

TARTU ÜLIKOOL
Loodus- ja täppisteaduste valdkond
Arvutiteaduse instituut
Informaatika õppekava

Oliver Savolainen

Konvolutsiooniliste neurovõrkude kalibreerimine järkjärgulise külmutamise meetodiga

Bakalaureusetöö (9 EAP)

Juhendaja: Meelis Kull, PhD

Tartu 2023

Konvolutsiooniliste neurovõrkude kalibreerimine järkjärgulise külmutamise meetodiga

Lühikokkuvõte: Neurovõrke on väga edukalt kasutatud mitmete valdkondade ja ülesannete juures, neist üks parimaid näiteid on konvolutsiooniliste neurovõrkude kasutamine piltide klassifitseerimiseks. Samas ei saa tihti neid mudeleid siiski usaldada, sest nad on üldjuhul liiga enesekindlad oma ennustustes ehk määravad liiga kõrge tõenäosuse ennustatavale klassile.

Selle tõttu on vaja mudeleid kalibreerida. Probleemi lahendamiseks on loodud mitmeid kalibreerimismeetodeid, neist üks tõhusamaid on temperatuuri skaleerimine. Selles töös on kalibreerimiseks, täpsemalt liigse enesekindluse vähendamiseks uurimise all järkjärguline külmutamine. Järkjärgulise külmutamise all mõeldakse meetodit neurovõrgu treenimiseks, kus mingitel hetkel lõpetatakse valitud kihtide kaalude muutmine.

Töös kasutati Kängsepa 2018. aasta magistritöö lähtekoodi ja tulemusi. Kõigepealt leiti ühe mudeli ja andmestiku põhjal parimad viisid meetodi rakendamiseks. Valiti kaks järkjärgulise külmutamise skeemi ning seejärel implementeeriti Kängsepa töö abil mitmeid konvolutsioonilisi neurovõrke, rakendati neile järkjärgulist külmutamist ning treeniti neid. Saadud mudelite tulemusi võrreldi Kängsepa töö tulemustega. Kuigi ei saa väita, et uuritud mõõdikute põhjal aitas järkjärguline külmutamine liigset enesekindlust vähendada, siis vähendas külmutamine ajaliste ressursside kasutamist ning saavutas samas paljude mudelite juures sarnaseid tulemusi.

Võtmesõnad: konvolutsioonilised neurovõrgud, kalibreerimine, piltide klassifitseerimine, residuaalne neurovõrk, järkjärguline külmutamine

CERCS: P176 Tehisintellekt

Calibration of Convolutional Neural Networks with Gradual Freezing

Abstract: Neural networks have been successfully used for various tasks in various fields, with one of the best examples being the use of convolutional neural networks for image classification. However, these models often cannot be trusted, as they are frequently excessively confident in their predictions, assigning too high of a probability to the predicted class.

As a result, it is necessary to calibrate the models. Several calibration methods have been developed to address this problem, one of the most efficient ones is temperature scaling. This study investigates gradual freezing for calibration, specifically for reducing excessive confidence. Gradual freezing refers to a method of training in which the updating of weights for selected layers is stopped at certain moments.

The source code and results of Kängsepp's 2018 master's thesis were used in this work. First, the best ways to implement the method were found based on one model and dataset. Two gradual freezing schemes were chosen, and then several convolutional neural networks were implemented using Kängsepp's work and trained with gradual

freezing. The obtained models' results were compared with Kängsepp's results. Although it cannot be claimed that gradual freezing helped reduce excessive confidence based on the metrics used, freezing did reduce the time required for training while still achieving similar results for many models.

Keywords: convolutional neural networks, calibration, image classification, residual networks, gradual freezing

CERCS: P176 Artificial intelligence

Sisukord

1	Mõisted ja definitsioonid	6
1.1	Neurovõrk	6
1.2	Konvolutsiooniline neurovõrk	7
1.3	Tõenäosuslik klassifikaator	9
1.4	Määramatuse kalibreerimine	9
1.5	Kalibreerimise mõõtmine ja meetodid	11
1.6	Järkjärguline külmutamine	12
2	Järkjärgulise külmutamise meetodi implementeerimine ja uurimine	15
2.1	Ülevaade järkjärgulise külmutamise implementeerimisest	15
2.1.1	Kasutatud tarkvara	15
2.1.2	Kasutatud andmestikud	15
2.1.3	Järkjärgulise külmutamise implementeerimine	16
2.2	Järkjärgulise külmutamise eksperimendid	17
2.2.1	Eksperimenteerimise esimese etapi kirjeldus	17
2.2.2	Eksperimenteerimise esimese etapi tulemuste analüüs	19
2.2.3	Eksperimenteerimise teise etapi kirjeldus	22
3	Eksperimenteerimise teise etapi tulemuste analüüs	24
3.1	Veamäära võrdlus	24
3.2	NLLi võrdlus	25
3.3	ECE võrdlus	26
3.4	Ajaline võrdlus	26
3.5	Järeldused	27
	Viidatud kirjandus	30
	Lisad	32
	II. Litsents	35

Sissejuhatus

Sügavõpe on juba aastaid olnud väga tähtis osa paljude ülesannete ja valdkondade juures. Osade ülesannete juures, mis kasutavad sügavõpet, on jõutud inimestega sarnase soorituseni, heaks näiteks on piltide klassifitseerimine. Kasutades neurovõrke, täpsemalt konvolutsioonilisi neurovõrke, on mitmetel suurtel pildiandmestikel jõutud väga kõrge täpsuseni. [1]. Kuid tihti ei saa neid tulemusi usaldada. Kui mudel ennustab mingite piltide põhjal, mis klassi see pilt kuulub, siis tihti ta väljastab ennustades liiga suure tõenäosuse ühe klassi kohta ning liiga madalad tõenäosused teiste klasside kohta [2]. Sellistes olukordades võib tulemuste tõlgendaja arvata, et ennustus on kindlasti õige, aga tihti on ennustatud tõenäosus kõrgem kui reaalselt sellise tõenäosusega saadud täpsus.

Seda nimetatakse liigseks enesekindluseks, mis on konvolutsiooniliste neurovõrkude puhul tihti probleemiks. Selle parandamiseks kasutatakse kalibreerimist. Loodud on mitmeid kalibreerimismeetodeid, neist üks tõhusamaid on temperatuuri skaleerimine [2]. Selles töös on kalibreerimiseks, täpsemalt liigse enesekindluse vähendamiseks uurimise all järkjärguline külmutamine. Kihi külmutamine tähendab, et selle kihi kaale enam tagasilevi algoritmis ei muudeta, need jäävad samaks kuni treeningprotsessi lõpuni või kuni need pole enam külmutatud. Järkjärgulisel külmutamisel külmutatakse treeningprotsessi käigus mingitel hetkedel valitud kihte. Selle kasutamine aitab mudeli keerukust vähendada, mis võiks viia ülesobitamise ning liigse enesekindluse vähendamiseni.

Siinse töö eesmärk on implementeerida järkjärguline külmutamine ning otsida vastust küsimusele, kas see aitab vähendada konvolutsiooniliste neurovõrkude liigset enesekindlust. Samuti tahetakse leida, millised on parimad viisid selle meetodi rakendamiseks. Töös on kasutatud Markus Kängsepa 2018. aasta magistritöö lähtekoodi ja tulemusi [3]. Kui esimeses sisulises peatükis selgitatakse mõisteid ning järkjärgulise külmutamise kasutamist kalibreerimiseks, siis teises peatükis kirjeldatakse meetodi ja Kängsepa töös kasutatud mudelite implementeerimist ja esimesi eksperimente. Kolmandas peatükis on võrreldud lõplikke saadud tulemusi Kängsepa töö omadega.

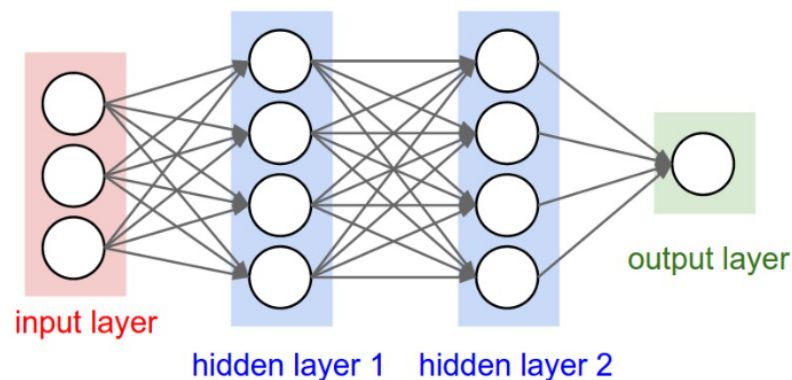
1 Mõisted ja definitsioonid

Selles töö osas selgitatakse peamiseid töös kasutatavaid mõisteid, definitsioone ning üldist teoreetilist tausta. Selle toel seletatakse peatüki lõpuks, miks kasutatakse järkjärgulist külmutamist liigse enesekindluse vähendamiseks.

1.1 Neurovõrk

Töös kasutatakse eesmärgi täitmiseks konvolutsioonilisi neurovõrke. Järgnevad kaks alampeatükki nende tutvustuseks on kirjutatud Goodfellowi jt [1] raamatu „Deep Learning“ põhjal.

Tehisnärvivõrk ehk neurovõrk (*neural network*) on arvutuslik arhitektuur, mis loodi ajast inspireerituna. Selle eesmärk on üldjuhul sisendi põhjal tagastada väljund millegi ennustamiseks. Tuntuimad neurovõrgud on pärilevi neurovõrgud (*feed-forward networks*), milles kõik neuronid on ühendatud suunaga väljundi poole moodustades suunatud tsükliteta graafi. Need koosnevad kihtidest, mis omakorda koosnevad omavahel ühendatud neuronitest ehk sõlmedest (*neurons*). Joonisel 1 on näha lihtne ülevaade pärilevi neurovõrgu arhitektuurist - väga tihti on neurovõrgul olemas sisendkiht (*input layer*) ja väljundkiht (*output layer*), üldjuhul ka nende vahel varjatud kihid (*hidden layers*). Igal kihil on sisendid, väljundid ja aktivatsioonifunktsioon (*activation function*), mis üldjuhul piirab neuroni väljundi vajalikku vahemikku. Neuronitevahelistel ühendustel on kaalud, millest sõltub sisendi tähtsus. Just kaalude väärtust muudetakse mudeli treeningprotsessi käigus.



Joonis 1. Neurovõrgu arhitektuuri ülevaade.

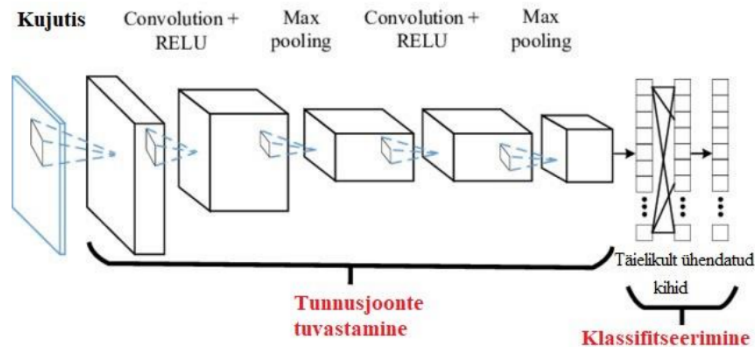
Nendes võrkudes rakendatakse tagasilevi algoritmi (*back-propagation*). Tagasilevi algoritmis leitakse kaofunktsiooni (*loss function*) väärtus ehk kui suur viga tehti ennustades ning vastavalt sellele muudetakse kaalusid alustades viimasest kihist kuni esimeseni. Eesmärk on minimeerida kaofunktsiooni väärtus ehk jõuda koondumiseni, mida tehakse muutuste suunda ja suurust näitava vektori abil, mida kutsutakse gradiendiks. Kui need muutused tehakse kõikide treeningandmete põhjal, on tegemist gradientlaskumisega. Kui aga ennustatakse väiksema osa ehk ploki põhjal gradiendi väärtust, siis on tegemist stohhastilise gradientlaskumisega (*stochastic gradient descent*).

