epochen

entsprechend.

4.4.3.2 Verify

```

1 double Network::verify(const std::vector<std::vector<double>> &inputs, const std::vector<double> &labels){
2     int correct = 0;
3     for (std::size_t i = 0; i < inputs.size(); ++i){
4         if (Network::predict(std::vector<double>(inputs[i])) == static_cast<int>(labels[i])){
5             correct++;
6         }
7     }
8     return static_cast<double>(correct) / static_cast<double>(inputs.size());
9 }

```

Die Funktion ist eigentlich gleich aufgebaut wie „Train“, jedoch mit dem markanten Unterschied, dass keine Back-Propagation zum Einsatz kommt. So werden die Bewertungen des neuronalen Netzwerks nur als richtig oder falsch bewertet und nicht korrigiert.

4.4.3.3 Predict

```

1 int Network::predict(const std::vector<double> &input){
2     std::vector<double> outputs = Network::forward_propagation(input);
3     return static_cast<int>(std::distance(outputs.begin(), std::max_element(outputs.begin(), outputs.end())));
4 }

```

Diese Funktion erhält nur ein einziges Bild als Attribut. Dieses Bild wird durch die Forward-Propagation bewertet. Die Bewertung wiederum besteht nur aus einer Liste von zehn Prozentzahlen (für die Ziffern von 0 bis 9) und damit diese besser für einen Menschen lesbar wird, muss die höchste Prozentzahl aus der Liste gesucht werden. Deren Index wird dann als errechnete Zahl zurückgegeben.

(Diese Funktion ist ein Zusammenzug aus der Predic-Funktion und Logic aus der main.cpp-Datei.)

4.4.3.4 Forward-Propagation

```
1 std::vector<double> Network::forward_propagation(const std::vector<double>& ↵  
    ↵ input){  
2     layers[0].calc_neuron_outputs(input);  
3  
4     for (int i = 1; i < layers.size(); ++i){  
5         std::vector<double> inputs = layers[i - 1].get_neuron_outputs();  
6         layers[i].calc_neuron_outputs(inputs);  
7     }  
8  
9     return layers[layers.size() - 1].get_neuron_outputs();  
10 }
```

In dieser Funktion wird eine Liste von grossen Fließkommazahlen übernommen. Diese sind die einzelnen Pixel, in dem zu bearbeiteten Bild, welche nur einen Helligkeitswert beinhalten. Diese Pixel werden an den ersten Layer im Netzwerk übergeben. Dieser kalkuliert durch die Sigmoid-Funktion neue Ausgabewerte. Die Ausgabewerte werden mit der For-Schleife durch alle Schichten des Netzwerks weitergereicht, bis diese am Ende ausgelesen und als Rückgabewert dieser Funktion genutzt werden. Die Ausgabewerte haben zwar das gleiche Datenformat, unterscheiden sich jedoch anhand ihrer Grösse (Länge der Liste) und den tatsächlichen Werten.

4.4.3.5 Back-Propagation

```
1 double Network::backward_propagation(const int& target, double learning_rate){  
2     std::vector<double> inputs (10, 0);  
3     inputs[target] = 1;  
4     std::vector<double> error = layers[layers.size() - 1].calculate_error(inputs);  
5  
6     double total_error = 0;  
7     for(double& e : error) { //remove for better performance  
8         total_error += e;  
9     }  
10 }
```

```

11     total_error /= error.size();
12
13     //std::cout << total_error << std::endl;
14
15     for (int i = layers.size() - 1; i >= 0; --i){
16         error = layers[i].update_weights_and_biases(error, learning_rate);
17     }
18     return total_error;
19 }

```

Im Gegensatz zur Forward-Propagation wird in der Back-Propagation ein Wert nicht von der ersten Schicht zur letzten übergeben, sondern Rückwärts von dem letzten Layer bis zur Eingabe zurück.

Dabei muss zuerst die Abweichung von der erhaltenen Ausgabe (aus der Forward-Propagation) mit der erwarteten Ausgabe abgeglichen und der Fehler berechnet werden. Nun wird ähnlich zur Forward-Propagation dieser Fehler durch die einzelnen Schichten geschickt und dort verwendet, um die Gewichtungen und Biases für jedes Neuron anzupassen. Dies passiert in einer eigenen Funktion im Layer.

4.4.3.6 Anpassung von Gewichtungen und Biases

```

1  std::vector<double> Layer::update_weights_and_biases(const std::vector<double>& ↵
    ↵ error, double learning_rate){
2      std::vector<double> prev_layer_error(Layer::neurons[0].weights.size(), 0.0);
3
4      for (size_t i = 0; i < Layer::neurons.size(); ++i){
5          Neuron &neuron = Layer::neurons[i];
6          const double neuron_error = error[i];
7          const double activation_derivative = neuron.activation_derivative();
8
9          for (size_t j = 0; j < neuron.weights.size(); ++j){
10             prev_layer_error[j] += neuron.weights[j] * neuron_error;
11
12             const double delta_weight = neuron_error * activation_derivative * ↵
                ↵ neuron.inputs[j];

```

```
13     neuron.weights[j] -= learning_rate * delta_weight;  
14 }  
15     neuron.bias += learning_rate * neuron_error * activation_derivative;  
16 }  
17     return prev_layer_error;  
18 }
```

Dies ist die Funktion, die in der Back-Propagation aufgerufen wird.

