



Deep Learning

Lecture 6: Large-scale Training for Deep Neural Networks

University of Agder,
Kristiansand, Norway

Prepared by: Hadi Ghauch^{*}, Hossein S. Ghadikolaei[†]

^{*} Telecom ParisTech

[†] Royal Institute of Technology, KTH

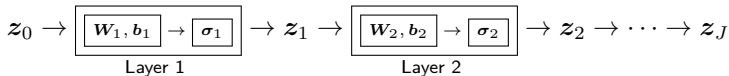
<https://sites.google.com/view/fundl/home>

Outline

1. Challenges for DNN optimization
2. Large-Scale methods for DNN training
3. Regularization in DNNs
4. Some tricks of the trade

Recap: Deep Neural Networks (DNNs)

DNN with J -layers: cascade of J filters + activation



In compact form: $z_j = \sigma_j(W_j z_{j-1} + b_j)$

$z_j \in \mathbb{R}^{d_j}$: input to layer j

$W_j \in \mathbb{R}^{d_j \times d_{j-1}}$: **Weights** for layer j (drop bias for simplicity)

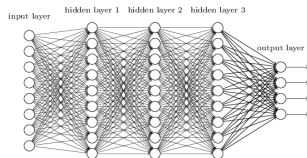
$\sigma_j : \mathbb{R}^{d_j} \mapsto \mathbb{R}^{d_j}$: **Activation func** for layer j

$[\sigma_j(u)]_i = \sigma_j([u]_i)$, $u \in \mathbb{R}^d$, $\forall i \in [N]$

Recap: DNN training and non-convex optimization

The loss for sample $(\mathbf{x}_i, \mathbf{y}_i)$:

$$\|\mathbf{y}_i - \sigma_J(\mathbf{W}_J \cdots \sigma_2(\mathbf{W}_2 \sigma_1(\mathbf{W}_1 \mathbf{x}_i)))\|_2^2, i \in [N]$$



Let $\mathbf{w} \in \mathbb{R}^d$ be the total number of weights in DNN (from all layers)
 $d \geq 10^6$ in current deep learning application

The resulting DNN training problem:

$$\min_{\mathbf{w} \in \mathcal{D}} \frac{1}{N} \sum_{i=1}^N f_i((\mathbf{x}_i, \mathbf{y}_i); \mathbf{w}) + \lambda \|\mathbf{w}\|_2^2 := f(\mathbf{w})$$

non-convex optimization problem

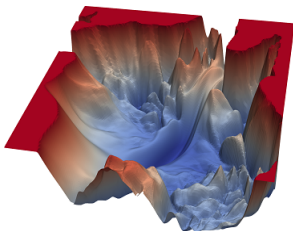
Outline

1. Challenges for DNN optimization
2. Large-Scale methods for DNN training
3. Regularization in DNNs
4. Some tricks of the trade

Loss Surface of DNNs

DNN training in fully-parametric form:

$$\arg \min_{\mathbf{w}} f(\mathbf{w}) := \sum_{i=1}^N f_i(\mathbf{w}) \quad (1)$$

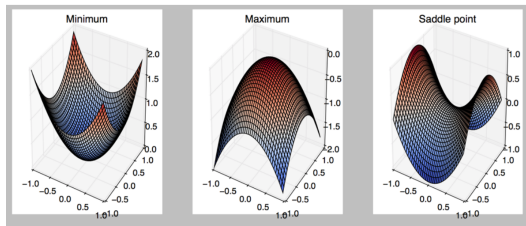


- $\mathbf{w} \in \mathbb{R}^d$: **all weights in DNN** ($d \gg 1$)
- $f()$: **regularized loss function**, $f_i()$: loss func for sample $i \in [N]$
- $f()$: non-convex for general DNN training problems
- $\mathbf{g}^{(k)} := \nabla_{\mathbf{w}} f(\mathbf{w}^{(k)}) \in \mathbb{R}^d$: gradient at $\mathbf{w}^{(k)}$
- $\mathbf{H}^{(k)} := \nabla_{\mathbf{w}}^2 f(\mathbf{w}^{(k)}) \in \mathbb{R}^{d \times d}$: Hessian at $\mathbf{w}^{(k)}$
- $\{ \lambda_i[\mathbf{H}^{(k)}] \}_{i=1}^d$: eigenvalues of $\mathbf{H}^{(k)}$

Loss Surface of DNNs

Only **first-order methods** used in training
due to complexity of optimization

- first-order methods converge to a **stationary point**
- Recall definition of **stationary point**: w is stationary $\Leftrightarrow \nabla f(w) = 0$
- Problematic in non-convex setting



All positive eigenvalues

All negative eigenvalues

Some positive
and some negative

- Particularly problematic for DNN training
Because loss surface exponentially many saddle points

Loss Surface of DNNs

How many saddle points or local minima ?