See on kõige tavapärasem neurovõrkude optimeerimise algoritm, mida kasutatakse ka siin töös. Algoritmi töö sõltub paljuski hüperparameetrite väärtustest. Hüperparameetrid on väärtused, mida mudelid ei õpi treenides, need on enne kaasa antud. Gradiendi ennustused toimuvad iga miniploki (*batch*) järel, mille suurus üldjuhul on määratud hüperparameetriga. Kui käia läbi kõik treeningandmed miniplokkide kaupa, on läbitud üks epohh (*epoch*). Epohhide arv on üldjuhul hüperparameeter, mille järgi määratakse treeningu pikkus. Stohhastilist gradientlaskumist kasutades saab ette anda veel hüperparameetreid, näiteks õpisammu (*learning rate*). Õpisamm määrab, kui suur samm astutakse gradiendi suunas kaalusid muutes. Lisaks on siin töös kasutuses inertsitegur (*momentum*), mis akumulereib eelnevate gradientide väärtuseid, et leida nende abil, kui palju tuleks kaalusid muuta, et kiiremini minimeerida kaofunktsiooni. Täpsemalt on siin töös tegemist Nesterovi inertsiteguriga (*Nesterov momentum*), kus akumulereimine toimub pärast inertsiteguri rakendamist. Lisaks kõigile nendele meetoditele ja hüperparameetritele on treeningu kiirendamiseks kasutuses ka tiiptasemel riistvara, eelkõige GPUd ehk graafikaprotsessorid (*Graphics processing units*) - graafikakaartidel asuvad mikroprotsessorid, mis suudavad paralleelselt palju arvutusi sooritada ning tänu sellele sobivad väga hästi neurovõrkude treenimiseks.

Kasutades kirjeldatud arhitektuuri ja meetodeid, saab neurovõrgule ette anda lihtsalt kõik andmepunktid ning tihti jõuda väga hea tulemuseni. Võrreldes mitmete teiste masinõppe mudelitega pole vaja ise tunnuseid eraldi valida või luua. Piisab lihtsalt andmepunktidest, millel on vastav väljund ette antud ja mudel leiab ise, milliseid tunnuseid tasub kasutada. Veidi erinevat arhitektuuri kasutatakse aga siis, kui on palju sisendeid, kuid sisendid võivad vaid osade väljunditega seotud olla. Kõige tavalisemalt on sisendiks 2D pildi andmed, kus kolmes järjendis on iga pildi piksli kohta vastavalt sinise, punase või rohelise värvi intensiivsuse väärtus. Sellisel juhul kasutatakse enamasti konvolutsioonilisi neurovõrke.

1.2 Konvolutsiooniline neurovõrk

Konvolutsioonilistel neurovõrkudel on erinevalt tavalistest pärilevi neurovõrkudest lisaks konvolutsioonilised kihid, mis koosnevad kolmest etapist. Joonisel 2 on toodud abstraktne joonis, kuidas tunnuseid nende abil leitakse, selles sisalduvaid mõisteid kirjeldatakse järgmistes lõikudes.



Joonis 2. Konvolutsioonilise neurovõrgu tööd selgitav abstraktne joonis [4].

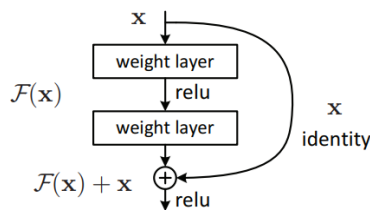
Esimeses etapis toimuvad matemaatilised operatsioonid nimega konvolutsioonid (*convolutions*), kus kahe funktsiooni põhjal luuakse kolmas funktsioon, mis kirjeldab, kuidas üks funktsioon mõjutab teist. Sisendiks olevad kaks funktsiooni konvolutsioonilistes neurovõrkudes on mitmedimensionaalne järjend andmeid ja mitmedimensionaalne järjend õpitavaid parameetreid, väljundiks olev funktsioon kannab üldjuhul tunnuste kaardi (*feature map*) nimetust.

Teises etapis rakendatakse mittenegatiivset lineaarfunktsiooni ehk ReLU, mis asendab tunnuste kaardi negatiivsed väärtused nullidega. See toob mittelineaarsust võrku ehk aitab leida andmetest keerulisemaid tunnuseid.

Kolmandas etapis tuleb kasutusse ahendus (*pooling*), mis asendab üksikute neuro- nite väljundid nende naabrite kokkuvõtva väljundiga. See vähendab väikeste muutuste tähtsust. Konvolutsioonikihte ja kirjeldatud etappe esineb kõikides konvolutsioonilistes neurovõrkudes, aga mudelid võivad muude osade tõttu siiski selgelt erineda.

Selles töös on kasutuses neli konvolutsiooniliste neurovõrkude arhitektuuri: ResNet, ResNet SD, Wide ResNet ja DenseNet. Kõik need on tuntud mudelid, mida kasutati Kängsepa 2018. aasta magistritöös [3].

ResNet ehk residuaalne neurovõrk (Residual Network) on konvolutsiooniline neurovõrk, mis kasutab lisaks tavalisele konvolutsioonilise neurovõrgu arhitektuurile jääk-ühendusi ehk annab varasemate kihtide sisendi otse hilisematele kihtidele samasusteisendusega (*identity function*) [5]. Näide sellest on toodud joonisel 3.



Joonis 3. Residuaalse neurovõrgu jääkühenduste kasutus [5].

ResNet SD ehk residuaalne neurovõrk stohhastilise sügavusega (*ResNet with Stochastic Depth*) on ResNeti modifikatsioon, mis kasutab stohhastilist sügavust ehk lülitab juhuslikult välja osad kihid treeningprotsessis, mis suurendab kiirust ja vähendab üleso-bitamist [6]. Wide ResNet on ResNeti variant, kus kihid on laiemad, aga neid on vähem ehk mudeli sügavus on väiksem [7]. DenseNet (*Dense Convolutional Network*) omab erinevalt tavalistest konvolutsioonilistest neurovõrkudest $L(L + 1)/2$ ühendust iga kihi ja sellele järgneva kihi vahel, L tähistab mitmenda kihiga on tegemist [8]. Neid kõiki mudeleid kasutatakse tõenäosuslike klassifikaatoritena.

1.3 Tõenäosuslik klassifikaator

Flachi 2012. aasta ülevaate kohaselt [9] on klassifikaator funktsioon, mis andmepunkti X tunnuste väärtuste põhjal tagastab hulka Y kuuluva klassimärgendi. Binaarne klassifikaator ennustab vaid kahe erineva kahe klassi korral. Üldjuhul toimub see sündmuse toimumise kohta. Mitmeklassiline klassifikaator ennustab seda rohkem kui kahe klassi korral. Näiteks binaarne klassifikaator ennustab vaid, kas sademeid tuleb või mitte, mitmeklassiline võib aga ennustada, kas tuleb vihma, lund, rahet või ei tule sademeid.

Tõenäosuslik klassifikaator on funktsioon, mis iga $x \in X$ korral tagastab tõenäosusvektori (p_1, \dots, p_k) , kus p_i tähistab ennustatud tõenäosust, et nende väärtuste põhjal vastab x -le i -nda klassi märgend [9]. Tõenäosusvektor tagastatakse neurovõrkude korral üldjuhul softmax funktsiooniga, mis viib mudeli arvutatud reaalarvulised väljundid vahemikku 0-st 1-ni, kus kõigi väärtuste summa on kokku täpselt 1 [2]. Töös kasutataksegi klassifikaatori mõistet tõenäosusliku klassifikaatori tähenduses ja proovitakse just selle määramatust kalibreerida.

1.4 Määramatuse kalibreerimine

Olgu kasutuses binaarne klassifikaator, kus ennustatakse sisendi põhjal, kas märgend on 0 või 1. Enesekindlus olgu siin ennustatud tõenäosus, et märgend on 1. Määramatust kalibreerides tahetakse vähendada enesekindluse ja reaalse tõenäosuse vahelist eksimust.

Reaalne tõenäosus olgu siin märgendiga 1 andmepunktide osakaal nendest andmepunktidest, mille korral klassifikaator on vaadeldava enesekindlusega. Näiteks tahetakse, et andmepunktid, mille korral on klassifikaator 70% enesekindlusega, oleksid 70% korral juhtudest märgendiga 1. Kui see on täpselt nii ehk eksimus on 0, on mudel perfektselt kalibreeritud. Samas ei tohi kalibreerides täpsus ideaalis üldse kehvemaks minna, sest täpsus on klassifitseerides ikkagi kõige tähtsam mõõdik. Mitmeklassilise klassifikaatori korral aga nii lihtsalt kalibreeritust hinnata ei saa, sest kalibreerida on vaja üldjuhul mitme klassi ennustusi [2].

Kull jt [10] on defineerinud klassifikaatori kalibreerimise kolmel erineval viisil. Mitmeklassiliselt kalibreeritud või lihtsalt kalibreeritud on klassifikaator \hat{p} siis, kui iga tõenäosusvektori $q = (q_1, \dots, q_k)$ korral vastab väljundiks olev vektor $\hat{p}(x) = (\hat{p}_1(x), \dots, \hat{p}_k(x))$ tegelikule klassijaotusele. Valemis (1) on kirjas, kuidas see väljendub:

$$P(Y = i \mid \hat{p}(X) = q) = q_i \text{ iga } i = 1, \dots, k \text{ korral.} \quad (1)$$

Enesekindluse järgi kalibreerituks nimetatakse klassifikaatorit, kui kõikide ennustavate maksimaalsete tõenäosuste c korral vastab tegelikult sellesse klassi kuuluvate andmepunktide hulk c väärtusele. Seda väljendab valem (2):

$$P(Y = \arg \max (\hat{p}(X)) \mid \max (\hat{p}(X)) = c) = c. \quad (2)$$

Klassifikaatorit nimetatakse klassikaupa kalibreerituks, kui iga klassi i korral vastab iga ennustatav tõenäosus $\hat{p}_i(X)$ tegelikult sellesse klassi kuuluvate andmepunktide hulgale. Välja on toodud see valemis (3):

$$P(Y = i \mid \hat{p}_i(X) = q_i) = q_i. \quad (3)$$

Kui klassifikaator on mitmeklassiliselt kalibreeritud, siis on ta enesekindluse ja klassikaupa kalibreeritud, vastupidiselt ei pruugi see kehtida [10].

