

Optimization and Regularization

AI602: Recent Advances in Deep Learning
Lecture 2

Slide made by

Insu Han
KAIST EE

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

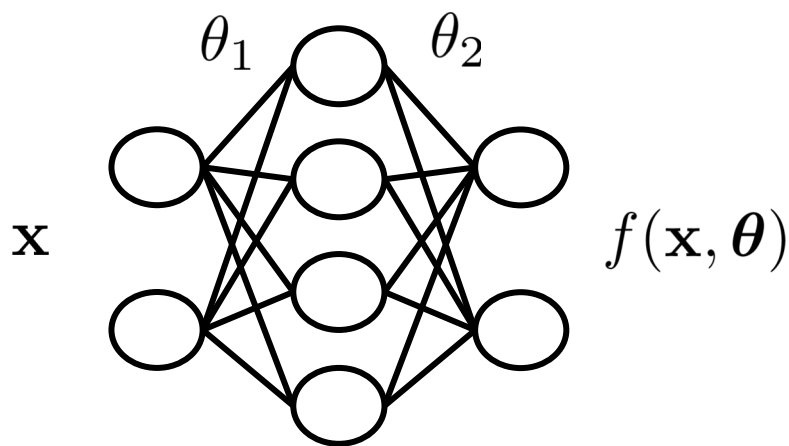
Empirical Risk Minimization (ERM)

- Given training set $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$
- Prediction function $f(\mathbf{x}_i, \boldsymbol{\theta}) \approx y_i$ parameterized by $\boldsymbol{\theta}$
- **Empirical risk minimization:** Find a parameter that minimizes the loss function

$$\min_{\boldsymbol{\theta}} \frac{1}{n} \sum_{i=1}^n \ell(f(\mathbf{x}_i, \boldsymbol{\theta}), y_i) := L(\boldsymbol{\theta})$$

where $\ell(\cdot, \cdot)$ is a loss function e.g., MSE, cross entropy,

- For example, neural network has $f(\mathbf{x}, \boldsymbol{\theta}) = \theta_k^\top \sigma(\theta_{k-1}^\top \sigma(\dots \sigma(\theta_1^\top \mathbf{x})))$



$$L(\boldsymbol{\theta}) = \frac{1}{n} \sum_i (f(\mathbf{x}_i, \boldsymbol{\theta}) - y_i)^2$$

Next, how to solve ERM?

Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

Gradient Descent (GD)

- **Gradient descent (GD)** updates parameters iteratively by taking gradient.

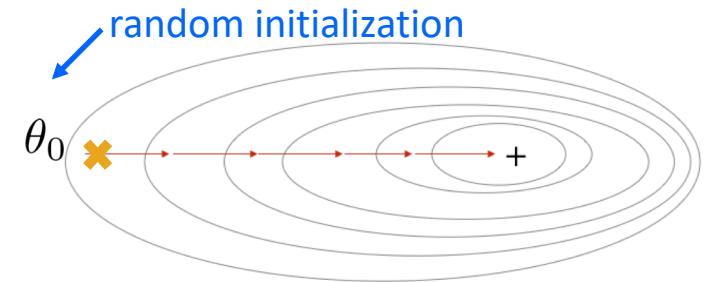
$$\theta_{t+1} = \theta_t - \underbrace{\gamma}_{\text{learning rate}} \underbrace{\nabla L(\theta_t)}_{\text{loss function}}$$

parameters

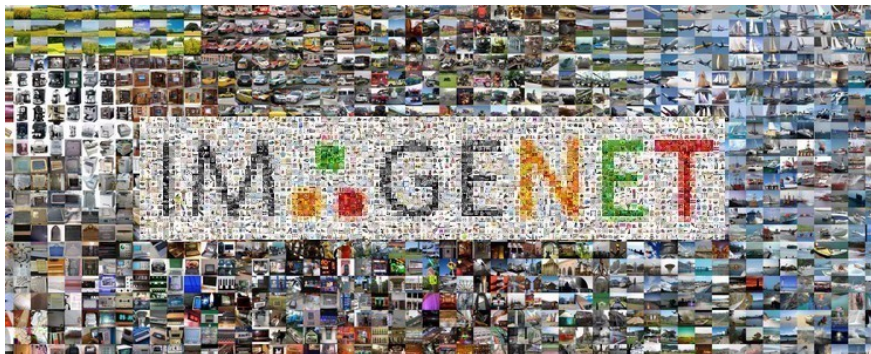
loss function

learning rate

$$:= \frac{1}{n} \sum_{i=1}^n \nabla \ell(\theta_t; \mathbf{x}_i, y_i)$$



- (+) Converges to global (local) minimum for convex (non-convex) problem.
- (−) Not efficient with respect to **computation time** and **memory space** for huge n .
- For example, ImageNet dataset has $n = \mathbf{1,281,167}$ images for training.



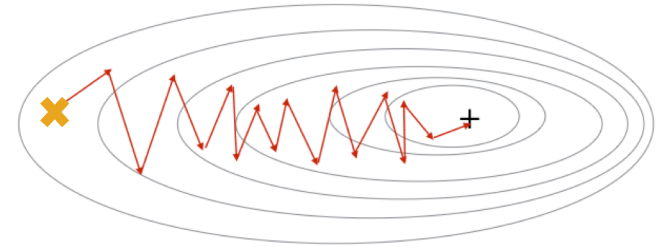
1.2M of 256x256 RGB images
≈ 236 GB memory

Next, efficient GD

Stochastic Gradient Descent (SGD)

- Stochastic gradient descent (SGD) use **samples** to approximate GD

$$\begin{aligned}\nabla L(\boldsymbol{\theta}) &= \frac{1}{n} \sum_{i=1}^n \nabla \ell(\boldsymbol{\theta}; \mathbf{x}_i, y_i) \\ &\approx \frac{1}{|\mathcal{B}|} \sum_{\text{sample } i \in \mathcal{B}} \nabla \ell(\boldsymbol{\theta}; \mathbf{x}_i, y_i)\end{aligned}$$