Dabei wird für jedes Neuron für jede Gewichtung eine eigene Anpassung gesucht. Um dies zu erreichen, werden mehrere Werte benötigt: Der Fehler des Neurons, der Input des Neurons, die Ableitung der Aktivierung und die Learningrate. Nun können die Änderungen bei der Gewichtung und dem Bias berechnet werden. Die Abweichung der Gewichtung ist ein Produkt des Inputs in diesem Neuron, dem Fehler, der Ableitung der Sigmoid-Funktion und der Learningrate. Die Korrektur für das Bias ist fast identisch zu der Abweichung der Gewichtung. Es ist ebenfalls ein Produkt der gleichen Faktoren, abgesehen von der Learningrate und dem Input.

Zusätzlich kann der Fehler für den nächsten Layer berechnet werden. Dieser wird benötigt, um die Backpropagation auf den nächsten Layer auszudehnen und dort dieselben Schritte auszuführen. Dabei setzt sich der Fehler für den nächst oberen Layer zusammen, aus dem Fehler und dem entsprechenden Gewicht für das Neuron. Dabei wird der Fehler neu gewichtet und zusammen mit allen anderen neu-gewichteten Fehlern für dieses eine Neuron auf der nächsten Ebene aufsummiert.

4.4.3.7 Sigmoid & Ableitung

```
1 double Util::sigmoid(double x){  
2     return 1.0 / (1.0 + exp(-x));  
3 }  
4  
5 double Util::sigmoid_derivative(double x){  
6     double s = sigmoid(x);  
7     return s * (1.0 - s);  
8 }
```

Das sind die mathematischen Implementierungen der Sigmoid-Funktion und deren mathematischen Abweichungen. Diese Funktionen sind zuständig für die einzelnen Aktivierungen der Neuronen und verschieben eine Zahl in den Wertebereich zwischen 0 und 1.

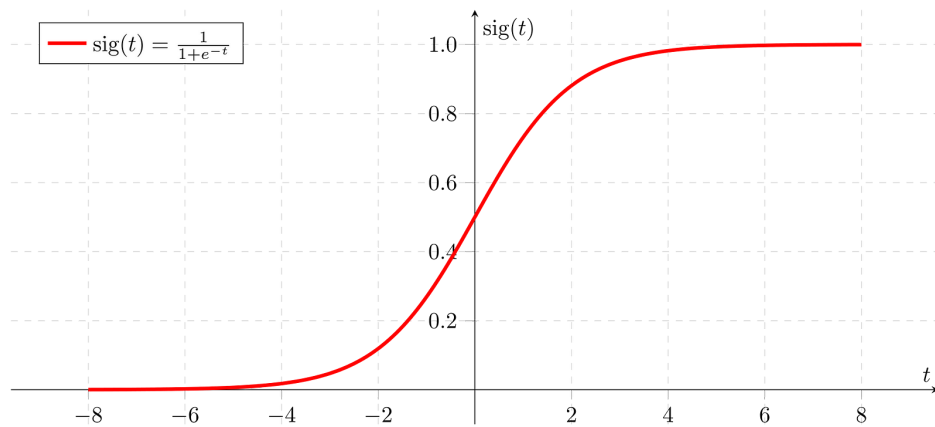


Abbildung 4.7: Sigmoid Funktion

4.4.3.8 MNIST-Reader

```

1 std::vector<std::vector<double>> Util::read_mnist_images(const std::string &↵
  ↵ filename){
2     std::ifstream file(filename, std::ios::binary);
3     if (file.is_open()){
4         int magic_number = read_int(file);
5         int number_of_images = read_int(file);
6         int number_of_rows = read_int(file);
7         int number_of_columns = read_int(file);
8
9         std::vector<std::vector<double>> images(number_of_images, std::vector<↵
  ↵ double>(number_of_rows * number_of_columns));
10
11         for (int i = 0; i < number_of_images; ++i){
12             for (int r = 0; r < number_of_rows * number_of_columns; ++r){
13                 unsigned char temp = 0;
14                 file.read(reinterpret_cast<char*>(&temp), sizeof(temp));
15                 images[i][r] = (double)temp / 255.0; // Normalizing pixel values↵
  ↵ to [0, 1]
16             }

```

```

17     }
18     return images;
19 } else {
20     throw std::runtime_error("Cannot open file : " + filename);
21 }
22 }

```

Um die Daten aus den von <http://yann.lecun.com/exdb/mnist/> stammenden UByte-Dateien auszulesen, müssen gewisse Byte-Werte aus den Dateien ausgelesen werden. Die hier aufgeführte Funktion dient dazu Bilder aus diesen Dateien auszulesen und ist beispielhaft auch für die Labels.