Olgu kasutuses klassifikaator, mis ennustab, kas tuleb vihma, lund, rahet või ei tule sademeid. Kogu andmestikus suurusega 10 olgu 4 andmepunkti, mille märgend vastab vihmale, 2 lumele vastava märgendiga andmepunkti, 1 rahele vastava märgendiga andmepunkt ning ülejäänud 3 puhul olgu tegemist sademeteta klassi kuuluvate andmepunktidega. Kui klassifikaatori ennustuseks iga andmepunkti korral on tõenäosusvektor (0,4; 0,2; 0,1; 0,3), on see kõikide definitsioonide järgi kalibreeritud. Mitmeklassiliselt on see mudel kalibreeritud, sest selle väljund vastab tegelikult klassijaotusele (0,4;0,2;0,1;0,3). Enesekindluse järgi on kalibreeritud, sest maksimaalne tõenäosus 0,4 väljundis vastab enim esineva klassi tõenäosusele. Klassikaupa on mudel kalibreeritud, sest 0,4 vastab vihma tähistava klassi kuuluvatele andmepunktide osakaalule, 0,2 lund tähistava klassi kuuluvatele andmepunktide osakaalule, 0,1 rahe klassi andmepunktide osakaalule ja 0,3 sademeteta klassi andmepunktide osakaalule. Definitsioonidele vastamine aga ei tähenda, et praktiliselt on käes suurepärase mudel, sest selle täpsus oleks selle näite

korral vaid 40%, sest iga kord ennustatakse, et tegemist on vihmaga. See näide näitab, et definitsioonidest ainult ei piisa, et leida, kui hästi kalibreeritud klassifikaator on, vaja on teistsuguseid mõõdikuid.

1.5 Kalibreerimise mõõtmine ja meetodid

Selle töö aluseks on võetud Kängsepa 2018. aasta magistritöö [3], mis omakorda toetub Guo 2017. aasta tööle [11]. Selles alampeatükis on kirjeldatud Kängsepa töös valitud kalibreerimise mõõdikuid ning meetodeid, mida käsitletakse ka selles töös.

Naeni jt [12] kohaselt on eeldatav kalibreerimisviga ehk ECE mõõdik, mis hindab mudeli kalibreerimise täpsust, näidates prognoositud tõenäosuste ja tegelike sageduste vahelist erinevust leides kaalutud keskmise täpsuse ja enesekindluse vahel M tulba kohta. Täpsemalt toimub see n ennustuse korral valem (4) järgi:

$$ECE = \sum_{m=1}^M \frac{|B_m|}{n} |acc(B_m) - conf(B_m)|. \quad (4)$$

Valemis (5) on toodud suurima sellise vahe leidmine, mis on ühtlasi ka suurima kalibreerimisvea ehk MCE leidmine:

$$MCE = \max_{m \in \{1, \dots, M\}} |acc(B_m) - conf(B_m)|. \quad (5)$$

Neurovõrkude mõõtmisel kasutatakse tihti negatiivset log-tõepära ehk NLLi (*Negative Log Likelihood*). Täpsemalt hinnatakse mudeli ebakindlust leides prognoositud tõenäosuste ja tegelike tulemuste vahelised erinevused. Olgu y vektor, kus y_i on väärtusega 1 vaid siis, kui klassi i korral on tegemist tõese klassimärgendiga, vastasel juhul on selle väärtus 0. Sel juhul saab NLLi leidmise panna kirja valemiga (6) [2]:

$$NLL = - \sum_{i=1}^k y_i \log(\hat{p}_i(X)) \quad (6)$$

Kull jt 2019. aasta artikli kohaselt [10] on mudelite puhul probleem eelkõige liigse enesekindlusega, mis tähendab, et väljundiks olevas tõenäosusjaotuses on ennustatavate klasside tõenäosused liiga kõrged. Liigne enesekindlus tuleneb üldjuhul treeningandmestikule ülesobitamisest. Enamus kalibreerimismeetodid proovivad seda probleemi lahendada pärast tavalise treeningprotsessi lõppu. Üks levinumaid viise selleks on Guo jt loodud [11] temperatuuri skaleerimine. Temperatuur on reaalarvuline väärtus, mis õpitakse treeningandmestikust eraldi oleva andmestiku peal, minimeerides mõõdikuna NLLi. Seda kasutatakse uute enesekindluste leidmiseks valemiga (7), kus T on temperatuur, z_i on viimaselt kihilt saadud reaalarvuliste väärtuste vektor ja σ_{SM} softmax funktsioon:

$$\hat{c}_i = \max_k \sigma_{SM} \left(\frac{z_i}{T} \right) (k) \quad (7)$$

Temperatuuriga jagatakse läbi mudeli viimaselt kihilt saadud reaalarvulised väärtused enne kui need tõenäosusteks teisendatakse. See ei muuda mudeli täpsust, aga peaaegu garanteeritult parandab kalibreeritust.

Kängsepa töös on käsitletud veel teisigi meetodeid. Beeta kalibreerimisel kasutatakse NLLi mõõdikuna, et leida kolme parameetri väärtused paremaks kalibreerimiseks. Histogrammimeetod ja isotooniline regressioon kasutavad andmepunktide tulpadeks jagamist ja minimeerivad kindla kaofunktsiooni väärtust. Need kolm meetodit jäid aga Kängsepa töös enamuse mõõdikute osas selgelt alla temperatuuri skaleerimisele, kui võrreldi väärtusi samadel mudelitel ja andmestikel, mida ka selles töös kasutatakse. Vaid täpsuse, täpsemalt veamäära ehk ebakorrektsete ennustuse osa juures olid beeta kalibreerimise tulemused paremad. Selles töös võrreldakse järkjärgulise külmutamise tulemusi nende nelja kirjeldatud meetodiga, rõhuasetus on võrdlusel temperatuuri skaleerimise tulemustega.

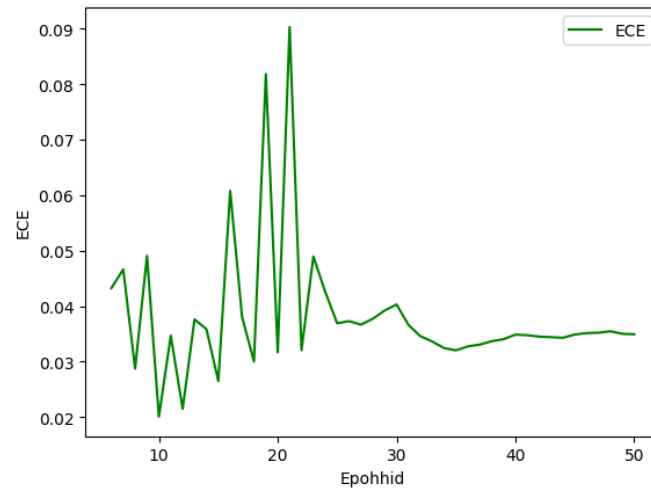
1.6 Järkjärguline külmutamine

Kihi külmutamine tähendab, et selle kihi kaale enam tagasilevi algoritmis ei muudeta, need jäävad samaks kuni treeningprotsessi lõpuni või kuni need pole enam külmutatud. Järkjärgulisel külmutamisel külmutatakse treeningprotsessi käigus mingitel hetkedel valitud kihte. Kuigi seda ja sarnaseid mõisteid on erinevates töödes kasutatud ka muu järjekorraga külmutamisel, siis siin töös mõistetakse seda kui külmutamist, mis hakkab ühel hetkel esimesest kihist ja toimub järjestikku kuni viimase kihini, mis külmutatakse. Kõik külmutatud kihtide kaalud jäävad muutumatuks kuni treeningprotsessi lõpuni. Kirjeldust kõikidest külmutamistest, sealhulgas nende külmutamise hetkedest ja külmutatavate kihtide arvust, nimetatakse siin töös järkjärgulise külmutamise skeemiks.

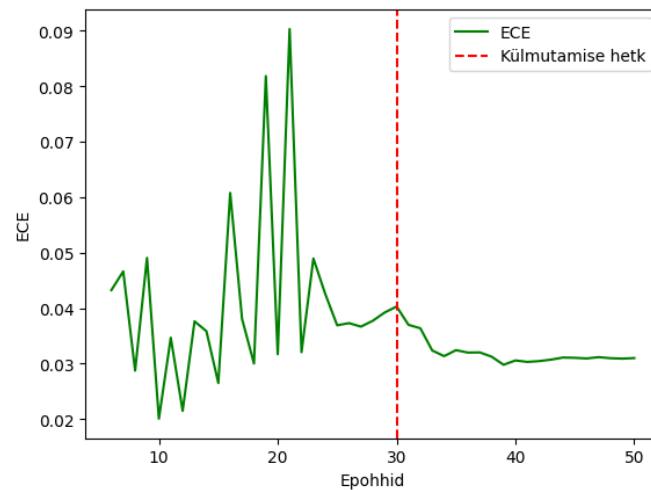
Üldjuhul on seda meetodit kasutatud treeningu kiirendamiseks. Seda on kasutanud näiteks Xiao jt [13] oma uurimuses, kus erinevalt sellest tööst arvutati gradiendi muutumise järgi, millal ning millised kihid külmutada. Uurimuses selgus, et selle tehnikaga on võimalik saavutada populaarse pildiandmestiku CIFAR-10 klassifitseerimisel peaaegu sama täpsus kui ilma külmutamata. Küll aga läks treening selgelt kiiremaks. Romero jt [14] kasutasid siirdeõppe raames tehnikat viisil, kus alustasid hilisemate kihtide treenimisega ja hiljem treeningprotsessis tõid kaasa esimeste kihtide treenimise. Seda tegid nad eeldusega, et varasemate kihtide kaalud väljendavad rohkem üldiseid tunnuseid ja hilisematega saab paremini peenhäälestada mudelit kindla andmestiku jaoks. Seda eeldust saab nende töö põhjal lugeda enam-vähem tõeseks, sest selle meetodiga suutsid mudelid olla sama täpsed kui külmutamata mudeli tulemustega.

Romero jt töös kasutatud eelduse põhjal võiks ka arvata, et järkjärguline külmutamine saab olla efektiivne kalibreerimismeetod. Kui varasemate kihtide optimaalsemate kaalude ni on võimalik jõuda varem kui hilisemate kihtide juures, siis võib nende muutmine kogu treeningprotsessi käigus olla ebavajalik ja isegi viia liigse enesekindluseni. Näide sellest on toodud joonistel 4 ja 5. Joonistel on näidatud ECE muutumine testandmestikul

50 epohhi kehtvas treeningus. Joonis 4 on külmutamiseta treening, joonisel 5 on tulemused külmutamise skeemi korral, kus 30. epohhi juures külmutatakse 80% kihtidest. Täpsemalt on kirjeldatud nende jooniste saamiseks kasutatud mudelit, treeningut ja skeemide leidmist järgmises peatükis.



Joonis 4. Külmutamata treenitud konvolutsioonilise neurovõrgu ECE muutumine.



Joonis 5. Konvolutsioonilise neurovõrgu ECE muutumine pärast 30. epohhi juures külmutamist.

Jooniste võrdluses on näha, et ECE väärtus tõuseb umbes viimase 20 epohhi juures külmutamata treenides. Külmutamise järel aga trend on pigem ühtlane ning ECE väärtus on lõpuks madalam kui 30. ja 35. epohhi juures. See on näide sellest, et järkjärguline külmutamine võiks liigse enesekindluse vastu aidata. Kuna külmutades on tagasilevi algoritmis vähem kaalusid, mida muudetakse, läheb mudel vähem keerukaks ning võib vähem ülesobituda treeningandmestikule, mis ühtlasi võiks ka aidata vähendada liigset enesekindlust. Võrdluse saab ka tuua temperatuuri skaleerimisega, kus sisuliselt külmutatakse kõik peale viimase kihi ning treenitakse seda. Selge erinevus on küll see, et see viimane kiht on üks parameeter ning selle väärtus leitakse teise andmestiku peal, üldjuhul valideerimisandmestikul. Varem pole järkjärgulist külmutamist kalibreerimise eesmärgiga uuritud, järgnevas peatükis kirjeldataksegi, kuidas seda esmakordselt tehakse.