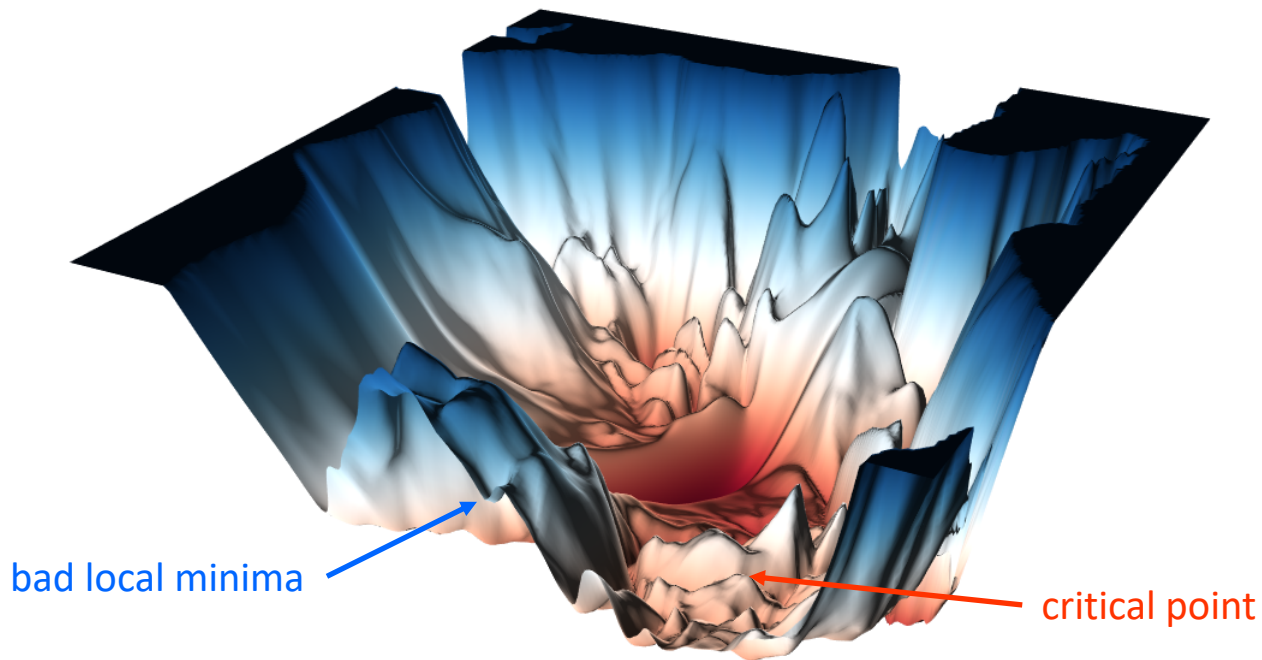


- In practice, minibatch sizes $|\mathcal{B}|$ can be 32/64/128.
- SGD can find the global solution when
 - loss function is convex
 - bounded variance
 - decreasing learning rate
- But, in many practical problems, SGD has some challenges

Hard to optimize practical problems

- Main practical challenges and current solutions:

1. Loss function is nonconvex and includes local minima/critical points
2. SGD can be too noisy and might be unstable —————→ momentum
3. hard to find a good learning rate —————→ adaptive learning rate



loss surface of neural net (ResNet-50)

Next, momentum

Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

Momentum Methods

1. Momentum gradient descent

- Add **decaying previous gradients (momentum)**.

$$\theta_{t+1} = \theta_t - \mathbf{m}_t$$

↓
momentum

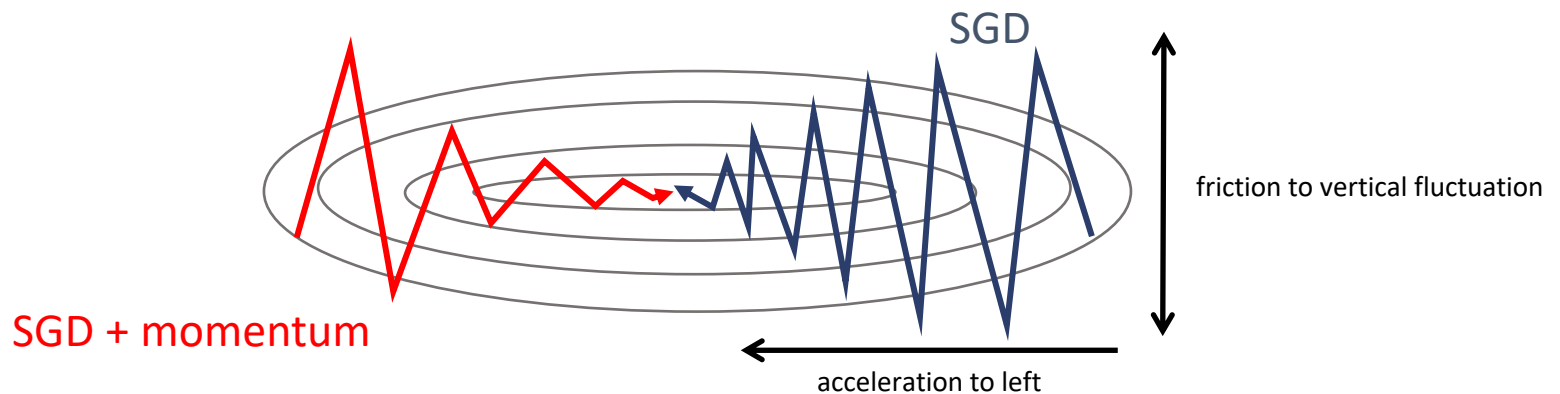
$$\mathbf{m}_t = \mu \mathbf{m}_{t-1} + \gamma \nabla L(\theta_t)$$

↓
preservation ratio $\mu \in [0, 1]$

- Equivalent to **moving average** with the fraction μ of previous update.

$$\theta_{t+1} = \theta_t - \gamma (\nabla L(\theta_t) + \mu \nabla L(\theta_{t-1}) + \mu^2 \nabla L(\theta_{t-2}) + \dots)$$

- (+) Momentum reduces the oscillation and accelerates the convergence.