Die Implementierung ist nach dem Muster von <http://yann.lecun.com/exdb/mnist/> implementiert. So sind in den ersten vier Byte gespeichert, wie viele Bilder/Beschriftungen im Datensatz zu finden sind. Um dies auszulesen, ist in einer eigenen Funktion implementiert (`read_int`). Nach den ersten vier Byte müssen die restlichen Bytes in eine Liste eingetragen, in eine ganze Zahl interpretiert (zwischen 0 und 255) und auf den Wert zwischen 0 und 1 gebracht werden. Man spricht auch von der Normalisierung der Pixel-Werte.

4.4.3.9 Ersten 4 Bytes

Dies ist die oben angesprochene „`read_int`“-Funktion.

```

1 int Util::read_int(std::ifstream &file){
2     unsigned char bytes[4];
3     file.read(reinterpret_cast<char*>(bytes), sizeof(bytes));
4     return (int)((bytes[0] << 24) | (bytes[1] << 16) | (bytes[2] << 8) | bytes[3]);
5 }

```

Die Funktion wurde mit der Hilfe von Stack Overflow unter dem Link <https://stackoverflow.com/questions/604431/c-reading-unsigned-char-from-file-stream> erstellt.

Zuerst wird ein Array in der Grösse von vier Bytes erstellt und vier Bytes aus einer Datei eingelesen. Danach werden die vier Bytes zu einer einzigen ganzen Zahl zusammengesetzt. Dies geschieht, indem die Bytes um jeweils 8 Bits nach links verschoben werden und somit das erste Byte die acht höchsten Stellen der ganzen Zahl symbolisiert, danach das zweite Byte

und so folgend die restlichen Bytes.

4.4.3.10 Bild Vorverarbeitung

```
1 std::vector<double> process_image(const char* filename) {
2     int width, height, channels;
3     unsigned char* img = stbi_load(filename, &width, &height, &channels, 0);
4     if (img == nullptr) {
5         std::cerr << "Error in loading the image" << std::endl;
6         exit(1);
7     }
8
9     unsigned char* resized_img = new unsigned char[28 * 28 * channels];
10    stbir_resize_uint8(img, width, height, 0, resized_img, 28, 28, 0, channels);
11
12    std::vector<double> image_data(28 * 28);
13    for (int i = 0; i < 28 * 28; ++i) {
14        int j = i * channels;
15        double gray = 0.299 * resized_img[j] + 0.587 * resized_img[j + 1] + ↵
            ↵ 0.114 * resized_img[j + 2];
16        image_data[i] = gray / 255.0; // Normalize to [0, 1]
17    }
18
19    stbi_image_free(img);
20    delete[] resized_img;
21
22    return image_data;
23 }
```

Dies Funktion lädt ein Bild von der angegebenen Datei, skaliert es auf 28x28 Pixel herunter und konvertiert es in Graustufen. Zuerst wird das Bild geladen, wobei Breite, Höhe und Kanäle des Bildes erfasst werden. Falls das Laden fehlschlägt, wird eine Fehlermeldung ausgegeben und das Programm beendet. Anschließend wird das Bild auf die Größe 28x28 Pixel verkleinert, wobei die Anzahl der Kanäle beibehalten wird. Danach wird für jeden Pixel ein Graustufenwert berechnet, indem eine gewichtete Summe seiner RGB-Werte gebildet und durch 255 geteilt wird, um den Wert auf den Bereich [0, 1] zu normalisieren. Diese Werte werden gespeichert

und am Ende der Funktion zurückgegeben. Zum Schluss werden der Speicher des ursprünglichen und des skalierten Bildes freigegeben, um Speicherlecks zu verhindern.

Diese Funktion wird verwendet, um in der Ausführung von Predict das Bild des Nutzers zu übernehmen und auf das geforderte Format zu verkleinern.

5 Testing

5.1 Unittests

Beim Unit-Testing wird Code durch Code überprüft. So kann die Integrität des Projektes von einer einzelnen Funktion zu ganzen Klassen vom Grund auf überprüft werden.

In CKi werden die Klassen für das Neuron, den Layer und das Netzwerk getestet. Zusätzlich gibt es Tests zu den meisten Funktionen der Utility-Klasse.

Dabei werden Unit-Tests in CKi mit der „Google-Tests“ Bibliothek von Google implementiert.

5.1.1 Util (Utility-Klasse)

5.1.1.1 Sigmoid-Funktionstests

... prüfen, ob die Sigmoid- und ihre Ableitungsfunktion die erwarteten Werte zurückgeben.