2 Järkjärgulise külmutamise meetodi implementeerimine ja uurimine

Töö eesmärk on leida, kas järkjärgulise külmutamise meetod aitab vähendada konvolutsiooniliste neurovõrkude liigset enesekindlust tuntud pildiandmestikel ning milline on parim viis selle meetodi rakendamiseks. Selleks kirjeldatakse järgnevalt, kuidas järkjärgulise külmutamise meetod implementeeriti ja eksperimendid selle uurimiseks läbi viidi.

2.1 Ülevaade järkjärgulise külmutamise implementeerimisest

Selles alampeatükis on välja toodud järkjärgulise külmutamise implementeerimiseks kasutatud tarkvara, kasutatud andmestikud ning implementeerimise protsess. Eksperimentide kood on kirjas repositooriumis, mille link on toodud lisas 1.

2.1.1 Kasutatud tarkvara

Töös on kasutatud Markus Kängsepa 2018. aasta magistritöö "Calibration of Convolutional Networks" [3] lähtekoodi, sest selles on varasemate artiklite põhjal implementeeritud tuntud konvolutsiooniliste neurovõrkude arhitektuurid, nende treeningprotsess ja lisaks rakendatud ka tuntud kalibreerimismeetodeid. See teeb üsna lihtsaks kõikide nende mudelite loomise ja treenimise ning annab võimaluse võrrelda, kuidas järkjärguline külmutamine tulemusi mõjutab. Kogu lähtekood on kirjutatud Pythoni programmeerimiskeeles ning mudelite loomiseks ja treenimiseks on kasutatud Kerase raamistikku [15]. Lisaks on kasutatud NumPy, pandase ja scikit-learn raamistikke erinevate arvutuste ja masinõppe protseduuride tegemiseks.

Samuti on kasutatud ChatGPT abi nii järkjärgulise külmutamise implementeerimisel kui ka mudelite treenimisel. Üldjuhul anti mudelile sisendiks veateade ning vastuse abil prooviti vigu lahendada või tahtava abimeetodi kirjeldus ning vastuseks saadud koodi abil loodi vastavaid meetodeid. Lisaks on ChatGPTd kasutatud valemite, tabelite ning viidete vormistamiseks.

2.1.2 Kasutatud andmestikud

Nagu eelnevalt kirjeldatud, implementeeriti Kängsepa töös [3] konvolutsioonilisi neurovõrke mõnedel tuntud andmestikel ja rakendati nendele mudelitele mitmeid kalibreerimismeetodeid [3]. Võimalikult heaks võrdluseks nende tulemustega on kasutatud ka siin samu arhitektuure ja andmestikke, mille treeningprotsessi on vaid lisatud järkjärguline külmutamine.

Eelnevas peatükis välja toodud mudelitega leiti täpsus ja kalibreerimistulemused neljal erineval pildiandmestikul. Nendeks on CIFAR-10, CIFAR-100 ja SVHN nimelised andmestikud ning lisaks esimeses eksperimentide etapis FMNISTi andmestik.

CIFAR-10 koosneb 60000 pildist, mis on võrdselt jagatud 10 erinevasse klassi: lennuk, auto, lind, kass, põder, koer, konn, hobune, laev ja veoauto [16]. CIFAR-100 andmestik on sarnane, aga 100 erineva klassiga [16]. SVHN (*Street View House Numbers*) on kogumik värvilisi pilte majanumbritest, kus on vaja klassifitseerida üks keskne number, kuigi numbreid võib pildil veel olla [17]. Kui need andmestikud koosnevad kõik värvilistest piltidest 32x32 dimensioonidega, siis FMNIST (*Fashion-MNIST*) sisaldab must-valgeid 28x28 pilte, mis jagunevad kümnesse klassi riideesemete järgi [18]. Need asjaolud teevad selle andmestiku peal treenimise kiiremaks kui teistel andmestikel, mis oli vajalik esimese etapi eksperimentide tegemiseks.

Andmestikud jaotati treening-, valideerimis- ning testandmestikeks Kängsepa töö kirjelduste järgi [3]. CIFAR-10 ja CIFAR-100 korral oli treeninguks 45000 pilti, valideerimiseks 5000 pilti ning testimiseks 10000 pilti. SVHN andmestikul treenides oli 598388 pilti treeningandmestikus, 6000 pilti valideerimisandmestikus ning 26032 testandmestiku pilti. Lisaks oli FMNISTi kasutades treeninguks 6000 pilti, valideerimiseks 5400 pilti ning testimiseks 10000 pilti.

Andmestike juures tuleb mainida ka andmete eeltöötlust ning rikastamist. SVHN ja FMNIST andmetele tehakse eeltöötlust lahutades iga piksli kolme põhivärvi väärtustest ehk kanalitest selle kanali treeningandmestiku keskmine väärtus ning jagatakse see tulemus samamoodi leitud standardhälvega [6]. Sarnaselt tehakse CIFAR-10 ja CIFAR-100 andmestikel, ainult ResNeti mudeli korral leitakse keskmine ja standardhälve kõikide kanalite peale kokku [11]. Selline eeltöötlus viib kõik väärtused vahemikku, mis sobib konvolutsioonilisele neurovõrgule paremini [1]. Lisaks kasutatakse CIFAR-10 ning CIFAR-100 andmestike juures andmete rikastamist, kus pööratakse osasid pilte ning lisatakse igasse äärde juurde 4 pikslit, mille järel kärbitakse juhuslikult pilte algsesse suurusesse tagasi [11]. See lisab treeningandmestikku rohkem andmeid, mis võiks mudelit treenimisel paremaks teha [1].

2.1.3 Järkjärgulise külmutamise implementeerimine

Kängsepa repositooriumis oli juba olemas mudelite arhitektuur ning nende treenimis- ja hindamisprotsess, sealhulgas kalibreerituse hindamine varem mainitud nelja mõõdikuga [3]. Lisada tuli aga järkjärgulise külmutamise osa treenimisele. Külmutamine oli alati alates esimesest kihist järjest kuni viimase külmutatavani. See tähendas, et treenimisele on ette antud iga mudeli kihi jaoks epohhi arv, mille juures seda enam ei treenita ehk selle kaalud enam ei muutu. Pärast ühe või mitme kihi külmutamist tuleb mudel uuesti kompileerida ja alles siis saab seda edasi treenida kuni järgmise epohhini, kus järgmiste kihtide külmutamine toimub.

Mudelitel oli Kängsepa lähtekoodis iga andmestiku juures määratud õpisammu

väärtuse muutumine epohhi arvust olenevalt. Tavaliselt saab need muutused ette anda Kerase vaikemeetodiga, mis loeb epohhe alates viimasest mudeli kompileerumisest. Kuna aga järkjärgulise külmutamisega oli vaja kompileerida mudelit korduvalt, oli vaja luua uus meetod, mis loeks kõik epohhid kokku ja muudaks õpisammu sellest olenevalt. Võrreldes Kängsepa lähtekoodiga oli vaja ka muuta ja lahendada üksikuid probleeme, mis olid tekkinud nii Kerase kui ka teiste raamistike muutustest pärast algse koodi kirjutamist. Näiteks oli vaja muuta osade meetodite nimesid mudelite loomisel või muuta andmestruktuure, mida ette anda NLLI arvutamiseks. Pärast probleemide lahendamist sai alustada eksperimenteerimisega.

2.2 Järkjärgulise külmutamise eksperimentid

Kui järkjärguline külmutamine oli implementeeritud ning kood korduvalt testitud erinevate mudelite ja andmetike puhul, sai alustada protsessi eksperimentide läbi viimiseks.

Algne plaan oli lokaalsete või muude lihtsalt kättesaadavate ressurssidega teha algseid teste, mida kutsutakse siin töös esimeseks etapiks. Teises etapis oli plaan viia edasi 1 kuni 3 parimat järkjärgulise külmutamise skeemi täisandmetike peal testimiseks. Algupäraste eksperimentide tulemused pole aga välja toodud. Need eksperimentid koosnesid tihti mitmetest erinevatest külmutamisest, mis tähendas, et ei saanud selgeks teha, kuidas mõjutasid erinevad muutused treeningprotsessis tulemust. Edasiselt muudeti veidi lähenemist, selle kirjeldus on toodud järgmises alampeatükis.

2.2.1 Eksperimenteerimise esimese etapi kirjeldus

Esimeses etapis sooritati teste lihtsama ja väiksema andmestiku peal arvutusressursside kokkuhoidmiseks, et leida parimad variandid külmutamiseks, mida hiljem rakendada Kängsepa töös kirjeldatud treeningutele.

Selles etapis tehtud eksperimentid on tehtud FMNISTi andmestikul ResNet SD mudeliga, millel on 152 konvolutsioonikihti. Vaja oli leida andmestik, mida oleks võimalik ühe mudeliga treenida koondumiseni 5 kuni 10 minutiga, et kõik eksperimentid oleks võimalik mõistliku aja sees ära teha. Kängsepa töös kasutatutest oli ehk kiireimaks treeninguks parim SVHN andmestik, sest see sisaldas vaid numbreid [3], kuid ka vaid 5% treeningandmeid kasutades ei olnud soovitud ajavahemik piisav treenimiseks. Seega tuli proovida veel lihtsamat andmestikku, aga ideaalis vältida nii lihtsat andmestikku nagu näiteks vaid käsikirjas numbreid sisaldav MNIST, sest selle peal on võimalik väga häid tulemusi saada üsna lihtsate mudelitega. Valikusse tuli FMNISTi andmestik, mille peal ei saada nii häid täpsuse tulemusi nagu MNISTi peal, aga tegemist on siiski must-valgete 28x28 dimensioonidega piltidega ehk treenimine on kiirem kui kõikidel teiste mainitud andmetike korral, kus kasutuses värvilised 32x32 pildid. Kahjuks ka selle korral sai kasutada vaid 10% treeningandmetest, et 10 minutiga koondumise lähedale jõuda treenides. Kuna kogu see proovimine algas SVHN andmestiku peal, mida treeniti

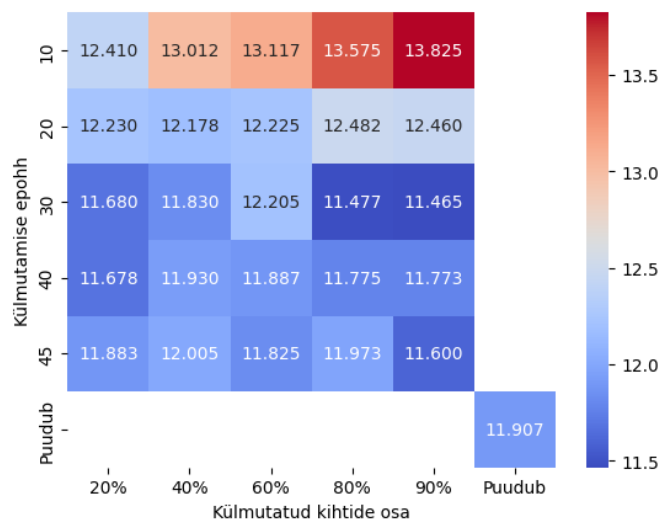
ResNet SD mudeliga, siis sai ka see mudel valitud nendeks eksperimentideks. See on küll otsus, mida võiks muuta tagasi vaadates, sest kuigi oli soov kasutada mudelit, mida kasutatakse ka teises etapis, siis oleks võinud ehk valida väiksema mudeli, mis oleks jõudnud soovitud ajavahemikus koonduda suurema andmestiku korral.