Momentum Methods: Nesterov's Momentum

1. Momentum gradient descent

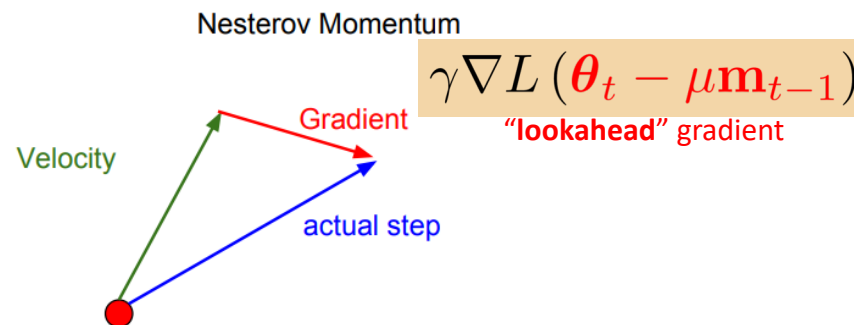
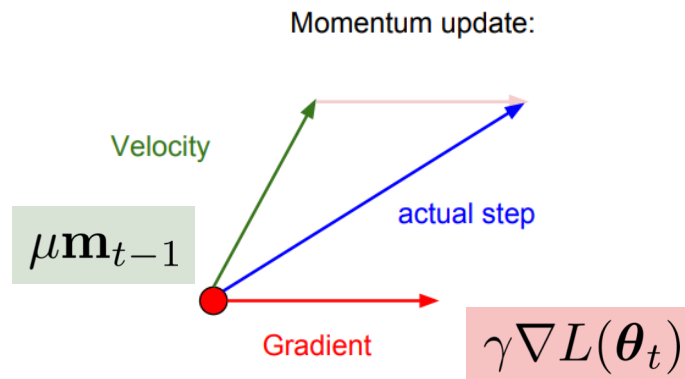
- Add **decaying previous gradients (momentum)**.

$$\theta_{t+1} = \theta_t - \underset{\substack{\downarrow \\ \text{momentum}}}{\mathbf{m}_t}$$

$$\mathbf{m}_t = \underset{\substack{\downarrow \\ \text{preservation ratio } \mu \in [0, 1]}}{\mu} \mathbf{m}_{t-1} + \gamma \nabla L(\theta_t)$$

- (–) Momentum can fail to converge even for simple convex optimizations.
- **Nesterov's accelerated gradient (NAG)** [Nesterov' 1983] use gradient for **approximate future position**, i.e.,

$$\mathbf{m}_t \leftarrow \mu \mathbf{m}_{t-1} + \gamma \nabla L(\theta_t - \mu \mathbf{m}_{t-1})$$



Momentum Methods: Nesterov's Momentum

1. Momentum gradient descent

- Add **decaying previous gradients (momentum)**.

$$\theta_{t+1} = \theta_t - \mathbf{m}_t$$

↓
momentum

$$\mathbf{m}_t = \mu \mathbf{m}_{t-1} + \gamma \nabla L(\theta_t)$$

↓
preservation ratio $\mu \in [0, 1]$

- **Nesterov's accelerated gradient (NAG)** [Nesterov' 1983] use gradient for **approximate future position**, i.e.,

$$\mathbf{m}_t \leftarrow \mu \mathbf{m}_{t-1} + \gamma \nabla L(\theta_t - \mu \mathbf{m}_{t-1})$$

```
while True:
    dtheta = compute_gradient(theta, batch_size)
    theta = theta - learning_rate * dtheta
```

SGD

```
m = 0
while True:
    dtheta = compute_gradient(theta, batch_size)
    m = mu * m + learning_rate * dtheta
    theta = theta - m
```

SGD + momentum

```
m = 0
while True:
    dtheta = compute_gradient(theta, batch_size)
    m_old = m
    m = mu * m + learning_rate * dtheta
    x =
```

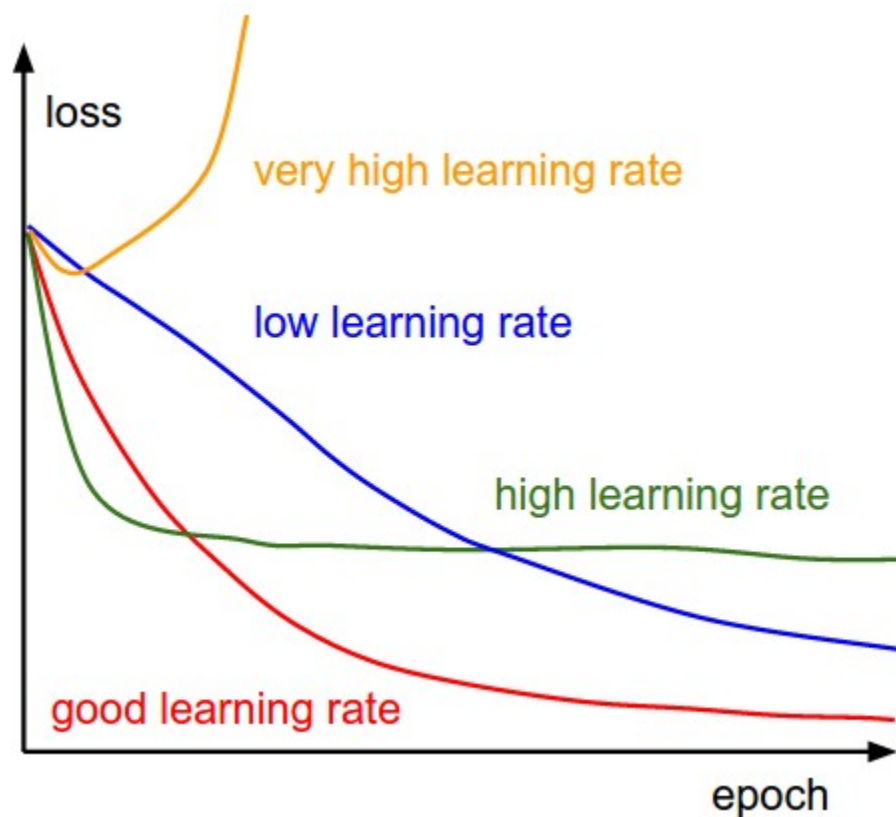
NAG

Quiz: fill in the pseudo code of Nesterov' accelerated gradient

Adaptive Learning Rate Methods

2. Learning rate scheduling

- Learning rate is critical for minimizing loss !