5.1.1.2 Dateilese-Tests

... überprüfen das korrekte Lesen eines Integers aus einer Binärdatei und das korrekte Verhalten beim Versuch, nicht vorhandene Dateien zu lesen.

5.1.1.3 MNIST-Daten-Tests

...überprüfen das Einlesen von MNIST-Bild- und Label-Daten, einschliesslich der Überprüfung der Anzahl und der Normalisierung von Bildern sowie der Überprüfung der Gültigkeit von Labels.

5.1.2 Network

5.1.2.1 Index in akzeptiertem Bereich

... überprüft, ob die predict-Methode einen gültigen Index innerhalb des erwarteten Bereichs zurückgibt.

5.1.2.2 Training beeinflusst das Verhalten.

... testet, ob das Training des Netzwerks dessen Verhalten ändert, indem es die Ausgaben vor und nach dem Training vergleicht. Es wird festgestellt, dass sich die Ausgaben verändern sollten, was auf eine Modifikation des Netzwerks hinweist.

5.1.2.3 Netzwerkspeicher-Tests

... überprüfen die Funktionalität zum Speichern und Laden des Netzwerks. Nach dem Speichern und Laden wird erwartet, dass das Netzwerk wie vorgesehen funktioniert.

5.1.3 Layer

5.1.3.1 Initialisierung

... überprüft die korrekte Initialisierung der Neuronen in der Schicht mit der richtigen Anzahl von Gewichten.

5.1.3.2 Gewichtungen setzen

... testet das Setzen und Abrufen der Gewichte eines Neurons.

5.1.3.3 Output kalkulieren

... überprüft die Berechnung der Ausgänge der Neuronen basierend auf gegebenen Eingaben und gesetzten Gewichten.

5.1.3.4 Fehler kalkulieren

... testet die Fehlerberechnung der Schicht basierend auf den Ausgängen und den erwarteten Werten.

5.1.3.5 Änderung der Gewichtungen und Biases

... überprüft das Aktualisieren der Gewichte und Biases der Neuronen basierend auf einem berechneten Fehler.

5.1.4 Neuron

5.1.4.1 Aktivierungsfunktion

... überprüft, ob das Neuron die richtige Ausgabe liefert, wenn eine bestimmte Eingabe gegeben wird.

5.1.4.2 Ableitung der Aktivierungsfunktion

... untersucht, ob die Ableitung der Aktivierungsfunktion korrekt berechnet wird.

5.1.4.3 Fehlerberechnung

Der Fehler wird durch den Unterschied zwischen dem erwarteten und dem tatsächlichen Ausgang des Neurons bewertet.

5.2 Integrationstests

Es gibt keine Integrationstests, da keine (externen) Komponenten verwendet werden, auf denen das Projekt beruht und deren Verwendung nicht über Unittests abgedeckt werden.

5.3 Deployment-Tests

Es gibt keine Deployment-Tests, da dieses Produkt nicht ausgelegt wurde an die breite Öffentlichkeit weitergegeben zu werden. Somit werden keine Deployment-Tests benötigt, da jeder, der dieses Projekt verwenden möchte, dieses selbst „bauen“ muss.

6 Deployment

6.1 Building

Die Erstellung des Projektes CKI lässt sich in drei Schritte unterteilen. Es gibt jedoch einige Anforderungen.

6.1.1 Anforderungen

Das Projekt CKI ist zwar plattformunabhängig, kann also unter allen gängigen Betriebssystemen verwendet werden, muss jedoch erst für die entsprechende Plattform erstellt werden. Für diesen Erstell-Prozess sind folgende Anforderungen (Programme) unabdingbar:

- CMake
- C++ Compiler
- Internet

Je nach Betriebssystem oder IDE können schon alle oder einige dieser Anforderungen erfüllt sein.

6.1.1.1 CMake

CMake ist ein plattformübergreifendes Tool zur Automatisierung des Build-Prozesses, das es ermöglicht, Makefiles und Projekte für verschiedene Entwicklungsumgebungen zu generieren. Es verwendet eine Konfigurationsdatei „(*CMakeLists.txt*)“, um den Build-Prozess zu steuern. Unter Windows kann CMake von der offiziellen Website heruntergeladen und entweder über einen grafischen Installer oder über die Kommandozeile installiert werden.

Unter Linux und Mac kann CMake über den Package Manager der Wahl installiert werden. Die häufigsten sind Pacman, APT oder RPM unter Linux und Brew unter Mac.

Für das Projekt CKI wird CMake benötigt, um die Build-Konfiguration zu erstellen, externe Abhängigkeiten zu verwalten und das Projekt für verschiedene Entwicklungsumgebungen vorzubereiten, wodurch eine konsistente und effiziente Entwicklung ermöglicht wird. Dabei ist bei der Installation notwendig, dass die CMake Version kompatibel mit der **Version 3.26**, wie in der CMakeLists.txt spezifiziert, ist.