Treening kestis 50 epohhi ning selle käigus muutus õpisamm sarnase tempoga nagu paljude Kängsepa töös tehtud treeningute korral [3], kus õpisammu väärtus langes 10 korda 50% epohhide juures ja veel 10 korda, kui 75% epohhidest oli möödunud. Täpsemalt oli õpisammu väärtus siin esimesed 25 epohhi 0,01, järgmised 15 epohhi 0,001 ning viimased 10 epohhi 0,0001. Tuleb küll märkida, et kõik need väärtused on 10 korda madalamad kui tavaliste täisandmestikega treeningute korral, sest vastasel juhul ei saanud mudel koonduda.

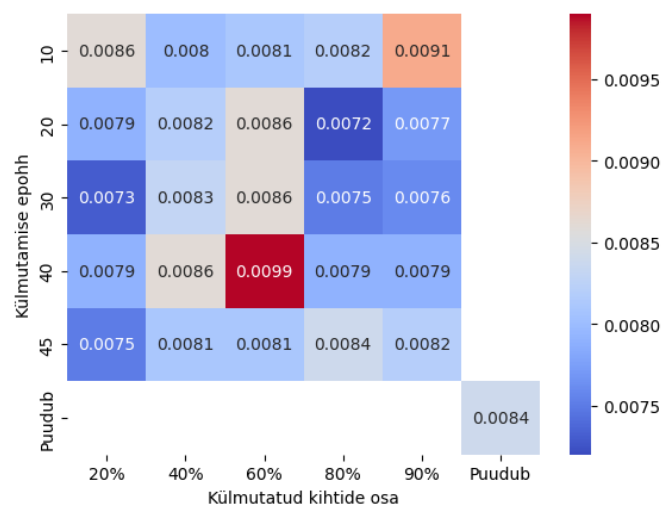
Esimeses eksperimentide etapis valiti 5 eri epohhide arvu, millal külmutada, et võrrelda külmutamist olenevalt epohhide arvust ning külmutatud kihtide arvust. Epohhide arvude valikuteks olid 10, 20, 30, 40 ja 45, see andis võimaluse võrrelda külmutamise hetki kõigi erinevate õpisammu väärtuste juures. Nende epohhide juures külmutati 20%, 40%, 60%, 80% või 90% kihtidest. Igat kombinatsiooni jooksutati 8 korda, et saada iga kombinatsiooni kohta võimalikult usaldatav tulemus. Usaldatavuse suurendamiseks arvestati neist 4 keskmisega ehk eemaldati on 2 parimat ja 2 halvimat tulemust, et juhuslikult väga head või väga kehvad tulemused eemaldada. Lisaks on kõikidele mudelitele rakendatud temperatuuri skaleerimist, sest see oli kalibreerimise jaoks Kängsepa töös parim meetod ning ka siin töös parandas ECE ja NLLi tulemusi [3]. Eelkõige oli eesmärk nii esimeses kui teises eksperimentide etapis võrrelda, kas järkjärguline külmutamine vähendab liigset enesekindlust määral, mida on osade mõõdikute arvestuses näha pärast kalibreerimismeetodi rakendamist. Esimese etapi tulemused on toodud järgmises alampeatükis.

2.2.2 Eksperimenteerimise esimese etapi tulemuste analüüs

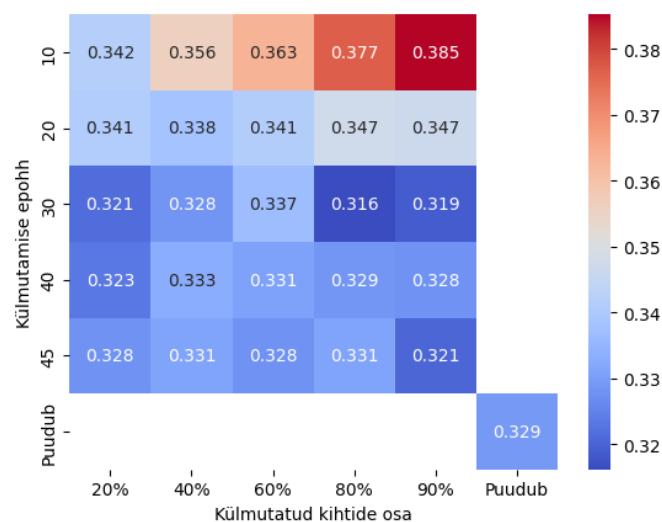
Kombinatsioonide võrdlust veamäära, ECE ja NLLi arvestuses on näha vastavalt joonistel 6, 7 ja 8. Tulemused on leitud pärast temperatuuri skaleerimist. Kõikide näitajate puhul kehtib reegel, et madalam tulemus on parem. Seega on nendel soojuskaartidel (*heatmaps*) tumesinised lahtrid parimate tulemustega ning tumepunased halvimate tulemustega. "Puudub" rea ja veeru nimega ruut tähistab tulemust, kus külmutamist ei rakendatud.



Joonis 6. Soojuskaart kombinatsioonide veamäära võrdlusest.



Joonis 7. Soojuskaart kombinatsioonide ECE võrdlusest.



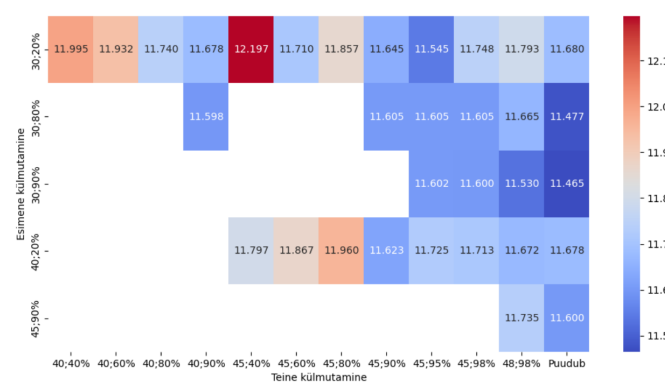
Joonis 8. Soojuskaart kombinatsioonide NLLi võrdlusest.

Pärast temperatuuri skaleerimist olid ECE tulemused väga sarnased. Seega ei olnud mõtet selle mõõdiku järgi siin etapis otsuseid teha. Pigem tuli lähtuda veamäära ja NLLi väärtustest, eriti võrreldes ilma külmutamata saadud tulemusega, mida võib siin kutsuda etaloniks. Nendes arvestustes jäid kõik enne 30. epohhi tehtud külmutamised alla etalonile. Samas enamus hiljem tehtud külmutamistest suutis parandada algset tulemust mõlema mõõdiku osas. Mõõdikutest enam peaks arvestama NLLi, sest antud töö on

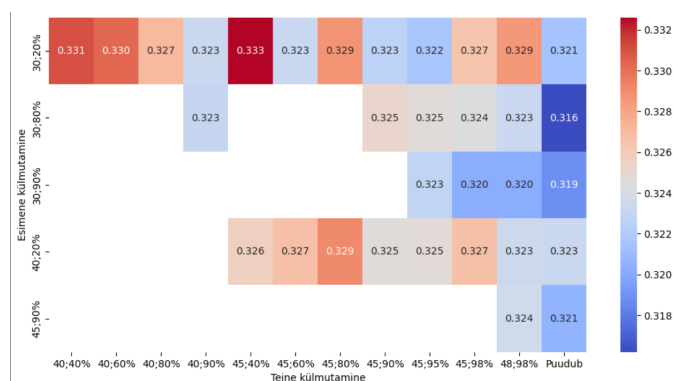
kalibreerimise kohta.

Just viis parimat kombinatsiooni NLLi väärtuse järgi sai valitud järgmisesse alametappi, kus nendele külmutamistele lisati teine külmutamine. Need olid 30. epohhi juures 20% kihtide külmutamine, 30. epohhi juures 80% kihtide külmutamine, 30. juures 90% külmutamine, 40. juures 20% külmutamine ja 45. juures 90% külmutamine. Teise külmutamisena tehti läbi kõik variandid algsetest kombinatsioonidest, mis olid selle esimese külmutamise korral võimalikud ning lisati ka osad hilisemad ja veel rohkemate kihtidega külmutamised, mis olid 45. juures 95% kihtide külmutamine, 45. juures 98% külmutamine, 48. juures 95% külmutamine ja 48. juures 98% külmutamine.

Kui esimesed külmutamised näitasid, et vähemalt selle mudeli ja andmestiku korral on mitmeid kombinatsioone, mis edestavad tavalist treeningut kõikide mõõdikute osas, siis teiste külmutamistega ei leidunud ühtegi kombinatsiooni, mis oleks toonud täpsuse ja NLLi korral edasisi paranemisi nagu joonistel 9 ja 10 on näha. Neist küll 30. juures 90% külmutamine ja hilisemad külmutamised tõid kaasa väga sarnased tulemused esimestele külmutamistele, aga üldjuhul teised külmutamised tõid kaasa veamäära ja NLLi suurenemise.



Joonis 9. Soojuskaart kahe külmutamise kombinatsioonide veamäära võrdlusest.



Joonis 10. Soojuskaart kahe külmutamise kombinatsioonide NLLi võrdlusest.

Seega tuli vastu võtta otsus kolmandaid külmutamisi enam mitte teha ning liikuda edasi teise etappi parimate skeemidega. Valitud sai kaks. Esiteks parima NLLi saavutanud 30. epohhi juures 80% kihtide külmutamine, mis sai tähistuse S1. Teiseks 30. epohhi juures 90% külmutamise ja 48. epohhi juures 98% külmutamise kombinatsioon, sest see saavutas ainult 0,001 kehvema NLLi kui lihtsalt 30. epohhi juures 90% külmutamine, aga erineb esimesest valitud skeemist rohkem ning läheb kokku paremini töö esimeses peatükis mainitud temperatuuri skaleerimise võrdlusega. Teine skeem sai tähistuse S2. Järgnevas alampeatükis on kirjeldatud läbi viidud treeninguid, kus enamus treeningutel on epohhide arv erinev kui 50. Seega on esimeses etapis leitud külmutamise hetked võetud suhtarvuna kogu epohhide arvust ehk 30. epohhi puhul 60% epohhide arvust ja 48. epohhi puhul 96% epohhide arvust.

2.2.3 Eksperimenteerimise teise etapi kirjeldus

Teises etapis treenitakse kõik varem mainitud mudelid vastavalt Kängsepa töö kirjeldustele ja lähtekoodile rakendades juurde kahte valitud järkjärgulise külmutamise skeemi [3]. Selles alampeatükis on kirjeldatud, kuidas need treeningud läbi viidi.

Ressurssidena kasutati siin Tartu Ülikooli High Performance Computing Centre ressursse [19], täpsemalt Nvidia Tesla V100 GPUd ja Kaggle ressursse [20], täpsemalt Nvidia Tesla P100 GPUd. Esimest kasutades esines tõrkeid ja ajalimiite, mis nõudsid meetodi loomist, mis suudaks eelmist treeningut jätkata kohast, kus viimane kaalude fail oli salvestatud.

Kõik treeningud tehti Kängsepa töö [3] lähtekoodi ja kirjelduste järgi. Järgnevas lõikudes on välja toodud ka viited varasematele artiklitele, mille järgi Kängsepa töösse hüperparameetrite väärtused ja implementatsioon valiti. Optimeerimise algoritm oli iga treeningu juures stohhastiline gradientlaskumine, kus Nesterovi inertsitegur oli väärtusega 0,9. Miniploki suurus oli 128, ainus erand oli DenseNeti treenimine. Kõi-

kide treeningute juures leiti tulemused pärast kogu epohhide arvu ja lisaks testiti ka pärast igat epohhi saadud kaalusid valideerimisandmestikul ning lisati parim mudel valideerimisandmestikul tulemuste võrdlusesse.

