# Optimization

## Deep Learning KU (DAT.C302UF)

---

Simon Hitzginger

Nov 12, 2025

Institute of Machine Learning and Neural Computation
Graz University of Technology, Austria

# OPTIMIZERS

- As soon as we have a gradient w.r.t. a mini-batch

$$\mathbf{g} = \nabla_{\boldsymbol{\theta}} \mathcal{L}_{\mathcal{B}}(\boldsymbol{\theta})$$
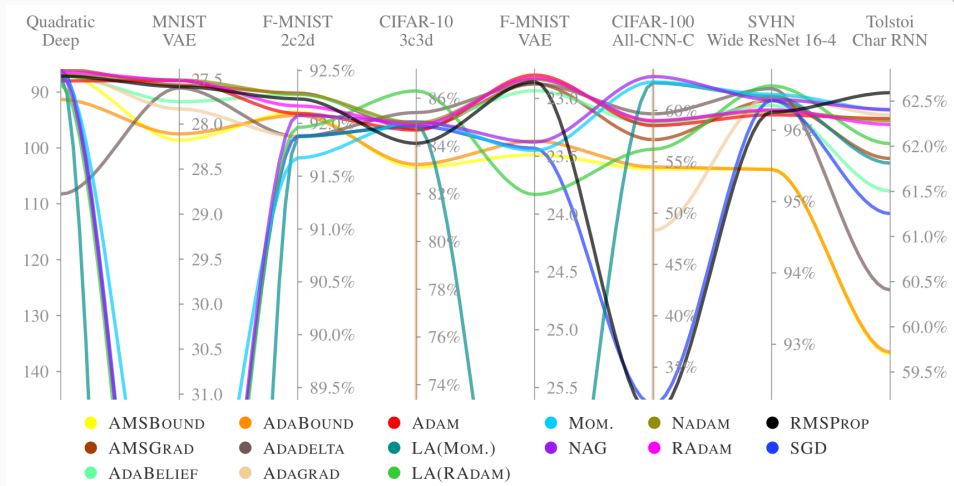
  we can pass it to an **optimizer**

- The optimizer uses the gradient to **update our current estimate** of $\boldsymbol{\theta}$

- For example, if the optimizer is **Stochastic Gradient Descent** (SGD), it implements the update rule
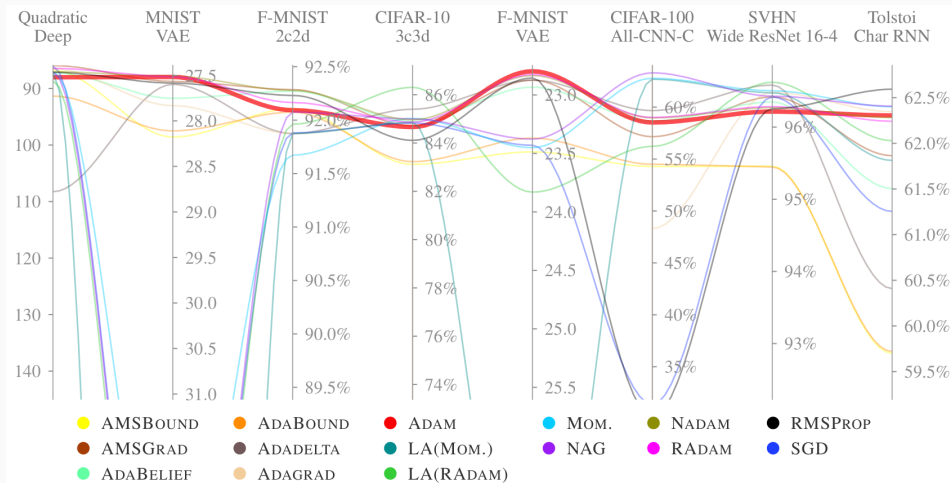
$$\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \mathbf{g}$$

  where $\eta > 0$ is the **learning rate**.