6.1.1.2 C++-Compiler

Ein C++-Compiler ist ein Software-Tool, das C++-Code in maschinenlesbaren Code übersetzt, sodass Programme ausgeführt werden können.

Unter Windows, Mac und Linux kann ein C++-Compiler durch die Installation einer Entwicklungsumgebung wie Visual Studio oder durch direkte Installation von GCC oder Clang eingebunden werden.

Für das Projekt CKI ist ein Compiler notwendig, da CMake, das für das Erstellen des Projektes verwendet wird, auf einen Compiler angewiesen ist, um den C++-Code in ausführbare Dateien zu übersetzen. CMake generiert Build-Konfigurationen, aber der eigentliche Kompilierungsprozess benötigt einen C++-Compiler. Dabei ist bei der Installation notwendig, dass der C++-Compiler den **C++-Standard 17** unterstützt, da das Projekt CKI auf diesen Standard aufgebaut wurde und moderne C++-Features nutzt.

6.1.1.3 Internet

Das Projekt CKI benötigt Internet, um externe Abhängigkeiten wie *googletest* für UnitTests und *nlohmann_json* für JSON-Verarbeitung automatisch herunterzuladen und zu integrieren. Diese Bibliotheken sind essenziell für das Funktionieren und Testen des Projektes. CMake verwaltet diesen Prozess, indem es die benötigten Pakete aus dem Internet lädt, was eine effiziente und konsistente Set-up-Umgebung über verschiedene Entwicklungsplattformen hinweg ermöglicht.

6.1.2 Ausführung

Um das Projekt CKI zu bauen und auszuführen, gibt es verschiedene Wege, die wegen unterschiedlicher Entwicklungsumgebungen stark variieren können. Dementsprechend ist hier nur beschrieben, wie CKI mithilfe von CMake, einem C++ Compiler und Internet in der Konsole gebaut werden kann.

Es folgen nun die notwendigen Schritte, um dies zu bewerkstelligen:

1. Erstellen Sie ein Build-Verzeichnis im Root-Verzeichnis des Projektes. Dies kann über ein grafisches Interface passieren oder über die Konsole. Der entsprechende Befehl unter Windows, Linux und Mac sollte „*mkdir*“ sein.
2. Wechseln Sie in das Build-Verzeichnis (Der entsprechende Befehl unter Windows, Linux und Mac sollte „*cd*“ sein.) und führen Sie „*cmake ..*“ aus, um die Build-Konfiguration zu generieren. Dieser Befehl variiert nicht unter den unterschiedlichen Betriebssystemen, setzt allerdings voraus, dass CMake installiert wurde. Siehe mehr unter Anforderungen 6.1.1.
3. Führen Sie anschliessend „*cmake --build .*“ aus, um das Projekt zu kompilieren. Dies setzt voraus, dass ein C++ Compiler installiert wurde. Siehe mehr unter Anforderungen 6.1.1.
4. Wenn Sie nun den Inhalt des Build-Verzeichnisses begutachten, sollten Sie unterschiedlichste Dateien und Verzeichnisse vorfinden, aber auch eine Datei namens CKI mit einer Endung einer ausführbaren Datei (unter Windows würde die Datei *CKI.exe* heissen). Diese ausführbare Datei ist das compilierte und fertige Produkt des Projektes CKI.

6.2 Installation

Dieser Abschnitt setzt voraus, dass die Schritte aus dem Abschnitt Building 6.1 erfolgreich ausgeführt wurden und das Produkt des Projekts CKI bereit zur Ausführung im Build-Verzeichnis zu finden ist.

6.2.1 Installation mit CMake

Nach dem Build-Vorgang können Sie das Projekt CKI installieren, indem Sie einen Installationsbefehl über CMake ausführen. Dieser Vorgang kopiert die erforderlichen ausführbaren Dateien, Bibliotheken und eventuell weitere Ressourcen in vordefinierte Verzeichnisse. Normalerweise wird dies durch den Befehl „*cmake --install .*“ im Build-Verzeichnis erreicht. Stellen Sie sicher, dass Sie über die erforderlichen Berechtigungen verfügen, um Installationen in den Zielverzeichnissen durchführen zu können. Die genaue Konfiguration des Installationspfades und anderer Parameter kann in der „*CMakeLists.txt*“ festgelegt werden. Diese Konfiguration muss jedoch vom Installierenden vorgenommen werden, da das Projekt CKI nicht für eine langfristige Installation gedacht ist, da es per se keinen Mehrwert für einen Endnutzer bringt.

6.2.2 Installation ohne CMake

Achtung! Dieser Vorgang wurde nur flüchtig unter Windows getestet und kann daher nicht für Linux oder Mac garantiert werden.

Nach dem Build-Vorgang kann die Datei, welche verantwortlich für die Ausführung ist (unter Windows *CKI.exe*), in ein beliebiges anderes Verzeichnis kopiert werden. Unter Windows sollte das Produkt des Projekts CKI (*CKI.exe*) ausführbar bleiben; für Linux und Mac kann aufgrund fehlender Versuchen keine Aussage getroffen werden.