ResNeti mudel oli 110 konvolutsioonikihti sisaldav ning seda treeniti CIFAR-10 ja CIFAR-100 andmestikel kasutades Li 2017. aasta artikli lähtekoodi [21] ning He jt 2015. aasta artikli kirjeldusi [5]. Treening kestis kokku 200 epohhi ehk esimese skeemi esimene külmutamine toimus pärast 120. epohhi. Õpisamm oli esimesed 80 epohhi 0,1, järgmised 70 epohhi 0,01 ning lõpus 0,001.

ResNet SD mudelit treeniti Huangi jt 2016. aasta artikli põhjal [6], implementatsioon pärineb Cheni 2016. aasta koodist [22]. CIFAR-10 andmestikul treenides ei hakanud see mudel alguses koonduma. Võimalus oleks olnud alustada väiksema õpisammuga, kuid see oleks võrdluse saamiseks nõudnud ka ilma külmutamata treenimist ja testimist, milline õpisamm toimiks. Seega toimus treening vaid CIFAR-100 ja SVHN andmestikel. CIFAR-100 andmestiku mudel oli 110-konvolutsioonikihiline ning treening kestis kokku 500 epohhi, õpisamm oli väärtusega 0,1 250 epohhi, 0,01 kuni 375. epohhini ning viimased 125 epohhi 0,001. SVHN andmestikul treenimine oli kokku vaid 50 epohhi, 30. epohhi juures langes õpisamm 0,1 pealt 0,01 peale ning 35. epohhi juures läks see veel 10 korda väiksemaks. Lisaks oli esimesed 2 epohhi õpisamm 0,04, et koondumine toimuks.

Wide ResNet oli siin 32 konvolutsioonikihiga neurovõrk, mille implementatsioon pärines Majumdari 2018. aasta lähtekoodist [23]. Hüperparameetrite väärtused treeninguks tulid Zagoruyko and Komodakise 2016. aasta artiklist [7]. Treening kestis CIFAR-10 ja CIFAR-100 andmestikel 200 epohhi. Õpisamm algas 0,1 peal ning langes 20 korda nii 60., 120. ja 160. epohhi juures.

DenseNeti mudel koosnes 40 konvolutsioonikihist, selle implementatsioon on võetud Majumdari 2018. aasta [24] ning Yu 2017. aasta repositooriumitest [21]. Treeningud toimusid CIFAR-10 ja CIFAR-100 andmestikel Huangi jt 2018. aasta töös kirjeldatud hüperparameetrite väärtuste põhjal [8]. Need olid ainsad treeningud, kus miniploki suurus oli 64. Õpisamm oli 0,1 esimesed 150 epohhi, 0,01 järgmised 75 ning 0,001 viimased 75 epohhi.

3 Eksperimenteerimise teise etapi tulemuste analüüs

Siin peatükis analüüsitakse teise etapi tulemusi veamäära, NLLi ja ECE osas. Kasutuses on sarnased tabelid nagu Kängsepa töös ning järkjärgulise külmutamise skeemide tulemusi on võrreldud Kängsepa töö tulemustega [3]. Nende võrdluste tabelites on toodud mudeli nimi, andmestik, ilma külmutamata tulemus ehk Algne ning seejärel viie kalibreerimismeetodi tulemused. Histo tähistab histogrammimeetodit, Iso isotoonilist regressiooni, Temp temperatuuri skaleerimist ning Beeta beeta kalibreerimist.

Nende järel on 4 veergu selle töö tulemusi. S1 tähistab esimese skeemi ehk 60% epohhide juures 80% kihtide külmutamise lõplike kaalude tulemust ehk mudelit, mis on saadud pärast kõikide epohhide möödumist. S2 tähistab teise skeemi ehk 60% epohhide juures 90% kihtide ning 96% epohhide juures 98% kihtide külmutamise lõplike kaalude tulemust. S1* tähistab parimat leitud mudelit valideerimisandmestiku põhjal esimese skeemiga treenides ning S2* parimat leitud mudelit valideerimisandmestiku põhjal teise skeemi korral. Kõik tulemused on leitud pärast temperatuuri skaleerimist, sest see oli Kängsepa töös kõige edukam kalibreerimismeetod enamus mõõdikute arvestuses [3]. Kalibreerimata tulemuste võrdlused NLLi ja ECE arvestuses on toodud lisades 2 ja 3. Iga tulemuse kohta on toodud selle järk vaadeldava treeningu kohta alaindeksina ning lõpuks leitud keskmine järk iga meetodi kohta. Lisaks nende kolme mõõdiku võrdlustele on peatükis toodud ka järkjärgulise külmutamise ajalise mõju analüüs.

3.1 Veamäära võrdlus

Tabelis 1 on toodud veamäära võrdlus. See on mõõdik, mille väärtust temperatuuri skaleerimine ei muuda. See oli ka mõõdik, kus beeta kalibreerimine suutis temperatuuri skaleerimist edestada, aga teised meetodid muutsid täpsuse kalibreerimata mudelitest kehvemaks [3].

Tabel 1. Veamäära (%) võrdlus mudelite ja andmestike kaupa.

Mudel	Andmestik	Algne	Histo	Iso	Temp	Beeta	S1	S2	S1*	S2*
DenseNet 40	CIFAR-10	7,58 ₄	7,93 ₈	7,65 ₇	7,58 ₄	7,59 ₆	7,57 ₃	7,99 ₉	7,49 ₂	7,37₁
DenseNet 40	CIFAR-100	30,0 ₃	32,49 ₇	30,22 ₆	30,0 ₃	29,81₁	30,02 ₅	29,87 ₂	41,0 ₉	36,88 ₈
ResNet 110	CIFAR-10	6,44 ₂	6,59 ₅	6,36₁	6,44 ₂	6,44 ₂	7,15 ₈	6,93 ₆	7,17 ₉	7,0 ₇
ResNet 110	CIFAR-100	28,52 ₂	31,26 ₉	29,31 ₄	28,52 ₂	28,36₁	29,98 ₆	30,32 ₈	29,91 ₅	30,31 ₇
ResNet SD 110	CIFAR-100	27,17 ₆	29,74 ₉	27,47 ₈	27,17 ₆	26,32 ₅	24,89₁	25,54 ₃	25,4 ₂	26,28 ₄
ResNet SD 152	SVHN	1,85 ₂	2,07 ₅	1,96 ₄	1,85 ₂	1,82₁	2,3 ₆	2,77 ₈	2,38 ₇	2,77 ₈
Wide ResNet 32	CIFAR-10	6,07 ₅	6,18 ₇	5,98 ₃	6,07 ₅	5,94 ₂	5,93₁	6,04 ₄	9,09 ₉	6,51 ₈
Wide ResNet 32	CIFAR-100	26,18 ₄	28,89 ₇	26,33 ₆	26,18 ₄	25,98 ₃	25,72₁	25,83 ₂	32,8 ₉	32,48 ₈
Keskmine järk		3,5 ₂	7,12 ₉	4,88 ₅	3,5 ₂	2,62₁	3,88 ₄	5,25 ₆	6,5 ₈	6,38 ₇

Mõlema skeemi lõplikute kaaludega mudelid on üsna selgelt paremad kui valideerimisandmestiku parimad mudelid, kuigi Kängsepa töös leiti nii DenseNeti kui ResNet SD

tulemused just valideerimisandmestiku parimate mudelite järgi. Väljaspool paari erandit oli S1* ja S2* tulemused ühed viimastest kõikide treeningute korral. Esines ka väga kehvasid tulemusi nende korral, näiteks DenseNet 40 mudelit CIFAR-100 andmestikul treenides olid S1* ja S2* tulemused külmutamata variandist vastavalt üle 10% ja 6% kehvemad.

S1 ja S2 keskmise leitud järgu põhjal temperatuuri skaleerimise ning beeta kalibreerimise tulemust ei edestanud. S1 suutis küll mõne treeningu juures kõiki teisi meetodeid edestada, kuid näiteks ResNet 110 mudeli osas oli üks kehvemaid. Sarnast kõikumist esines ka S2 korral, kuid üldkokkuvõttes oli see teine skeem esimesest veidi kehvem. Erinevused polnud aga kunagi väga suured ning jäid pea kõikide võrdluste osas alla 2%, tihti olid need veel väiksemad. Ühe erandina ResNet SD 110 mudeli CIFAR-100 treenimisel suutis esimene skeem kõiki Kängsepa töös leitud tulemusi üle 2% edestada.

3.2 NLLi võrdlus

NLL oli esimeses etapis mõõdik, mille järgi lõplikud otsused tehti. Selle tähtsus ei vähene ka selles etapis ehk selle tulemusi saab tõsiselt arvestada järelduste tegemisel. Kängsepa töös tulid parimad tulemused üsna selgelt temperatuuri skaleerimisega, ainult beeta kalibreerimine jõudis lähedale temperatuuri skaleerimise tulemustele [3]. Tabelis 2 on näha võrdlused selle töö tulemustega.

Tabel 2. NLLi võrdlus mudelite ja andmestike kaupa.

Mudel	Andmestik	Algne	Histo	Iso	Temp	Beeta	S1	S2	S1*	S2*
DenseNet 40	CIFAR-10	0,4282 ₈	0,5725 ₉	0,2773 ₇	0,2251 ₃	0,2392 ₆	0,2249 ₂	0,2391 ₅	0,2259 ₄	0,2187₁
DenseNet 40	CIFAR-100	2,0174 ₈	4,1829 ₉	1,6491 ₇	1,0571₁	1,1532 ₄	1,0705 ₂	1,0705 ₂	1,4647 ₆	1,3186 ₅
ResNet 110	CIFAR-10	0,3583 ₈	0,5472 ₉	0,2708 ₇	0,2093 ₂	0,2139 ₄	0,226 ₆	0,2101 ₃	0,2249 ₅	0,2079₁
ResNet 110	CIFAR-100	1,6937 ₇	4,2139 ₉	1,8926 ₈	1,0917 ₃	1,1318 ₄	1,1389 ₅	1,165 ₆	1,0621₁	1,0794 ₂
ResNet SD 110	CIFAR-100	1,3525 ₇	4,0187 ₉	1,6371 ₈	0,9421 ₅	0,9643 ₆	0,8745₁	0,8922 ₃	0,891 ₂	0,9108 ₄
ResNet SD 152	SVHN	0,0854 ₃	0,2887 ₉	0,1093 ₆	0,0786₁	0,0796 ₂	0,0993 ₄	0,1124 ₇	0,1006 ₅	0,1126 ₈
Wide ResNet 32	CIFAR-10	0,3817 ₈	0,5132 ₉	0,2327 ₆	0,1915 ₃	0,202 ₅	0,1854 ₂	0,1822₁	0,2653 ₇	0,192 ₄
Wide ResNet 32	CIFAR-100	1,8022 ₈	3,9352 ₉	1,5607 ₇	0,9445 ₂	1,0386 ₄	0,9466 ₃	0,9423₁	1,1557 ₆	1,1311 ₅
Keskmine järk		7,12 ₈	9,0 ₉	7,0 ₇	2,5₁	4,38 ₅	3,12 ₂	3,5 ₃	4,5 ₆	3,75 ₄

Esimese eksperimentide etapi tulemused näitasid, et pärast temperatuuri skaleerimist on NLLi tulemused selgelt veamääraga korreleeruvad. S1* ja S2* keskmised järgud olid paremad kui veamäär osas, kuid jäid lõplike kaalude tulemustele üldjuhul alla.