- Eigenvalues of \mathbf{H} are Bernoulli distributed:
- $\lambda_i[\mathbf{H}^{(k)}] > 0$ with prob p (independently of other ones).
- local min \Leftrightarrow all eigenvalues of Hessian are positive.
Probability of having a local min: $\propto p^d$ (where $d \gg 1$)
- local minima **exponentially less likely** to occur than saddle points
- argument based on random matrix theory

Almost all stationary pts are saddle

- Result validated empirically

However, most minima give a 'good enough' loss value

Analysis of Gradient Descent

Improvement in f with one GD step? $w^{(k+1)} = w^{(k)} - \alpha^{(k)} g^{(k)}$

- intractable in closed form: use 2nd order Taylor approx at $w^{(k)}$

$$f(w^{(k)} - \alpha^{(k)} g^{(k)}) \approx f(w^{(k)}) - \alpha^{(k)} \|g^{(k)}\|_2^2 + \frac{(\alpha^{(k)})^2}{2} (g^{(k)})^T H^{(k)} g^{(k)}$$

Use Taylor approx to find **decay after one GD step**:

$$\Delta^{(k)} := f(w^{(k+1)}) - f(w^{(k)}) \approx \underbrace{-\alpha^{(k)} \|g^{(k)}\|_2^2}_{\text{decay from grad. norm}} + \underbrace{\frac{(\alpha^{(k)})^2}{2} (g^{(k)})^T H^{(k)} g^{(k)}}_{\text{depends on alignment of } g^{(k)}, H^{(k)}}$$

Conditions on $\Delta^{(k)}$: $\Delta^{(k)} \leq 0$ and $|\Delta^{(k)}|$ as **large as possible**. why ?

- balance b/w term **decays cost** (depends on $\|g^{(k)}\|_2^2$), and **increase or decay** (depend on $\|g^{(k)}\|_2^2$ and $\lambda_{\min}[H^{(k)}]$)
- one way to satisfy conditions on $\Delta^{(k)}$ is pick **stepsize optimally**:
 $\alpha^{(k)*} = \|g^{(k)}\|_2^2 / (g^{(k)})^T H^{(k)} g^{(k)}$
not used in practice (need Hessian): sequence $\alpha^{(k)}$ selected in advance

III-conditioning

III-Conditioning: when curvature of f not even across dim

- large difference b/w eigenvalues of $\mathbf{H}^{(k)}$
- sensitivity of loss fn different across features

Different perspective on ill-conditioning for **non-convex** f

Mathematically ill-conditioning happens:

$$\underbrace{\frac{(\alpha^{(k)})^2}{2} (\mathbf{g}^{(k)})^T \mathbf{H}^{(k)} \mathbf{g}^{(k)}}_{\text{increase/decay from alignment of } \mathbf{g}^{(k)}, \mathbf{H}^{(k)}} \geq \underbrace{\alpha^{(k)} \|\mathbf{g}^{(k)}\|_2^2}_{\text{decay from grad. norm}}$$

recall $f(\mathbf{w})$ non-convex: $\lambda_i[\mathbf{H}^{(k)}]$ may take any value

$\Rightarrow \Delta^{(k)} \geq 0 \Rightarrow f(\mathbf{w}^{(k)}) \leq f(\mathbf{w}^{(k+1)})$: GD step **increases the cost**
ill-conditioning more harmful in non-convex setting

Preconditioning

Mitigate ill-conditioning ?

most DNN training methods find a **direction**, $d^{(k)}$, to move the solution
then they **scale each coordinate** by a single factor, $\alpha^{(k)}$

$$\begin{bmatrix} w_1^{(k+1)} \\ w_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} w_1^{(k)} \\ w_2^{(k)} \end{bmatrix} - \alpha^{(k)} \begin{bmatrix} d_1^{(k)} \\ d_2^{(k)} \end{bmatrix}$$

true for methods seen so far: GD, SGD, momentum, etc.