6.2.3 Hinzufügung von vortrainierten Gewichtungen und Biases

Um im Projekt CKI, einer Installation oder einem Build, vorhandene Gewichtungen und Biases für das neuronale Netzwerk hinzuzufügen, wird eine JSON-Datei benötigt, welche dem Netzwerk als Speicher für die Gewichtungen und Biases dient. Diese Datei liegt standardmässig im Root-Verzeichnis der installierten oder gebauten Applikation (also im gleichen Verzeichnis wie die ausführbare Datei/Applikation).

Somit können alle Gewichtungen und Biases zwischen unterschiedlichen neuronalen Netzwerken (welche für die gleiche Aufgabe trainiert wurden) ausgetauscht werden. Dabei gilt allerdings zu beachten, dass das Netzwerk gleich aufgebaut sein muss. Das heisst, die Netzwerke

benötigen die gleiche Anzahl an Layer und Neuronen pro Layer. Es muss die „load_network“-Funktion ausgeführt werden, vor der Verwendung des Netzwerkes und die JSON-Datei muss richtig benannt sein.

Sollte die Standard-Version der Main-Datei verwendet werden, heisst diese Datei „network.json“ und wird bei Nichtvorhandensein automatisch beim ersten Trainingsvorgang erstellt. Somit wird mindestens beim Erststart ein Fehler erscheinen, dass die JSON-Datei nicht gelesen werden konnte. Dieser kann allerdings ignoriert werden, wenn keine vorhandenen Gewichtungen und Biases manuell zum Netzwerk hinzugefügt worden sind.

6.3 Verwendung

Das Projekt CKI nutzt ein neuronales Netzwerk, um mit dem MNIST-Datensatz zu arbeiten, einem Standarddatensatz für Handschrifterkennung. Dabei bietet es eine flexible Anwendung für maschinelles Lernen mit Fokus auf Bilderkennung. Es ermöglicht Benutzern, ein Convolutional Neural Network (CNN) mit dem MNIST-Datensatz für Handschrifterkennung zu trainieren, zu verifizieren und Vorhersagen für einzelne Bilder zu treffen. Durch Kommandozeilenargumente kann der Benutzer wählen, ob das Netzwerk trainiert, dessen Genauigkeit überprüft oder ein Bild klassifiziert werden soll.

6.3.1 Interaktion

Um über die CLI mit CKI zu interagieren, muss der Benutzer das Programm mit spezifischen Kommandozeilenargumenten ausführen. Diese Argumente steuern die Aktionen des Programms, wie das Trainieren des Netzwerkes, die Überprüfung seiner Genauigkeit oder die Klassifizierung eines spezifischen Bildes. Eine genaue Definition der Befehle ist im Abschnitt Design 3.1 zu finden. Der Übersichtlichkeit sind diese aber hier nochmals zusammengefasst:

1. Training

- **Befehl:**

```
$ cki --train -l [labels] -i [images] -e [epochs] -lr [learningrate]
```

- **Verwendung:** Trainiert das Netzwerk mit Bildern und Labels. Die Optionen für Epochen (-e) und Lernrate (-lr) sind optional. -l und -i benötigen einen Dateipfad zu den entsprechenden UByte-Dateien.

2. Verifizierung

- **Befehl:**

```
1 $ cki --verify -l [labels] -i [images]
```

- **Verwendung:** Überprüft die Genauigkeit des trainierten Netzwerks mit einem Testdatensatz. -l und -i benötigen einen Dateipfad zu den entsprechenden UByte-Dateien.

3. Vorhersage

- **Befehl:**

```
1 $ cki [file]
```

- **Verwendung:** Macht eine Vorhersage für ein einzelnes Bild. Das Bild kann im JPG- oder PNG-Format sein.

4. Hilfe

- **Befehl:**

```
1 $ cki --help
```

- **Verwendung:** Zeigt eine Liste aller verfügbaren Befehle an.

Literaturverzeichnis

- [1] 3Blue1Brown. Backpropagation-kalkül | kapitel 4, deep learning. 3Blue1Brown, 03.11.2017. Accessed on 19 February 2024.
- [2] 3Blue1Brown. Aber was ist ein neuronales netzwerk? | kapitel 1, deep learning. 3Blue1Brown, 15.10.2017. Accessed on 19 February 2024.
- [3] AndroidKT. Extract images from mnist idx3 ubyte file format in python. AndroidKT, 11.06.2023. Accessed on 23 November 2023.
- [4] Yann LeCun; Corinna Cortes; Christopher J.C. Burges. The mnist database of handwritten digits. MNIST, 22.06.2002. Accessed on 23 November 2023.
- [5] databasecamp. Was ist die relu-funktion (rectified linear unit)? Database-Camp, 26.04.2023. Accessed on 19 February 2024.
- [6] IBM. What are convolutional neural networks? IBM, 16.12.2021. Accessed on 23 November 2023.
- [7] Paul Hudak; Mark P. Jones. Haskell vs. ada vs. c++ vs awk vs. ... an experiment in software prototyping productivity. Yale Edu, 10.1994. Accessed on 12 February 2024.
- [8] Andrej Karpathy. Yes you should understand backprop. Medium, 19.12.2016. Accessed on 19 February 2024.
- [9] Sven Nüesch. *Artificial Intelligence (AI)*. KS Frauenfeld, 07.04.2023. Only access over Sven Nüesch.
- [10] Amrita Pathak. Convolutional neural networks (cnns): Eine einföhrung. Geekflare, 24.09.2023. Accessed on 23 November 2023.