NLLi arvestuses on S1 ja S2 tulemused paremad ilma külmutamata temperatuuri skaleerimisest täpselt nende treeningute osas, kus tulemused olid ka paremad veamäär osas. Sarnaselt eelmisele alampeatükile on vahed üsna väikesed kahe skeemi lõplike kaalude tulemuste ning Kängsepa töö temperatuuri skaleerimise tulemuste vahel. NLLi võrdluses on aga kõikidel teistel kalibreerimismeetoditel selgelt kehvemad tulemused, sest järkjärgulise külmutamise mudelitele on rakendatud temperatuuri skaleerimist, seega on S1 ja S2 keskmised järgud madalamad.

3.3 ECE võrdlus

ECE numbrid on tabelis 3 toodud protsendina. Selle arvestuses oli Kängsepa töös selgelt parim meetod temperatuuri skaleerimine. Teised meetodid jäid selgelt alla.

Tabel 3. ECE (%) võrdlus mudelite ja andmestike kaupa.

Mudel	Andmestik	Algne	Histo	Iso	Temp	Beeta	S1	S2	S1*	S2*
DenseNet 40	CIFAR-10	5,5 ₉	2,13 ₈	1,68 ₆	0,95 ₅	1,7 ₇	0,93 ₄	0,63 ₁	0,87 ₃	0,63 ₁
DenseNet 40	CIFAR-100	21,16 ₉	11,97 ₈	5,25 ₆	0,9 ₁	6,03 ₇	0,91 ₂	0,91 ₂	1,34 ₅	1,09 ₄
ResNet 110	CIFAR-10	4,75 ₉	1,25 ₆	1,47 ₈	1,13 ₃	1,42 ₇	1,21 ₅	0,75 ₂	1,19 ₄	0,61 ₁
ResNet 110	CIFAR-100	18,48 ₉	9,06 ₈	6,54 ₇	2,38 ₄	4,6 ₆	1,77 ₃	2,55 ₅	1,62 ₁	1,7 ₂
ResNet SD 110	CIFAR-100	15,86 ₉	7,6 ₈	5,21 ₇	1,21 ₄	3,55 ₆	1,15 ₃	0,96 ₁	1,55 ₅	0,98 ₂
ResNet SD 152	SVHN	0,86 ₉	0,53 ₃	0,25 ₁	0,61 ₄	0,5 ₂	0,81 ₈	0,69 ₅	0,76 ₇	0,7 ₆
Wide ResNet 32	CIFAR-10	4,51 ₉	1,01 ₇	1,19 ₈	0,8 ₄	0,97 ₆	0,58 ₂	0,52 ₁	0,85 ₅	0,59 ₃
Wide ResNet 32	CIFAR-100	18,78 ₉	7,64 ₈	5,82 ₆	0,78 ₁	7,01 ₇	1,6 ₃	1,98 ₅	1,11 ₂	1,77 ₄
Keskmine järk		9,0 ₉	7,0 ₈	6,12 ₇	3,25 ₃	6,0 ₆	3,75 ₄	2,75 ₁	4,0 ₅	2,88 ₂

ECE arvestustes suutis S2* saavutada parema keskmise järgu kui S1, kuid jällegi olid mõlema skeemi korral lõplike kaalude keskmised järgud paremad kui valideerimisandmestiku prima tulemuse järgi valides.

S2 oli ECE võrdluses enamus treeningute korral parem kui tavaline temperatuuri skaleerimine, eriti selgelt paranes seda skeemi kasutades kalibreeritus ECE arvestuses CIFAR-10 andmestikul. S1 tulemused olid vaid kahe treeningu puhul paremad kui S2 omad, aga olid külmutamata temperatuuri skaleerimisest neljal puhul madalama ECE väärtusega.

3.4 Ajaline võrdlus

Tabelites 4 ja 5 on toodud vastavalt esimese ja teise skeemi ajalised võrdlused enne ja pärast külmutamisi, kõik ajad on sekundites. Esimeses veerus on keskmine epohhi ajaline pikkus enne külmutamist ning teises veerus summaarne aeg selle põhjal. Järgnevates veergudes on keskmised epohhide pikkused pärast külmutamisi, esimese skeemi puhul oli külmutamisi üks ning teise skeemi korral kaks. Eelviimases veerus on tegelik summaarne aeg. Lõpuks on toodud kui palju kogu aeg vähenes protsentides. Siin tuleb küll mainida, et arvutus on tehtud puhtalt epohhide aja järgi, välja on jäetud protsessid enne treenimist, epohhide vahel toimuvad protsessid nagu valideerimisandmestikul ennustamine ja failide salvestamine ning arvutused pärast treeningut. Nende sisse arvestamisel suureneksid mõlemad koguajad üsna sarnasel määral ehk vähenemise protsendid ilmselt veidi langetaksid. Kuna aga esines tõrkeid nagu näiteks internetiühenduse peatumine või koguni treeningu katkemine, siis tundus kõige mõistlikum võrrelda puhtalt kaua epohhid kokku aega võtsid.

Tabel 4. Esimese skeemi ajaline võrdlus sekundites

Mudel	Andmestik	Epohhi pikkus	Eeldatav koguaeg	Epohh pärast külmutamist	Tegelik koguaeg	Koguaja kiirenemine
DenseNet 40	CIFAR-10	40,17	12051	31,33	10990,2	8,8%
DenseNet 40	CIFAR-100	48,56	14568	29,04	12225,6	16,08%
ResNet 110	CIFAR-10	58,9	11780	38,85	10284	12,7%
ResNet 110	CIFAR-100	58,57	11714	38,56	10113,2	13,67%
ResNet SD 110	CIFAR-100	37,16	18580	31,56	17460	6,03%
ResNet SD 152	SVHN	454,65	22732,5	185,3	17345,5	23,7%
Wide ResNet 32	CIFAR-10	126,96	25392	55,58	19636,8	22,67%
Wide ResNet 32	CIFAR-100	126,96	25392	55,7	19691,2	22,45%

Tabel 5. Teise skeemi ajaline võrdlus sekundites

Mudel	Andmestik	Epohhi pikkus	Eeldatav koguaeg	Epohh pärast külmutamist	Epohh pärast teist külmutamist	Tegelik koguaeg	Koguaja kiirenemine
DenseNet 40	CIFAR-10	39,82	11946	30,38	29,68	10804,8	9,55%
DenseNet 40	CIFAR-100	48,45	14535	28,77	28,5	12170,16	16,27%
ResNet 110	CIFAR-10	44,02	8804	29,4	27,92	7622,56	13,42%
ResNet 110	CIFAR-100	35,9	7180	28,18	27,08	6553,6	8,72%
ResNet SD 110	CIFAR-100	44,13	22065	29,6	28,82	19143,4	13,24%
ResNet SD 152	SVHN	449,85	22492,5	163,15	147,03	16726,26	25,64%
Wide ResNet 32	CIFAR-10	126,51	25302	50,08	50,14	19154,88	24,29%
Wide ResNet 32	CIFAR-100	127,38	25476	49,78	45,99	19237,68	24,49%

Eeldada võis, et külmutamised kiirendavad treeningut, sest kaalusid muudetakse vähem. Tulemustest on näha, et nii on. Ainus kord, kui kiirenemine jäi alla 8%, oli ResNet SD 110 CIFAR-100 andmestikul treenimisel esimese skeemiga. Kõige enam aega võtnud mudelite osas jõudis kiirenemine mõlema skeemi korral üle 20%. Teine skeem külmutab rohkem kihte ja on seega kiirem enamuse treeningute osas. Kõik vähenemised on piisavalt suured, et kiirendada selgelt treenimist, isegi kui arvestada sisse järkjärgulise külmutamise protsessid.

3.5 Järeldused

Võrreldes mudeleid, mida on treenitud järkjärgulise külmutamise kahe valitud skeemiga, täpsemalt võrreldes tulemusi pärast temperatuuri skaleerimise rakendamist, ei saa väita, et järkjärguline külmutamine kindlasti aitab vähendada konvolutsiooniliste neurovõrkude liigset enesekindlust. Teine skeem oli küll keskmise järgu võrdluses ECED mõõtes veidi parem, kuid jäi täpsuses ja seetõttu ka NLLi arvestuses alla ilma külmutamata treenimisele. Esimene skeem sisaldas vähemate kihtide külmutamist sama epohhi juures ning seega üsna loogiliselt viis teisest skeemist parema täpsuseni, aga selle puhul ei saa väita, et kalibreeritus paranes. Kõige kindlamalt saab iga kasutatud mõõdiku juures väita, et järkjärgulise külmutamise meetodi osas on treeningu lõplikud kaalud paremad kui parimad leitud kaalud valideerimisandmestiku järgi.

See, et järkjärguline külmutamine selles töös ei toonud selgeid paranemisi nii täpsuse kui ka kalibreerimise osas, ei tähenda, et järkjärgulist külmutamist konvolutsiooniliste neurovõrkude osas ei võiks edasi uurida. Esiteks tuleb mainida, et selles töös olid piiratud ressursid ja aeg, et leida paremad skeemid. Kui uurida selgelt rohkem kombinatsioone, eelistatavalt rohkemate mudelite ja andmestike peal kui esimeses etapis tehtud eksperi-

mendid ühe andmestiku ja mudeliga, on ilmselt võimalik leida paremaid kombinatsioone kui siin töös leiti. Saaks ka teha paremaid valikuid mudeli osas, ResNet SD 152 polnud ilmselt parim valik andmestikku ja ressursse arvestades esimeses etapis. Saaks ka proovida teisi lähenemisi skeemide leidmiseks. Võiks kaaluda näiteks Xiao jt [13] meetodi kasutamist, mis suutis CIFAR-10 peal mitme erineva mudeliga saavutada pea sama või veidi parema täpsuse kui tavaline treenimine.

Teiseks on siin töös ikkagi leitud, et üsna sarnaseid tulemusi on võimalik saavutada nii kalibreerituse kui täpsuse osas võrreldes külmutamata variandiga vähendades samal ajal 6-25% treenimiseks kuluvat aega. See tähendab, et rakendades järkjärgulist külmutamist, saab sama aja ja ressursidega treenida suuremaid ja keerulisemaid mudeleid kui algselt plaanis. Kahe skeemi lõplike kaalude tulemuste peale leidis vaid üks kord, kus veamäär oli rohkem kui 1,5% suurem kui tavalise treeningu tulemus ja ka siis oli vahe 1,8%. NLLi tulemused olid veel lähedamal ning ECE pigem paranes teise skeemiga. Sellised muutused ei ole väga negatiivsed ning saab väita, et päris tihti on külmutamisega saadud aja säästmise väärtuslikum, eriti kui tänu sellele on võimalik treenida suurem ja parem mudel.

Kokkuvõte

Töös kirjeldati konvolutsiooniliste neurovõrkude tööd, kalibreerimist, järkjärgulist külmutamist ning põhjuseid selle kasutamiseks kalibreerimismeetodina. Eksperimente tehti kahes etapis. Esimeses leiti ühe mudeli ja andmestiku põhjal kaks parimat järkjärgulise külmutamise skeemi. Neid rakendati teises etapis mitmete mudelite treenigutele ning tulemusi võrreldi Kängsepa 2018. aasta töö tulemustega.

Eesmärgiks oli leida, kas järkjärgulise külmutamise meetod aitab konvolutsiooniliste neurovõrkude liigset enesekindlust vähendada tuntud pildiandmestikel ning milline on parim viis selle meetodi rakendamiseks. Esimeses etapis leiti kaks skeemi. Esimene ehk S1 sisaldas 60% epohhide juures 80% kihtide külmutamist. Teine ehk S2 koosnes 60% epohhide juures 90% kihtide külmutamise ja 96% epohhide juures 98% kihtide külmutamisest. Nendele rakendati veel temperatuuri skaleerimist. Kuigi S2 suutis kalibreerimismõõdiku ECE arvestuses edestada tavalist temperatuuri skaleerimist, oli see kehvem NLLi ja veamäära arvestuses. S1 oli viimase kahe mõõdiku arvestusest teisest skeemist parem, kuid jäi siiski üldjuhul tavalisele temperatuuri skaleerimisele alla. Samas olid mõlemad skeemid üsna selgelt kiiremad tavalisest treeningust.