Preconditioning: scale coordinates in direction differently

$$\begin{bmatrix} w_1^{(k+1)} \\ w_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} w_1^{(k)} \\ w_2^{(k)} \end{bmatrix} - \alpha^{(k)} \begin{bmatrix} p_1^{(k)} & p_2^{(k)} \\ p_3^{(k)} & p_4^{(k)} \end{bmatrix} \begin{bmatrix} d_1^{(k)} \\ d_2^{(k)} \end{bmatrix}$$

$$\begin{bmatrix} w_1^{(k+1)} \\ w_2^{(k+1)} \end{bmatrix} = \begin{bmatrix} w_1^{(k)} \\ w_2^{(k)} \end{bmatrix} - \alpha^{(k)} \begin{bmatrix} p_1^{(k)} \\ p_3^{(k)} \end{bmatrix} \circ \begin{bmatrix} d_1^{(k)} \\ d_2^{(k)} \end{bmatrix}$$

adapt stepsize using more info than gradient (w/out Hessian)

Outline

1. Challenges for DNN optimization
2. Large-Scale methods for DNN training
3. Regularization in DNNs
4. Some tricks of the trade

Stochastic Gradient Descent (SGD)

Recall that computing full-gradient (all samples) is too expensive

SGD: Do gradient step for one training sample

- Pick i uniformly at random from $[N]$
- Compute SG w.r.t. sample i : $\hat{\mathbf{g}}_i^{(k)} := \nabla f_i(\mathbf{w}^{(k)})$
- Update weights: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \hat{\mathbf{g}}_i^{(k)}$

where $\nabla f_i()$ gradient of f_i , k iteration index, $\alpha^{(k)}$ stepsize

Pros: low computational footprint

Cons: estimating full-gradient ∇f , using gradient of sample i , ∇f_i .
Too noisy!

Mini-batch SGD

Solution to noisy gradient estimate in SGD

Mini-batch SGD: gradient over batch of training samples

- Pick batch of samples, \mathcal{B} , uniformly at random from $[N]$
- Compute SG w.r.t. batch \mathcal{B} : $\hat{\mathbf{g}}_i^{(k)} := \sum_{i \in \mathcal{B}} \nabla f_i(\mathbf{w}^{(k)})$
- Update weights: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \hat{\mathbf{g}}_i^{(k)}$

Improves the gradient estimate wrt SGD

$\mathcal{B} = i$ recovers SGD, and $\mathcal{B} = [N]$ recovers full GD

Pros: Changing \mathcal{B} trades-off accuracy of gradient estimation for complexity

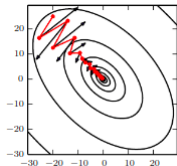
Momentum Methods

SGD with momentum

- Pick i uniformly at random from $[N]$
- Compute SG w.r.t. sample i : $\hat{\mathbf{g}}_i^{(k)} := \nabla f_i(\mathbf{w}^{(k)})$
- Update momentum: $\mathbf{t}_i^{(k+1)} = \eta \hat{\mathbf{g}}_i^{(k)}, \quad 0 < \eta < 1$
- Update weights: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \alpha^{(k)} \hat{\mathbf{g}}_i^{(k)} - \mathbf{t}_i^{(k+1)}$

Faster conv (wrt SGD) when loss has

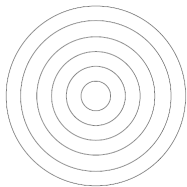
- high curvature and $\hat{\mathbf{g}}$ too noisy
- norm of $\hat{\mathbf{g}}$ vary across surface
- momentum 'smoothes-out' variation in $\hat{\mathbf{g}}$
- similar behavior as Nesterov and heavy ball methods (lec 2)



Choice of Stepsize

Stepsize, $\alpha^{(k)}$, aka, **learning rate** in backpropagation

- chosen apriori **independently of gradient/Hessian**
- examples: $\alpha^{(k)} = c/\sqrt{k}$, $\alpha^{(k)} = c/k$, $\alpha^{(k)} = c (c \ll 1)$
- Problematic when **curvature is uneven**



$\lambda_1 = \lambda_2$ (ideal)



$\lambda_1 \gg \lambda_2$ (bad for GD)

- Stepsize should be adapted to compensate for uneven curvature, **w/out using Hessian**:
idea behind **Adaptive rate methods**

Adaptive Gradient (ADAGRAD)

Adapt stepsize using norm of all previous gradients

Intuition: Shrink stepsize for sensitive features (high grad norm)