- [11] Abhijit Roy. An introduction to gradient descent and backpropagation. Towards Datas-
cience, 14.06.2020. Accessed on 19 February 2024.
- [12] Sumit Saha. A comprehensive guide to convolutional neural networks — the eli5 way.
Towards Datascience, 15.12.2018. Accessed on 23 November 2023.
- [13] Dave Jarvis; stimms. Development time in various languages. Stack Overflow,
12.12.2009. Accessed on 12 February 2024.
- [14] Fei-Fei Li; Justin Johnson; Serena Yeung. Backpropagation and gradients. Stanford,
13.04.2017. Accessed on 19 February 2024.
- [15] Ph.D. Stephen F. Zeigler. Comparing development costs of c and ada. Archive Adaic,
30.03.1995. Accessed on 12 February 2024.

Abbildungsverzeichnis

1.1	CNN Model Aufbau Grafik von researchgate.net	6
1.2	Use-Case-Diagramm	10
1.3	Ablaufdiagramm für die Interpretation von Nutzer-gelieferten Einzelbildern . . .	11
1.4	Ablaufdiagramm vom Training des CNN	12
1.5	Ablaufdiagramm von einem Test des CNN	13
1.6	Visuelle Darstellung der Funktionsweise eines CNN von blog.goodaudience.com	14
2.0	Das Gantt-Diagramm für das Projekt CKI	28
3.1	Darstellung einer Ziffer und deren Code von blog.csdn.net	31
3.2	Das Klassendiagramm, Grundaufbau der Applikation (überholt)	32
4.1	Klassendiagramm nach der Realisation	36
4.2	Abbildung von 3 isolierten Neuronen (Inspiriert von 3Blue1Brown)	41
4.3	Abbildung von einigen Neuronen (Inspiriert von 3Blue1Brown)	42
4.4	4 Neuronen für die Back-Propagation (Bild von 3Blue1Brown)	44
4.5	Kostenbaum in der Back-Propagation (Bild von 3Blue1Brown)	45
4.6	5 Neuronen in einem Netzwerk (Bild von 3Blue1Brown)	48
4.7	Sigmoid Funktion	54

Tabellenverzeichnis

1.0	Kostenanalyse	18
1.0	Vergleich bestehendes Produkt gegen eigene Kreation	19
1.0	Nutzwertanalyse	20
2.0	Rollenverteilung im Projekt CKI	21
2.0	Meilensteine	23

Codeverzeichnis

3.0	Konsolenbefehl Grundlage	28
3.0	Design Konsolenbefehl Trainieren	28
3.0	Design Konsolenbefehl Verifizieren	29
3.0	Design Konsolenbefehl Prediction	29
3.0	Design Konsolenbefehl für Hilfe	30
4.0	Train Funktion aus der Netzwerk-Klasse	49
4.0	Verify Funktion aus der Netzwerk-Klasse	50
4.0	Predict Funktion aus der Netzwerk-Klasse	50
4.0	Forward-Propagation Funktion aus der Netzwerk-Klasse	51
4.0	Backward-Propagation Funktion aus der Netzwerk-Klasse	52
4.0	Anpassung von Gewichtungen und Biases aus der Layer-Klasse	53
4.0	Sigmoid & Ableitung aus der Util-Klasse	53
4.0	MNIST-Reader aus der Util-Klasse	55
4.0	Utility Funktion für den MNIST-Reader aus der Util-Klasse	55
4.0	Bild Vorverarbeitung für Nutzer-gelieferte Bilder	56
6.0	Implementation des Trainingsbefehls	66
6.0	Implementation des Verifizierungsbefehls	67
6.0	Implementation des Predictionbefehls	67
6.0	Implementation des Befehls für Hilfe	67










Anhang

A Fortschrittsberichte

Sie sind in chronologischer Reihenfolge.