[Bosch et al., 2022]

| Name | Name | Name |
|---|---|---|
| AccelGrad | C-ADAM | PAGE |
| ACClip | CADA | PAL |
| AdaAlter | Cool Momentum | PolyAdam |
| AdaBatch | CProp | Polyak |
| AdaBayes/AdaBayes-SS | Curveball | PowerSGD/PowerSGDM |
| AdaBelief | Dadam | Probabilistic Polyak |
| AdaBlock | DeepMemory | ProbLS |
| AdaBound | DGNOpt | PStorm |
| AdaComp | DiffGrad | QHAdam/QHM |
| AdaDelta | EAdam | RAdam |
| Adafactor | EKFAC | Ranger |
| AdaFix | Eve | RangerLars |
| AdaFom | Expectigrad | RMSProp |
| AdaFTRL | FastAdaBelief | RMSterov |
| Adagrad | FRSGD | S-SGD |
| ADAHESSIAN | G-AdaGrad | SAdam |
| Adai | GADAM | Sadam/SAMSGrad |
| AdaLoss | Gadam | SALR |
| Adam | GOALS | SAM |
| Adam* | GOLS-I | SC-Adagrad/SC-RMSProp |
| AdamAL | Grad-Avg | SDProp |
| AdaMax | GRAPES | SGD |
| AdamNC | Gravilon | SGD-BB |
| AdaMod | Gravity | SGD-G2 |
| AdamP/SGDP | HAdam | SGDEM |
| AdamT | HyperAdam | SGDHess |
| AdamW | K-BFGS/K-BFGS(L) | SGDM |
| AdamX | KF-QN-CNN | SGDR |
| ADAS | KFAC | SHAdagrad |
| AdaS | KFLR/KFRA | Shampoo |
| AdaScale | L4Adam/L4Momentum | SignAdam++ |
| AdaSGD | LAMB | SignSGD |
| AdaShift | LaProp | SKQN/S4QN |
| AdaSqrt | LARS | SM3 |
| Adathm | LHOPT | SMG |
| AdaX/AdaX-W | LookAhead | SNCM |
| AEGD | M-SVAG | SoftAdam |
| ALI-G | MADGRAD | SRSGD |
| AMSBound | MAS | Step-Tuned SGD |
| AMSGrad | MEKA | SWATS |
| AngularGrad | MTAdam | SWNTS |
| ArmijoLS | MVRC-1/MVRC-2 | TAdam |
| ARSG | Nadam | TEKFAC |
| ASAM | NAMSB/NAMSG | VAdam |
| AutoLRS | ND-Adam | VR-SGD |
| AvaGrad | Nero | vSGD-b/vSGD-g/vSGD-l |
| BAdam | Nesterov | vSGD-fd |
| BGAdam | Noisy Adam/Noisy K-FAC | WNGrad |
| BPGrad | NosAdam | YellowFin |
| BRMSProp | Novograd | Yogi |
| BSGD | NT-SGD | |
| | Padam | |

- **Many** optimizers to choose from…

- Is there a **single best** general-purpose optimizer?

2

[Schmidt et al., 2021]



AMSBOUND · AdaBound · Adam · Mom. · Nadam · RMSPROP
AMSGRAD · Adadelta · LA(MOM.) · NAG · RADAM · SGD
AdaBelief · Adagrad · LA(RADAM)

3

[Schmidt et al., 2021]

- ADAM is a **good default choice**
  - Main hyperparameter: **learning rate**

- ADAM is a **good default choice**
  - Main hyperparameter: **learning rate**

- However, benchmarking optimizers is **hard...**
  - Choice of **hyperparameters**
  - Performance specific to problems

# DEMO: OPTIMIZERS IN `PyTorch`

📄 Bosch, N., Grosse, J., Hennig, P., Kristiadi, A., Pförtner, M., Schmidt, J., Schneider, F., Tatzel, L., and Wenger, J. (2022).
**Numerics of machine learning.**
Technical report, Tübingen AI Center.

📄 Schmidt, R. M., Schneider, F., and Hennig, P. (2021).
**Descending through a crowded valley - benchmarking deep learning optimizers.**
In Meila, M. and Zhang, T., editors, *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*, volume 139 of *Proceedings of Machine Learning Research*, pages 9367–9376. PMLR.