ADAGRAD:

- Pick i uniformly at random from $[N]$
- Compute SG w.r.t. sample i : $\hat{\mathbf{g}}_i^{(k)} := \nabla f_i(\mathbf{w}^{(k)})$
- Compute all previous norms of SG: $\mathbf{u}_i^{(k+1)} = \mathbf{u}_i^{(k)} + \hat{\mathbf{g}}_i^{(k)} \circ \hat{\mathbf{g}}_i^{(k)}$
// $\mathbf{u}_i^{(k)}$ is running a sum of all previous squared grad norms
- Update weights: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \left((\mathbf{u}_i^{(k)})^{1/2} + \delta \mathbf{1} \right)^{-1} \circ \hat{\mathbf{g}}_i^{(k)}$
// $\left((\mathbf{u}_i^{(k)})^{1/2} + \delta \mathbf{1} \right)^{-1}$: behaves like 'stepsize' (not strictly though)

$(-)^{1/2}$ and $(-)^{-1}$ are element-wise operations.

Stepsize inversely prop to **gradient norm**. why ?

weight update step is form of **preconditioning**. why?

Limitations: Stepsize decays too much due to scaling with sum of all prev grad : add exponential decay to gradients \Rightarrow **RMSProp**

RMSPProp

Intuition: add exponentially decaying sum of grad norm

RMSPProp:

- Pick i uniformly at random from $[N]$
- Compute SG w.r.t. sample i : $\hat{\mathbf{g}}_i^{(k)} := \nabla f_i(\mathbf{w}^{(k)})$
- Exp decaying sum of prev SGs: $\mathbf{u}_i^{(k+1)} = \beta \mathbf{u}_i^{(k)} + (1 - \beta) \hat{\mathbf{g}}_i^{(k)} \circ \hat{\mathbf{g}}_i^{(k)}$
// $\mathbf{u}_i^{(k)}$ is exponentially decaying sum of all previous squared grad norms
- Update weights: $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \left((\mathbf{u}_i^{(k)})^{1/2} + \delta \mathbf{1} \right)^{-1} \circ \hat{\mathbf{g}}_i^{(k)}$

$(-)^{1/2}$ and $(-)^{-1}$ are element-wise operations.

Difference with ADAGRAD: exponentially decreasing sum for all the gradients.

put **less emphasis on past values** of gradient norm

weight update step is a form of **preconditioning**

add **bias correction** for first/second moments \Rightarrow ADAM

Adaptive Moments (ADAM)

Intuition: $\hat{\mathbf{g}}_i^{(k)}$ may be 'bad' estimator of true gradient
approximate bias for FoM/SoM of $\hat{\mathbf{g}}_i^{(k)}$, and **compensate** for it

ADAM:

- pick i uniformly at random from $[N]$
- compute SG w.r.t. sample i : $\hat{\mathbf{g}}_i^{(k)} := \nabla f_i(\mathbf{w}^{(k)})$
- **approximate bias for FoM:** $\mathbf{u}_i^{(k+1)} = \beta_1 \mathbf{u}_i^{(k)} + (1 - \beta_1) \hat{\mathbf{g}}_i^{(k)}$
// exponentially decaying sum of prev SGD
- **approximate bias for SoM:** $\mathbf{v}_i^{(k+1)} = \beta_2 \mathbf{v}_i^{(k)} + (1 - \beta_2) \hat{\mathbf{g}}_i^{(k)} \circ \hat{\mathbf{g}}_i^{(k)}$
// exponentially decaying sum of prev SGD norms
- **compute bias term of FoM:** $\hat{\mathbf{u}}_i^{(k+1)} = \mathbf{u}_i^{(k+1)} / (1 - \beta_1^k)$
- **compute bias term of SoM:** $\hat{\mathbf{v}}_i^{(k+1)} = \mathbf{v}_i^{(k+1)} / (1 - \beta_2^k)$
- **update weights:** $\mathbf{w}^{(k+1)} = \mathbf{w}^{(k)} - \left((\hat{\mathbf{v}}_i^{(k)})^{1/2} + \delta \mathbf{1} \right)^{-1} \circ \hat{\mathbf{u}}_i^{(k)}$

Difference? move along **updated FoM bias**, $\hat{\mathbf{u}}_i^{(k)}$, not the SG.
precondition using updated SoM bias $\hat{\mathbf{v}}_i^{(k)}$.