Too **high** → May ignore the narrow valley, can diverge

Too **low** → May fall into the local minima, slow converge

Next, learning rate scheduling

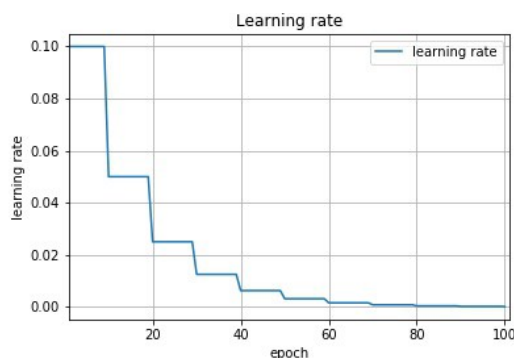
Adaptive Learning Rate Methods: Learning rate annealing

2. Learning rate scheduling : decay methods

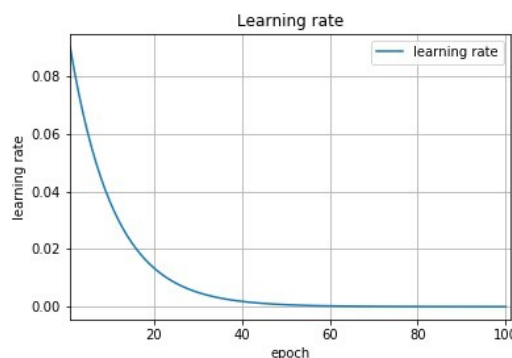
- A naive choice is the **constant** learning rate
- Common learning rate schedules include time-based/step/exponential decay

	Time-based	Exponential	Step (most popular in practice)
γ_t	$\frac{\gamma_0}{1 + kt}$	$\gamma_0 \exp(-kt)$	$\gamma_0 \exp(-k \lfloor \frac{t}{T_{\text{epoch}}} \rfloor)$

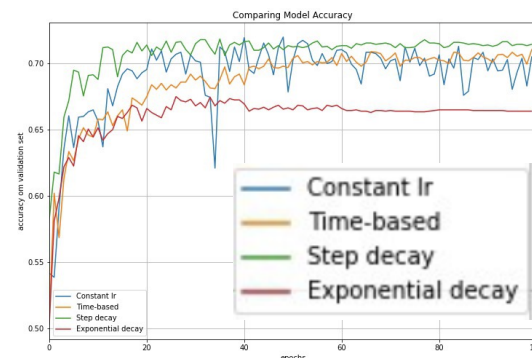
- “Step decay” decreases learning rate by a factor every few epochs
- Typically, it is set $\gamma_0 = 0.01$ and drops by half ever $T_{\text{epoch}} = 10$ epoch



step decay



exponential decay

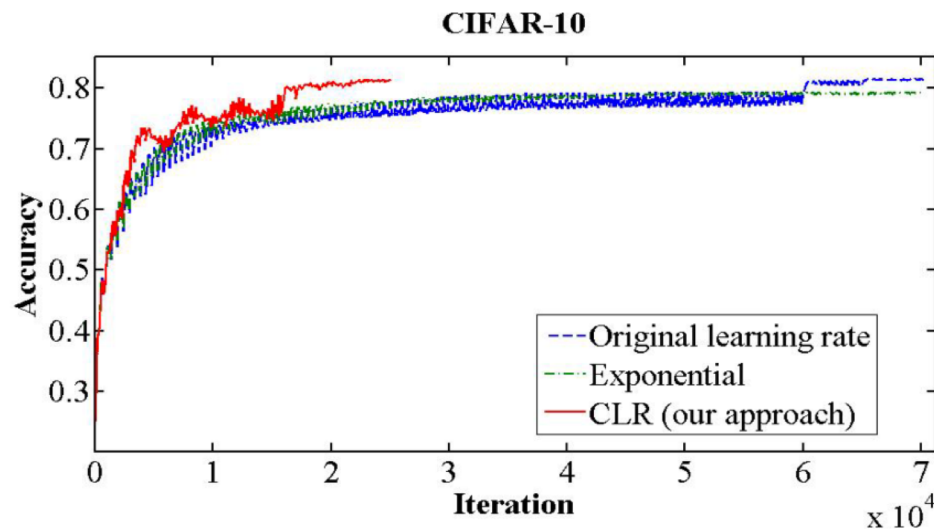
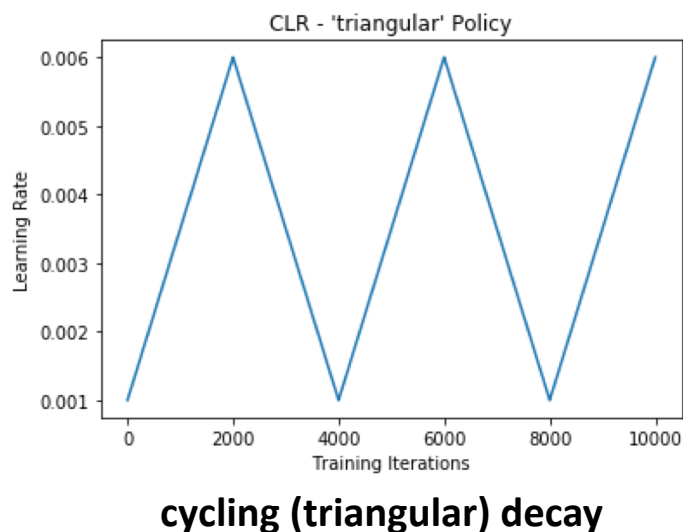


accuracy

Adaptive Learning Rate Methods: Learning rate annealing

2. Learning rate scheduling : cycling method

- [Smith' 2015] proposed **cycling** learning rate (triangular)
- Why “**cycling**” learning rate?
 - Sometimes, **increasing learning rate** is helpful to escape the saddle points
- It can be combined with exponential decay or periodic decay