Projektfortschrittsbericht				
Projekt: CKi		Arbeitspaket: Vorbereitungen		
Mitarbeiter: Simeon Stix		Kalenderwoche: 42	Datum: 20.10.2023	
Projektkurzbeschreibung: Eine Desktop-Applikation, die handschriftliche Zahlen in verwendbare digitale Zahlen umwandelt				
Projektfortschritt (in Prozent): 10%				
Arbeiten, Ziele des Projektschrittes: Abschluss der Planung, der Analyse und des Designsvorbereitungen				
Projektstand	Status			Begründung (nur bei gelb und rot)
Termine				
Kosten				
Produkt				
Störungen/Probleme/Risiken: Verzögerungen an Arbeit 30.09.23, gesundheitliche Gründe; Kompensierbar durch erhöhte Arbeitsgeschwindigkeit				
Notwendige Entscheidungen: Konsolen Interaktionen (Verwendbarkeit der Applikation)				
Anlagen:				

NÜE 16.5.2021

Projektfortschrittsbericht				
Projekt: CKi		Arbeitspaket: Design & Beginn Implementation		
Mitarbeiter: Simeon Stix		Kalenderwoche: 45	Datum: 10.11.2023	
Projektkurzbeschreibung: Eine Desktop-Applikation, die handschriftliche Zahlen in verwendbare digitale Zahlen umwandelt.				
Projektfortschritt (in Prozent): 15%				
Arbeiten, Ziele des Projektschrittes: Design Abschluss & Implementation der KI-Klassen				
Projektstand	Status			Begründung (nur bei gelb und rot)
Termine				unerwartet hohe Prüfungsdichte & Fehl kalkulation des Arbeitsaufwands des Designs & Implementation
Kosten				Mehr Arbeitszeit benötigt
Produkt				Produkt ist nicht am geplanten Stand
Störungen/Probleme/Risiken:				
Notwendige Entscheidungen: Dokumentation von Markdown zu LaTeX abgewandelt				
Anlagen:				

NÜE 16.5.2021

Projektfortschrittsbericht			
Projekt: CKi		Arbeitspaket: Design & Beginn Implementation	
Mitarbeiter: Simeon Stix		Kalenderwoche: 48	Datum: 02.12.2023
Projektkurzbeschreibung: Eine Desktop-Applikation, die handschriftliche Zahlen in verwendbare digitale Zahlen umwandelt.			
Projektfortschritt (in Prozent): 20%			
Arbeiten, Ziele des Projektschrittes: Dokumentation auf Stand des Realisationsbeginn.			
Projektstand	Status	Begründung (nur bei gelb und rot)	
Termine	<input type="radio"/> <input checked="" type="radio"/> <input type="radio"/>	Realisation ist noch nicht auf dem gewünschten Stand.	
Kosten	<input checked="" type="radio"/> <input type="radio"/> <input type="radio"/>		
Produkt	<input checked="" type="radio"/> <input type="radio"/> <input type="radio"/>		
Störungen/Probleme/Risiken: Misskalkulation des Arbeitsaufwandes			
Notwendige Entscheidungen:			
Anlagen:			

NÜE 16.5.2021

Projektfortschrittsbericht				
Projekt: CKi		Arbeitspaket: Abschluss Produkt		
Mitarbeiter: Simeon Stix		Kalenderwoche: 1	Datum: 05.01.2024	
Projektkurzbeschreibung: Eine Desktop-Applikation, die handschriftliche Zahlen in verwendbare digitale Zahlen umwandelt.				
Projektfortschritt (in Prozent): 55%				
Arbeiten, Ziele des Projektschrittes: Abschluss der Realisation				
Projektstand	Status			Begründung (nur bei gelb und rot)
Termine				
Kosten				
Produkt				Netzwerk ist noch nicht bereit. Back-Props funktioniert noch nicht (korrekt).
Störungen/Probleme/Risiken: Das Wetter war zu gut. Ich musste einfach Skifahren. So hatte ich keine Zeit für die Projektarbeit, sollte aber Rückstand kompensieren können				
Notwendige Entscheidungen: Skifahren oder Arbeiten.				
Anlagen:				

NÜE 16.5.2021

Projektfortschrittsbericht				
Projekt: CKi		Arbeitspaket: Abschluss		
Mitarbeiter: Simeon Stix		Kalenderwoche: 06	Datum: 10.02.2024	
Projektkurzbeschreibung: Eine Desktop-Applikation, die handschriftliche Zahlen in verwendbare digitale Zahlen umwandelt.				
Projektfortschritt (in Prozent): 99%				
Arbeiten, Ziele des Projektschrittes: Abschluss der Realisation & Dokumentation				
Projektstand	Status			Begründung (nur bei gelb und rot)
Termine				Die Dokumentation ist noch nicht abgeschlossen.
Kosten				
Produkt				
Störungen/Probleme/Risiken: Synchronisation zwischen GitHub & GitLab läuft nicht einwandfrei -> verzögerter Projektfortschritt auf GitLab, etc. -> noch zu lösen				
Notwendige Entscheidungen: Reduktion der Realisation auf die Must-Have's				
Anlagen: https://github.com/MrS-E/CKi				

NÜE 16.5.2021