Outline

1. Challenges for DNN optimization
2. Large-Scale methods for DNN training
3. Regularization in DNNs
4. Some tricks of the trade

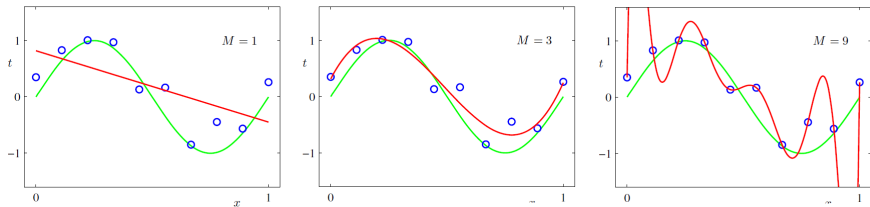
Regularization in DNNs

Regularization is critical in DDN:

- huge number of parameters: 1 Billion weights in modern DNN
- need to avoid *overfitting*

Reduce **test error** at the expense of higher **training error**

- trading increased bias for reduced variance of the estimator (Lec 2)
- reduce the model complexity to match the data-generating process



source: C. Bishop, 2006

Regularization: Parameter norm penalty

Regularization via **Parameter norm penalty**

- $\min_{\mathbf{w}} \tilde{f}(\mathbf{w}) := f(\mathbf{w}) + \lambda \|\mathbf{w}\|_p, p \in \{1, 2\}$
- $f()$: loss function (from ERM problem)
- $\tilde{f}()$: regularized loss function
- no regularization needed for biases, \mathbf{b}_j

In deep learning:

- common to start with **overparametrized model** (lots of weight)
- and gradually increase **regularization parameter**, λ
- **hyperparameter tuning** done using **cross-validation**. how?
grid search over set of possible values for λ
pick value for λ that minimizes **test error (over test set)**
see also K -fold cross validation

Regularization: Parameter norm penalty

Weight Decay: Tikhonov regularization
($p = 2$)

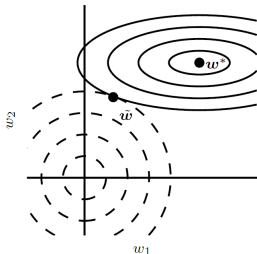
- $\mathbf{w}^* = \arg \min \mathbf{f}(\mathbf{w})$, $\tilde{\mathbf{w}} = \arg \min \tilde{\mathbf{f}}(\mathbf{w})$,

Effect of weight decay on solution ?

- if $\lambda_i[\nabla^2 f(\mathbf{w}^*)] \gg \lambda$: direction i sensitive
 $\tilde{w}_i \approx w_i^*$: no change in sol
- if $\lambda_i[\nabla^2 f(\mathbf{w}^*)] \ll \lambda$: direction i not sensitive
 $\tilde{w}_i \approx 0$: shrink component to zero

no change in solution along sensitive direction.

lower model complexity: $\|\tilde{\mathbf{w}}\|_2 \leq \|\mathbf{w}^*\|_2$



solid line = f , dashed line = r . $\tilde{\mathbf{w}}$ closer to origin (lower magnitude) than \mathbf{w}^* . Reg pushes optimal solution closer to 0 (lower norm)

Regularization: Parameter norm penalty

L_1 -regularization: (sparse regularization)

- $\mathbf{w}^* = \arg \min \mathcal{f}(\mathbf{w})$, $\tilde{\mathbf{w}} = \arg \min \tilde{\mathcal{f}}(\mathbf{w})$,

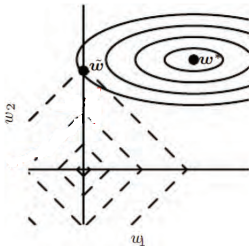
Effect on solution ?

- if $w_i^* \leq \lambda / [\nabla^2 \mathcal{f}(\mathbf{w}^*)]_{i,i}$: direction i not sensitive. set $\tilde{w}_i = 0$ to zero
- if $w_i^* \geq \alpha / [\nabla^2 \mathcal{f}(\mathbf{w}^*)]_{i,i}$: direction i sensitive. set $\tilde{w}_i = w_i^*$ (no change)

Round w_i^* to zero if below threshold

L_1 -regularization gives *sparse* solutions.

- Insignificant components of model (small entries in \mathbf{w}) set to zero



solid line = \mathcal{f} , dashed line = r .
Due to shape of ℓ_1 norm, $\tilde{\mathbf{w}}$ has one non-zero component. Reg makes solution sparse.