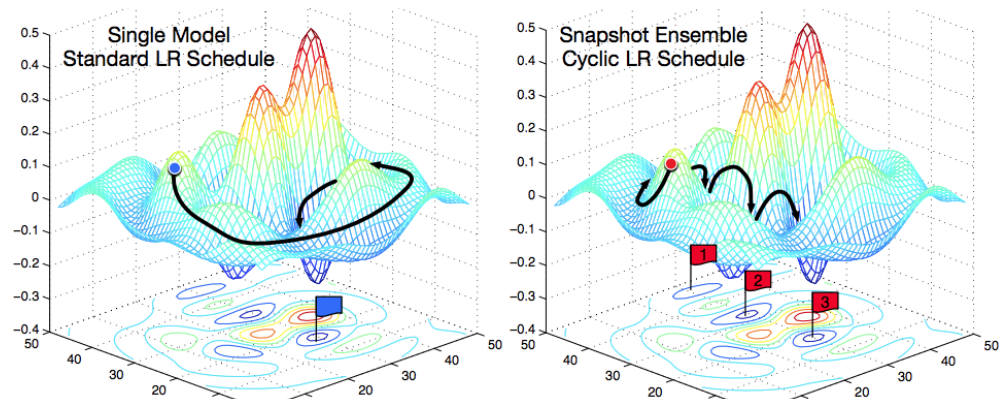
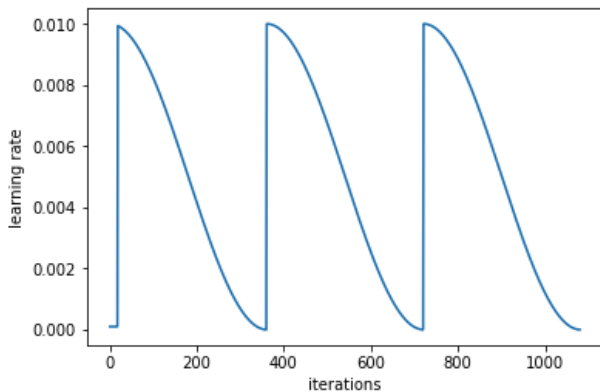


Adaptive Learning Rate Methods: Learning rate annealing

2. Learning rate scheduling : cycling method

- [Loshchilov' 2017] use **cosine cycling** and **restart** the maximum at each cycle
- Why “**cosine**” ?
 - It decays slowly at the half of cycle and drop quickly at the rest
- (+) can climb down and up the loss surface, thus can traverse several local minima
- (+) same as restarting at good points with an initial learning rate γ_{\max}

$$\gamma_t = \gamma_{\min} + \frac{1}{2} (\gamma_{\max} - \gamma_{\min}) (1 + \cos(\text{mod}(t, T)\pi)) \quad T : \text{period}$$



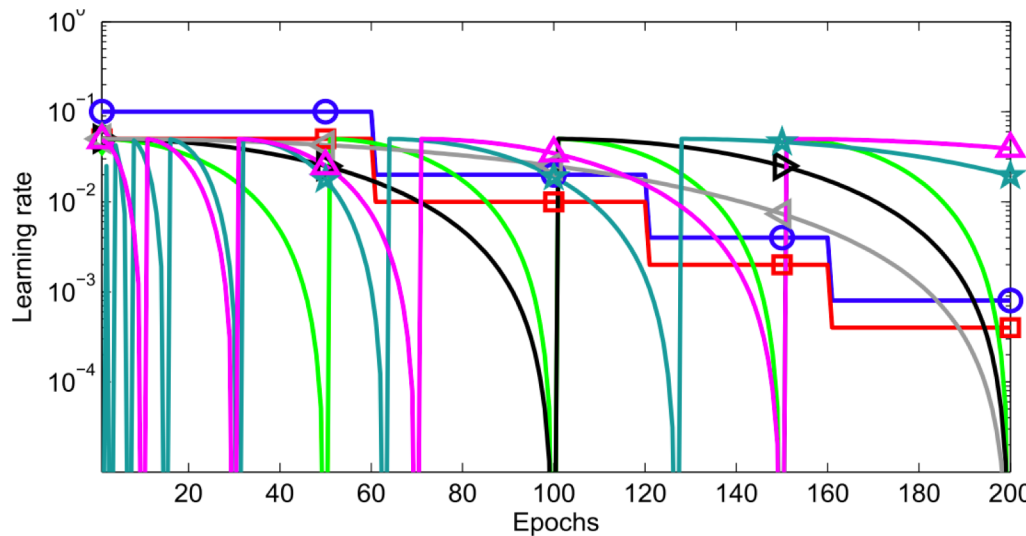
Adaptive Learning Rate Methods: Learning rate annealing

2. Learning rate scheduling : cycling method

- [Loshchilov' 2017] also proposed **warm restart** in cycling learning rate

*Warm restart : frequently restart in early iterations

- (+) It help to escape saddle points since it is more likely to stuck in early iteration



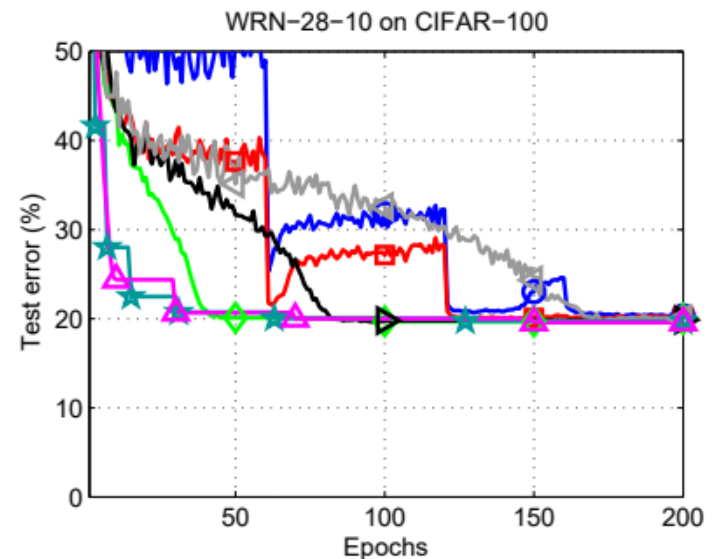
—○— : step decay
—□—



—▶— : cycling with no restart



—△— : cycling with restart



But, there is no perfect learning rate scheduling! It depends on specific task.