Edasiarendusena saaks proovida veel erinevaid skeeme, võrreldes neid rohkem kui ühe mudeli puhul. Samuti näitavad tulemused, et järkjärgulist külmutamist rakendades saab treeningu aega selgelt vähendada, mis annab võimaluse kasutada suuremaid mudeleid, mis võiksid täpsuse osas olla paremad.

Viidatud kirjandus

- [1] Goodfellow, I., Bengio, Y. ja Courville, A. Deep Learning. Cambridge, Massachusetts: The MIT Press, 2016.
- [2] Filho, T., Song, H., Nieto, M. P., Rodriguez, R. S., Kull, M. ja Flach, P. Classifier Calibration: How to assess and improve predicted class probabilities: a survey (2021). (2023-05-09). <https://arxiv.org/abs/2112.10327>.
- [3] Kängsepp, M. Calibration of Convolutional Networks. TÜ arvutiteaduse instituudi magistritöö, 2018.
- [4] Arusoo, E. Tehisnärvivõrkude kirjeldamine väikese maailma omaduse abil. Tallinna Tehnikaülikooli tervisetehnoloogiate instituudi magistritöö, 2020.
- [5] He, K., Zhang, X., Ren, S. ja Sun, J. Deep residual learning for image recognition (2015). (2023-05-09). <https://arxiv.org/abs/1512.03385>.
- [6] Huang, G., Sun, Y., Liu, Z., Sedra, D. ja Weinberger, K. Deep networks with stochastic depth (2016). (2023-05-09). <https://arxiv.org/abs/1603.09382>.
- [7] Zagoruyko, S. ja Komodakis, N. Wide residual networks (2015). (2023-05-09). <https://arxiv.org/abs/1605.07146>.
- [8] Huang, G., Liu, Z., Maaten, L. van der ja Weinberger, K. Q. Densely connected convolutional networks (2018). (2023-05-09). <https://arxiv.org/abs/1608.06993>.
- [9] Flach, P. A. Machine learning: the art and science of algorithms that make sense of data. Cambridge ; New York: Cambridge University Press, 2012.
- [10] Kull, M., Nieto, M. P., Kängsepp, M., Filho, T., Song, H. ja Flach, P. Beyond temperature scaling: Obtaining well-calibrated multiclass probabilities with Dirichlet calibration (2019). (2023-05-09). <https://arxiv.org/abs/1910.12656>.
- [11] Guo, C., Pleiss, G., Sun, Y. ja Weinberger, K. Q. On calibration of modern neural networks (2017). (2023-05-09). <https://arxiv.org/abs/1706.04599>.
- [12] Naeini, M. P., Cooper, G. ja Hauskrecht, M. Obtaining Well Calibrated Probabilities Using Bayesian Binning. *Proceedings of the AAAI Conference on Artificial Intelligence* 29.1 (veebruar 2015), lk. 2901–2907.
- [13] Xiao, X., Mudiyansele, T. B., Ji, C., Hu, J. ja Pan, Y. Fast Deep Learning Training through Intelligently Freezing Layers. *2019 International Conference on Internet of Things (iThings) and IEEE Green Computing and Communications (GreenCom) and IEEE Cyber, Physical and Social Computing (CPSCom) and IEEE Smart Data (SmartData)* (2019), lk. 1225–1232.

- [14] Romero, M., Interian, Y., Solberg, T. ja Valdes, G. Targeted transfer learning to improve performance in small medical physics datasets. *Medical Physics* 47.12 (detsember 2020). arXiv:1912.06761, lk. 6246–6256.
- [15] Chollet, F. Keras. <https://keras.io>. 2015.
- [16] Krizhevsky, A. Learning Multiple Layers of Features from Tiny Images (2012).
- [17] Netzer, Y, Wang, T, Coates, A, Bissacco, A, Wu, B ja Ng, A Y. Reading Digits in Natural Images with Unsupervised Feature Learning (2011). http://ufldl.stanford.edu/housenumbers/nips2011_housenumbers.pdf.
- [18] Xiao, M, Rasul, K ja Vollgraf, R. Fashion-MNIST: a Novel Image Dataset for Benchmarking Machine Learning Algorithms. <https://arxiv.org/abs/1708.07747>. (2023-05-09). 2017.
- [19] HPC. <https://hpc.ut.ee/>. (2023-05-09).
- [20] Kaggle. <https://www.kaggle.com/>. (2023-05-09).
- [21] Li, W. cifar-10-cnn: Play deep learning with cifar datasets. <https://github.com/BIGBALLON/cifar-10-cnn>. (2023-05-09). 2017.
- [22] Chen, L. Implementation of stochastic depth networks in keras. <https://github.com/transcranial/stochastic-depth>. (2023-05-09). 2016.
- [23] Majumdar, S. Wide residual networks in keras. <https://github.com/titu1994/Wide-Residual-Networks>. (2023-05-09). 2018.
- [24] Majumdar, S. Densenet implementation in keras. <https://github.com/titu1994/DenseNet>. (2023-05-09). 2018.

Lisad

Lisa 1. Repositoorium

Töös eksperimentide jooksutamiseks kirjutatud lähtekood on leitav aadressil
<https://github.com/OliverSavolainen/freezing-calibration>.

Lisa 2. NLLi võrdlus mudelite ja andmestike kaupa kasutades järkjärguliste skeemide kalibreerimata tulemusi.

Mudel	Andmestik	Algne	Histo	Iso	Temp	Beeta	S1	S2	S1*	S2*
DenseNet 40	CIFAR-10	0,4282 ₇	0,5725 ₉	0,2773 ₃	0,2251 ₁	0,2392 ₂	0,4238 ₆	0,445 ₈	0,3832 ₅	0,3568 ₄
DenseNet 40	CIFAR-100	2,0174 ₇	4,1829 ₉	1,6491 ₅	1,0571 ₁	1,1532 ₂	2,0277 ₈	1,9961 ₆	1,6168 ₄	1,4557 ₃
ResNet 110	CIFAR-10	0,3583 ₇	0,5472 ₉	0,2708 ₃	0,2093 ₁	0,2139 ₂	0,3587 ₈	0,3202 ₅	0,3463 ₆	0,3065 ₄
ResNet 110	CIFAR-100	1,6937 ₆	4,2139 ₉	1,8926 ₈	1,0917 ₁	1,1318 ₂	1,6912 ₅	1,7315 ₇	1,2053 ₃	1,248 ₄
ResNet 110 (SD)	CIFAR-100	1,3525 ₇	4,0187 ₉	1,6371 ₈	0,9421 ₁	0,9643 ₂	1,2317 ₅	1,252 ₆	1,1423 ₃	1,1787 ₄
ResNet 152 (SD)	SVHN	0,0854 ₃	0,2887 ₉	0,1093 ₆	0,0786 ₁	0,0796 ₂	0,1031 ₄	0,118 ₇	0,1063 ₅	0,1206 ₈
Wide ResNet 32	CIFAR-10	0,3817 ₈	0,5132 ₉	0,2327 ₃	0,1915 ₁	0,202 ₂	0,3621 ₇	0,3408 ₆	0,306 ₄	0,3091 ₅
Wide ResNet 32	CIFAR-100	1,8022 ₈	3,9352 ₉	1,5607 ₅	0,9445 ₁	1,0386 ₂	1,8008 ₇	1,7764 ₆	1,3175 ₄	1,2165 ₃
Keskmine järk		6,62 ₈	9,0 ₉	5,12 ₅	1,0 ₁	2,0 ₂	6,25 ₆	6,38 ₇	4,25 ₃	4,38 ₄

Lisa 3. ECE (%) võrdlus mudelite ja andmestike kaupa kasutades järkjärguliste skeemide kalibreerimata tulemusi

Model	Dataset	Uncal	Histo	Iso	Temp	Beeta	S1	S2	S1*	S2*
DenseNet 40	CIFAR-10	5.5 ₇	2.13 ₄	1.68 ₂	0.95 ₁	1.7 ₃	5.58 ₈	5.8 ₉	5.17 ₆	5.02 ₅
DenseNet 40	CIFAR-100	21.16 ₉	11.97 ₅	5.25 ₂	0.9 ₁	6.03 ₃	20.82 ₇	21.11 ₈	13.33 ₆	11.86 ₄
ResNet 110	CIFAR-10	4.75 ₆	1.25 ₂	1.47 ₄	1.13 ₁	1.42 ₃	5.05 ₉	4.83 ₇	4.88 ₈	4.69 ₅
ResNet 110	CIFAR-100	18.48 ₇	9.06 ₄	6.54 ₃	2.38 ₁	4.6 ₂	18.58 ₈	18.65 ₉	11.45 ₅	12.37 ₆
ResNet 110 (SD)	CIFAR-100	15.86 ₉	7.6 ₄	5.21 ₃	1.21 ₁	3.55 ₂	13.91 ₇	14.26 ₈	12.43 ₅	13.25 ₆
ResNet 152 (SD)	SVHN	0.86 ₅	0.53 ₃	0.25 ₁	0.61 ₄	0.5 ₂	1.66 ₆	1.95 ₈	1.94 ₇	2.33 ₉
Wide ResNet 32	CIFAR-10	4.51 ₉	1.01 ₃	1.19 ₄	0.8 ₁	0.97 ₂	4.36 ₆	4.46 ₇	3.99 ₅	4.49 ₈
Wide ResNet 32	CIFAR-100	18.78 ₉	7.64 ₄	5.82 ₂	0.78 ₁	7.01 ₃	18.38 ₈	18.22 ₇	12.64 ₆	9.76 ₅
Keskmine järk		7.62 ₈	3.62 ₄	2.62 ₃	1.38 ₁	2.5 ₂	7.38 ₇	7.88 ₉	6.0 ₅	6.0 ₅

II. Litsents

Lihtlitsents lõputöö reprodutseerimiseks ja üldsusele kättesaadavaks tegemiseks

Mina, **Oliver Savolainen**,

1. annan Tartu Ülikoolile tasuta loa (lihtlitsentsi) minu loodud teose
Konvolutsiooniliste neurovõrkude kalibreerimine järkjärgulise külmutamise meetodiga,
mille juhendaja on Meelis Kull,
reprodutseerimiseks eesmärgiga seda säilitada, sealhulgas lisada digitaalarhiivi DSpace kuni autoriõiguse kehtivuse lõppemiseni.
2. Annan Tartu Ülikoolile loa teha punktis 1 nimetatud teos üldsusele kättesaadavaks Tartu Ülikooli veebikeskkonna, sealhulgas digitaalarhiivi DSpace kaudu Creative Commons'i litsentsiga CC BY NC ND 3,0, mis lubab autorile viidates teost reprodutseerida, levitada ja üldsusele suunata ning keelab luua tuletatud teost ja kasutada teost ärieesmärgil, kuni autoriõiguse kehtivuse lõppemiseni.
3. Olen teadlik, et punktides 1 ja 2 nimetatud õigused jäävad alles ka autorile.
4. Kinnitan, et lihtlitsentsi andmisega ei riku ma teiste isikute intellektuaalomandi ega isikuandmete kaitse õigusaktidest tulenevaid õigusi.

Oliver Savolainen

09.05.2023