Overfitting

- *Data Augmentation*: Avoid overfitting with more data in training set
apply transformations on training set to get more data
- *Early stopping*: evaluate training and test error with each epoch.
return model with lower test error
- *Dropout training*: drop weights in DNN following a distribution

see chapter 7 [Goodfellow, 2016] for detailed treatment of regularization

Outline

1. Challenges for DNN optimization
2. Large-Scale methods for DNN training
3. Regularization in DNNs
4. Some tricks of the trade

Batch normalization

Input normalization to bring all the features on the same scale
improves learning

Why not applying this idea to every layer? two new parameters to learn per layer (γ_j and β_j), same training pipeline, almost same backprop.

Batch normalization for layer j [Ioffe-Szegedy-2015]

Consider a minibatch of size N_j , $\mathcal{B} = \{z_j^{[i]}\}_{i \in [N_j]}$

Compute batch mean $m_j^{[\mathcal{B}]}$ and batch standard deviation $\Sigma_j^{[\mathcal{B}]}$

$$m_j^{[\mathcal{B}]} = \frac{1}{N_j} \sum_{i \in \mathcal{B}} z_j^{[i]}, \quad \Sigma_j^{[\mathcal{B}]} = \frac{1}{N_j} \sum_{i \in \mathcal{B}} \left(z_j^{[i]} - m_j^{[\mathcal{B}]} \right)^2$$

For every $i \in \mathcal{B}$, normalize the outputs by $z_{\text{norm},j}^{[i]} = \frac{z_j^{[i]} - m_j^{[\mathcal{B}]}}{\sqrt{\Sigma_j^{[\mathcal{B}]} + \epsilon}}$

Use the new output: $\hat{z}_j^{[i]} = \gamma_j z_{\text{norm},j}^{[i]} + \beta_j$

Tricks of the trade

Initialization of algo critical in *non-convex opt* (unlike convex). Why ?
- *Heuristics* (verified empirically). No principled approach for DNNs

Parameter Initialization Strategies for Weight

Recall model for layer j : $\mathbf{z}_j = \sigma_j(\mathbf{W}_j \mathbf{z}_{j-1})$, where $\mathbf{W}_j \in \mathbb{R}^{d_j \times d_{j-1}}$

- 1 *Random Init*: $[\mathbf{W}_j]_{i,k} \sim U(\frac{-1}{\sqrt{d_{j-1}}}, \frac{1}{\sqrt{d_{j-1}}})$
- 2 *Xavier Init (2015)*: $[\mathbf{W}_j]_{i,k} \sim U(\frac{-6}{\sqrt{d_j d_{j-1}}}, \frac{6}{\sqrt{d_j d_{j-1}}})$
 - takes into account matrix mult (w/out non-linearity)
- 3 *Saxe Init (2013)*: Select weight as random orthogonal matrices.
 - increase norm of activation fnc (avoid vanishing grad)
- 4 *Martens Init (2010)*: Initial weights are sparse.

Conclusions

large scale training for DNNs

- Challenges:
non-convex loss, exponential number of saddle points, ill and poor conditioning
- SoA Methods:
SGD, Nesterov, minibatch SGD, ADAGRAD, RMSProp, ADAM
- Regularization: weight decay, sparse
- Heuristics: batch normalization, initialization
- Many SoA methods not covered:
Nesterov acceleration, Polyak averaging, Dropout training,

Some references

- A. Goodfellow, Y. Bengio, A Courville, “Deep Learning”, MIT Press, Chap. 7,8
- Y. LeCun, L Bottou, G. Orr, K-B Muller, “Efficient Backprop”, Neural Networks: Tricks of the trade
- C. Bishop, “Pattern recognition and machine learning”, 2006, Springer, Chap. 1
- S. Trivedi, R. Kondor, CMSC 35246 Deep Learning, spring 2017
- S. Bach, Stochastic gradient descent and robustness to ill-conditioning, May 2016
- L. Bottou, F. E. Curtis, J Norcedal, “Optimization Methods for large-scale machine learning ”, Sec. 2-3



Deep Learning

Lecture 6: Large-scale Training for Deep Neural Networks

University of Agder,
Kristiansand, Norway

Prepared by: Hadi Ghauch^{*}, Hossein S. Ghadikolaei[†]

^{*} Telecom ParisTech

[†] Royal Institute of Technology, KTH

<https://sites.google.com/view/fundl/home>