Next, adaptive learning rate

3. Adaptively changing learning rate (AdaGrad, RMSProp)

- **AdaGrad** [Duchi' 11] downscales a learning rate by magnitude of previous gradients.

$$\theta_{t+1} = \theta_t - \frac{\gamma}{\sqrt{v_t}} \nabla L(\theta_t) \quad v_{t+1} = v_t + \nabla L(\theta_t)^2$$

↓
sum of all previous squared gradients

- (–) the learning rate strictly decreases and becomes too small for large iterations.
- **RMSProp** [Tieleman' 12] uses the moving averages of squared gradient.

$$v_{t+1} = \mu v_t + (1 - \mu) \nabla L(\theta_t)^2$$

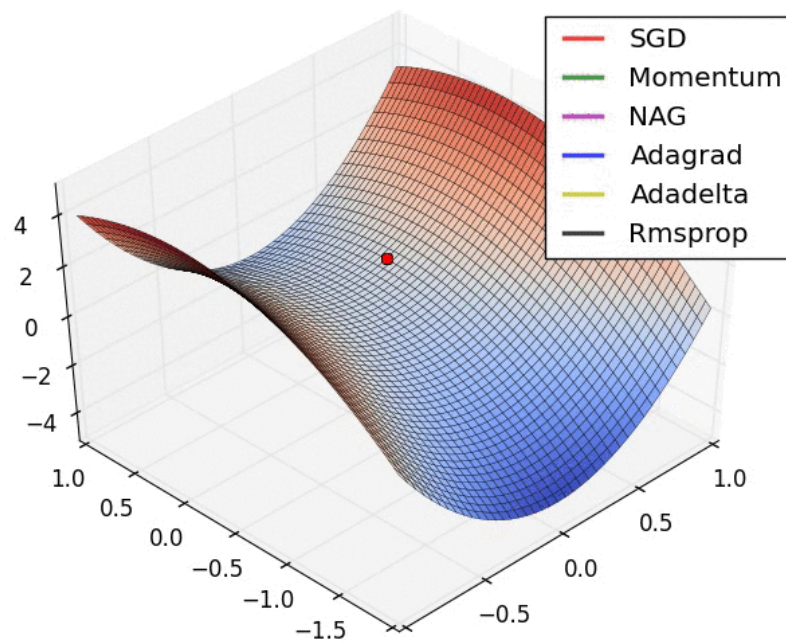
↓
preservation ratio $\mu \in [0, 1]$

- Other variants also exist, e.g., Adadelta [Zeiler' 2012]

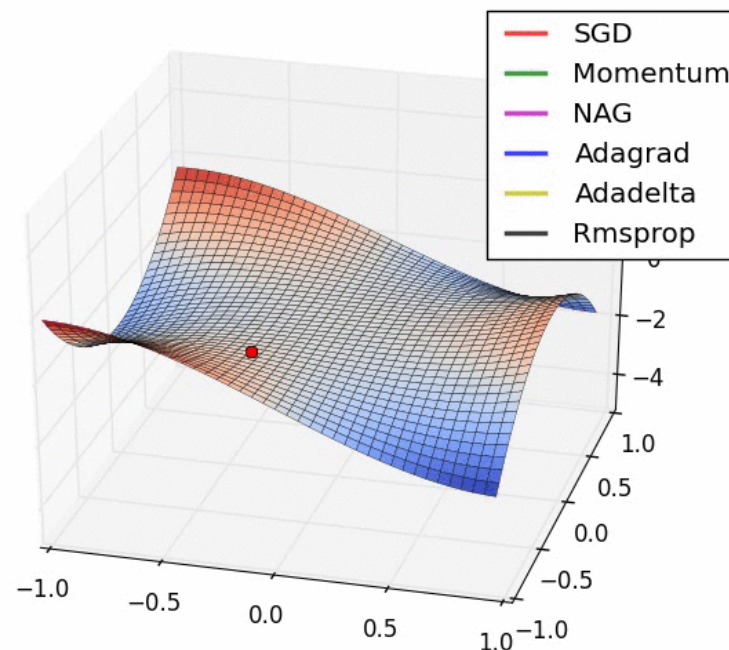
Adaptive Learning Rate Methods

- Visualization of algorithms

optimization on saddle point



optimization on local optimum



- Adaptive learning-rate methods, i.e., Adadelta and RMSprop are most suitable and provide the best convergence for these scenarios

Next, momentum + adaptive learning rate

Adaptive Learning Rate Methods: ADAM

1 + 2. Combination of momentum and adaptive learning rate

- **Adam** (ADAPtive Moment estimation) [Kingma' 2015]

$$\begin{aligned}\theta_{t+1} &\leftarrow \theta_t - \frac{\gamma}{\sqrt{v_t}} m_t \\ m_{t+1} &\leftarrow \mu_1 m_t + (1 - \mu_1) \nabla L(\theta_t) \\ v_{t+1} &\leftarrow \mu_2 v_t + (1 - \mu_2) \nabla L(\theta_t)^2\end{aligned}$$

momentum
average of squared gradients

- Can be seen as momentum + RMSprop update.
- Other variants exist, e.g., Adamax [Kingma' 14], Nadam [Dozat' 16]

