

Deep Learning Szeminárium

5. előadás: Mély hálózatok

Csiszárík Adrián

MTA, Rényi Alfréd Matematikai Kutatóintézet

- A 2000-es évekig 1-2 rejtett rétegű, sekély hálózatokkal kísérleteztek
- Több réteg növeli a reprezentáció rugalmasságát
 - Kevesebb neuron is elég lehet
 - Hierarchikus adatrepresentációk tudnak létrejönni
 - Gazdagabb reprezentációs struktúra, jobb általánosítás
- Sajnos mélyebb hálózatokat nehezebb tanítani!

Instabil gradienssek problémája

- Vanishing gradients: bemenet felé egyre kisebb gradienssek
- Exploding gradients: bemenet felé egyre nagyobb gradienssek

A korai rétegek gradiensze egy hosszú szorzat eredménye $W_i^T \sigma'(z_i)$ alakú tényezőkkel

Ahhoz, hogy a különböző rétegek gradienszei hasonló nagyságrendűek legyenek, valamilyen mechanizmusnak gondosan ki kell(ene) egyensúlyoznia ezeket a tényezőket.

Instabil gradiensek problémája

Sok trükk segíthet ennek kezelésében:

- Más aktivációs függvények keresése: ReLU
- Súlyok és eltolások gondos inicializálása
- Több tanító adat, augmentáció
- Különféle SGD variánsok átméretezhetik a gradienst
- Konvolúciós hálók: kevesebb paraméter, gazdagabb gradiens jel
- **Batch Normalization**: stabil eloszlás a rétegek bemenetein
- **Reziduális hálózatok**

Batch Normalization

- A rétegek bemenetének eloszlása változik a tanulás során
- Törékenyebb és lassabb a tanulási folyamat
- A rétegeknek folyamatosan adaptálódniuk kell a megváltozott adatkörnyezethez
- Ötlet:
 - Normalizáljuk a rétegek bemeneteit
 - Mini batch alapján, komponensenként normalizálunk
 - A normalizált eloszlás paraméterei tanulhatóak

Batch Normalization

- Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift
Sergey Ioffe, Szegedy Krisztián (2015)
<https://arxiv.org/abs/1502.03167>
- Google Scholar szerint 4020 9122 hivatkozás

Input: Values of x over a mini-batch: $\mathcal{B} = \{x_{1\dots m}\}$; Parameters to be learned: γ, β
Output: $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$
$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$
$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$
$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$
$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$

Algorithm 1: Batch Normalizing Transform, applied to activation x over a mini-batch.

Deep Residual Learning for Image Recognition

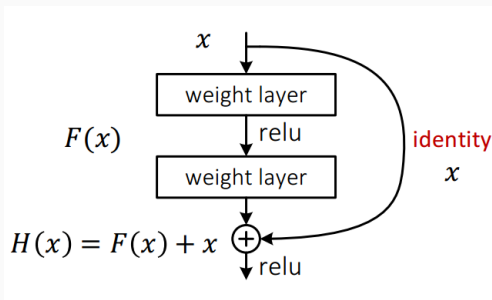
He et al.

<https://arxiv.org/abs/1502.03167>

- 19236 hivatkozás
- Degradation problem: Egy hálózathoz új rétegeket adva egy ideig nő a pontosság, de gyorsan telítődik és aztán elkezdi csökkenni.
- Pedig a hálózat reprezentációs képessége nő!
- Hiszen az identitás függvényt kellene megtanulni, hogy az új réteg ne ártson.

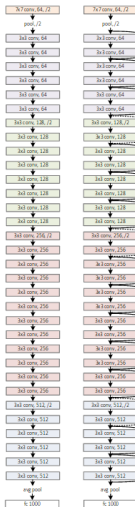
Reziduális hálózatok

- Nem könnyű SGD-vel identitás függvényt tanulni.
- Lényesen egyszerűbb elérni, hogy a hálózat konstans nulla függvényt tanuljon: csupa nulla súly.
- Ha az optimális leképezés az identitás közelében van: könnyebb megtalálni a konstans nullához viszonyítva.



Reziduális hálózatok

plain net



ResNet



er)

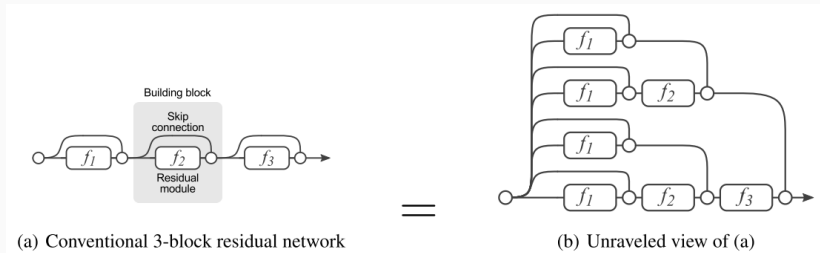
- Nem jön létre új reprezentációs tér
- A gradiens exponenciálisan sok úton terjed
- A közvetlen kapcsolatokon keresztül erős gradiens tud áramlani

Residual Networks Behave Like Ensembles of Relatively Shallow Networks

Andreas Veit, Michael Wilber, Serge Belongie

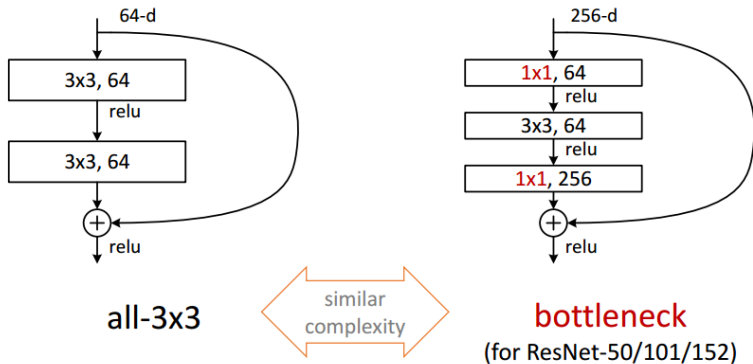
<https://arxiv.org/abs/1605.06431>

- Tekinhetünk egy reziduális hálóra úgy, mint sok különböző mélységű hálózat együttesére



Reziduális hálózatok

- Egy praktikus megfontolás: paraméterszám csökkentése
- Bottleneck reziduális blokk



- Akár 1000 felett is tudnak javulást hozni új rétegek.
- Egy idő után nagyon kicsit növekedés nagyon sok paraméterrel
- A háló figyelmen kívül tudja hagyni a reziduális blokkokat, ahelyett, hogy hasznosítaná őket.
- A szélesség növelésével (a mélység rovására), ugyanannyi paraméterrel jobb eredményt tudunk elérni
- Wide Residual Networks
- Densenets

Wide Residual Networks

Sergey Zagoruyko, Nikos Komodakis

<https://arxiv.org/abs/1605.07146>

- Egyszerű implementálni
- https://github.com/zsoltzombori/keras_fashion_mnist_tutorial/blob/master/fashion_mnist_resnet.py

An overview of gradient descent optimization algorithms

Sebastian Ruder

<http://ruder.io/optimizing-gradient-descent/>

Visualizing the Loss Landscape of Neural Nets

Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, Tom Goldstein
NIPS 2018

<https://arxiv.org/abs/1712.09913>

Cyclical Learning Rates for Training Neural Networks

Leslie N. Smith

WACV 2017

<https://arxiv.org/abs/1506.01186>

On the importance of initialization and momentum in deep learning

Sutskever et al.

ICML 2013

<http://proceedings.mlr.press/v28/sutskever13.pdf>