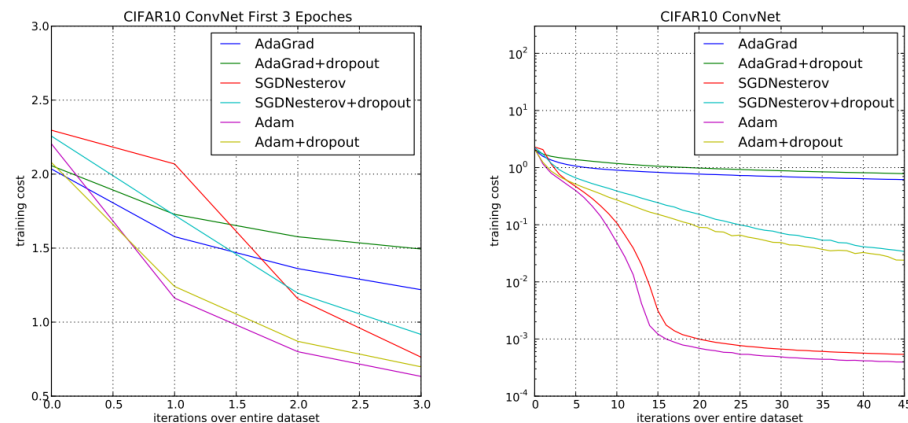
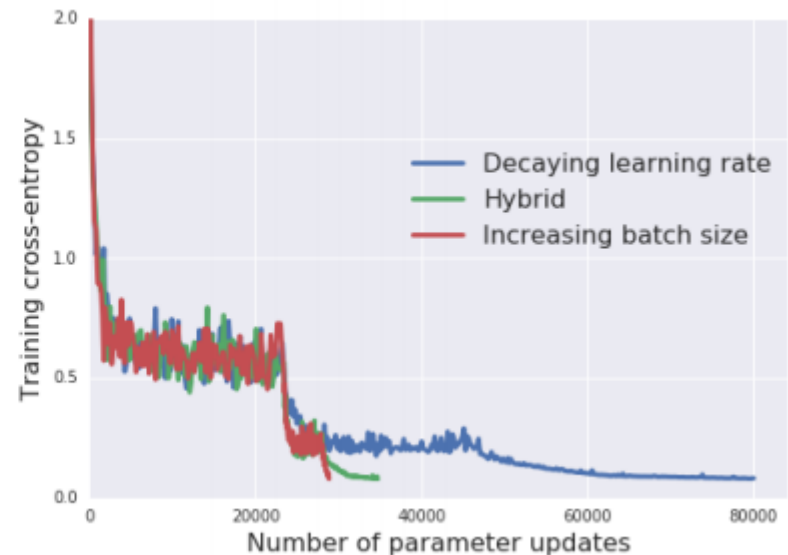
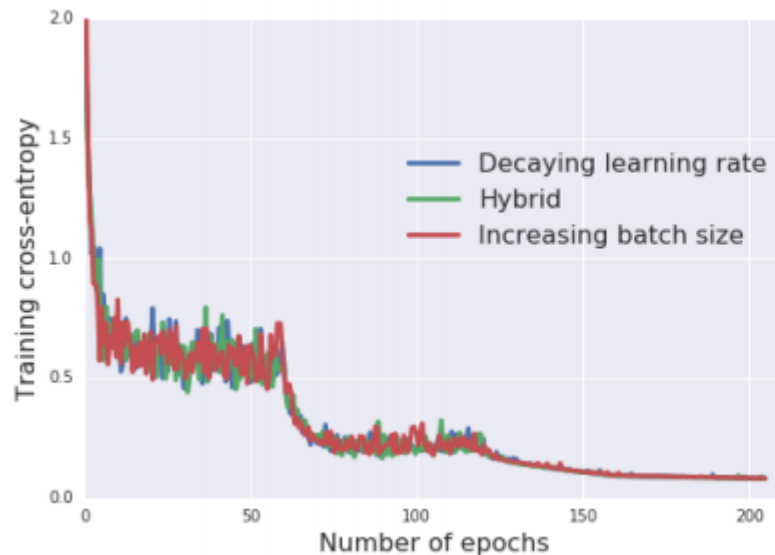


Figure 3: Convolutional neural networks training cost. (left) Training cost for the first three epochs. (right) Training cost over 45 epochs. CIFAR-10 with c64-c64-c128-1000 architecture.

Decaying the Learning Rate = Increasing the Batch Size

- In practice, **SGD + Momentum** and **Adam** works well in many applications.
- But, scheduling learning rates is still critical! (should be decay appropriately)
- [Smith' 2017] shows that decaying learning rate = increasing batch size,
 - (+) A large batch size allows fewer parameter updates, leading to parallelism!



Next, decoupled SGD with momentum

SGD/Adam with decoupled weight decay

- Many learning problems optimize the loss with **L^2 norm penalty** (details in later)

$$\tilde{L}(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2^2$$

- It is sometimes called “weight decay” since its gradient decays weight:

$$\underbrace{\boldsymbol{\theta} - \eta \nabla (L(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2^2)}_{\text{SGD on L2-norm penalty}} \quad \Leftrightarrow \quad \underbrace{(1 - 2\eta\lambda)\boldsymbol{\theta} - \eta \nabla L(\boldsymbol{\theta})}_{\text{weight decay}} \quad \nabla \|\boldsymbol{\theta}\|_2^2 = 2\boldsymbol{\theta}$$

- However, both are different when SGD with momentum methods (check!)
- [Loshchilov 2019'] proposes **decoupled** weight decay with momentum.
- For example, decoupled SGD with momentum iterates (also applicable to Adam)

$$\mathbf{m}_{t+1} \leftarrow \mu_1 \mathbf{m}_t + (1 - \mu_1) (\underbrace{\nabla L(\boldsymbol{\theta}_t) + \lambda \boldsymbol{\theta}_t}_{\text{gradient of loss with L2 penalty}})$$

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \mathbf{m}_t - \underbrace{2\eta\lambda\boldsymbol{\theta}_t}_{\text{weight decay}}$$

SGD/Adam with decoupled weight decay

- Many learning problems optimize the loss with **L^2 norm penalty** (details in later)

$$\tilde{L}(\boldsymbol{\theta}) = L(\boldsymbol{\theta}) + \lambda \|\boldsymbol{\theta}\|_2^2$$

- It is sometimes called “weight decay” since its gradient decays weight.
- However, both are different when SGD with momentum methods
- [Loshchilov 2019'] proposes **decoupled** weight decay with improved results than standard Adam

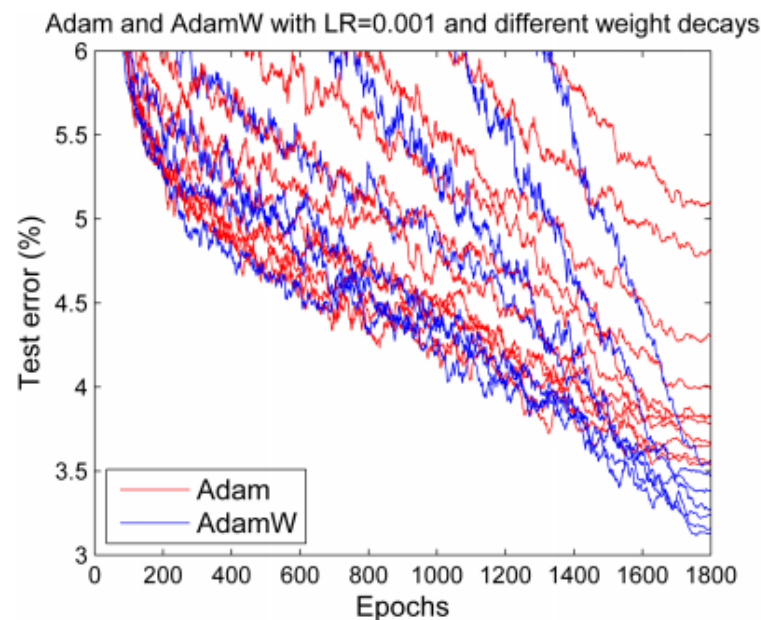
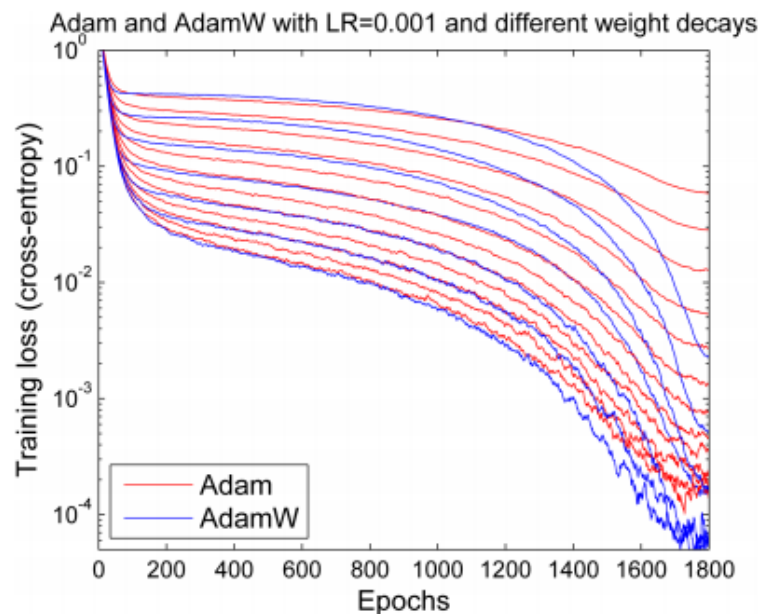


Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

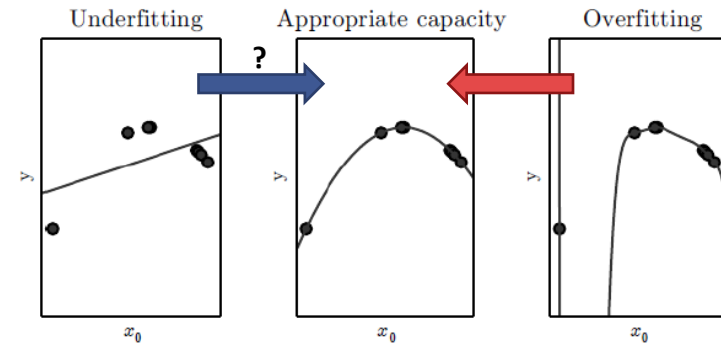
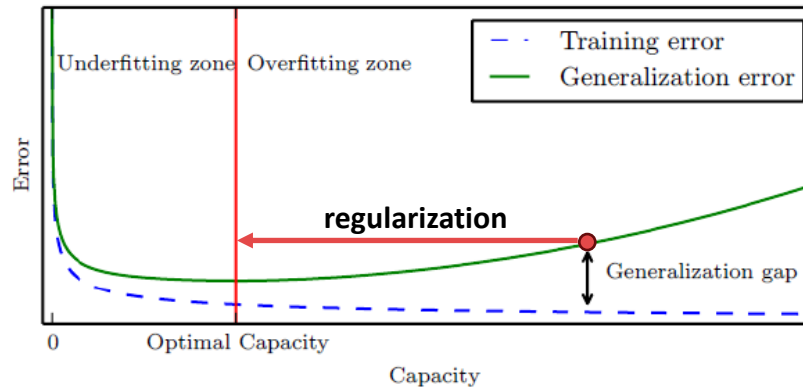
3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

Regularization

- **Overfitting** is a central problem in machine learning



- Why overfitting? model capacity (number of parameters) is too large
- **Regularization:** any modification to reduce the generalization error
 - The main challenge is to find **a right model complexity** for a given task
 - There is no universal model working for all tasks (no free lunch theorem)

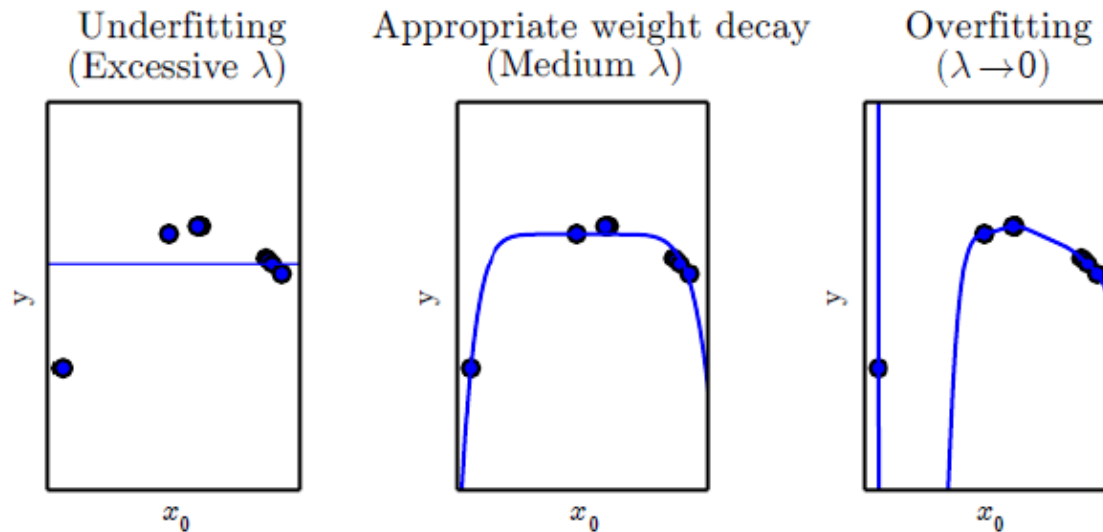
- **Practical regularizations in Neural Networks**

	Regularizations
Loss penalty	<ul style="list-style-type: none">• Parameter norm penalty (L2/L1/L0-norm decay)• Directly approximately regularizing complexity• ...
Parameter sharing	<ul style="list-style-type: none">• Convolutional neural networks• Skip connections• ...
Noise robustness	<ul style="list-style-type: none">• Noises on hidden units (Dropout)• Noises on gradients (Shake-shake)• ...
Data augmentation	<ul style="list-style-type: none">• Making new data by local masking (CutOut)• Mixing two samples in dataset (mixup)• ...

- Adding a parameter penalty $\Omega(\boldsymbol{\theta}) \geq 0$ to the objective loss

$$\tilde{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \lambda \Omega(\boldsymbol{\theta})$$

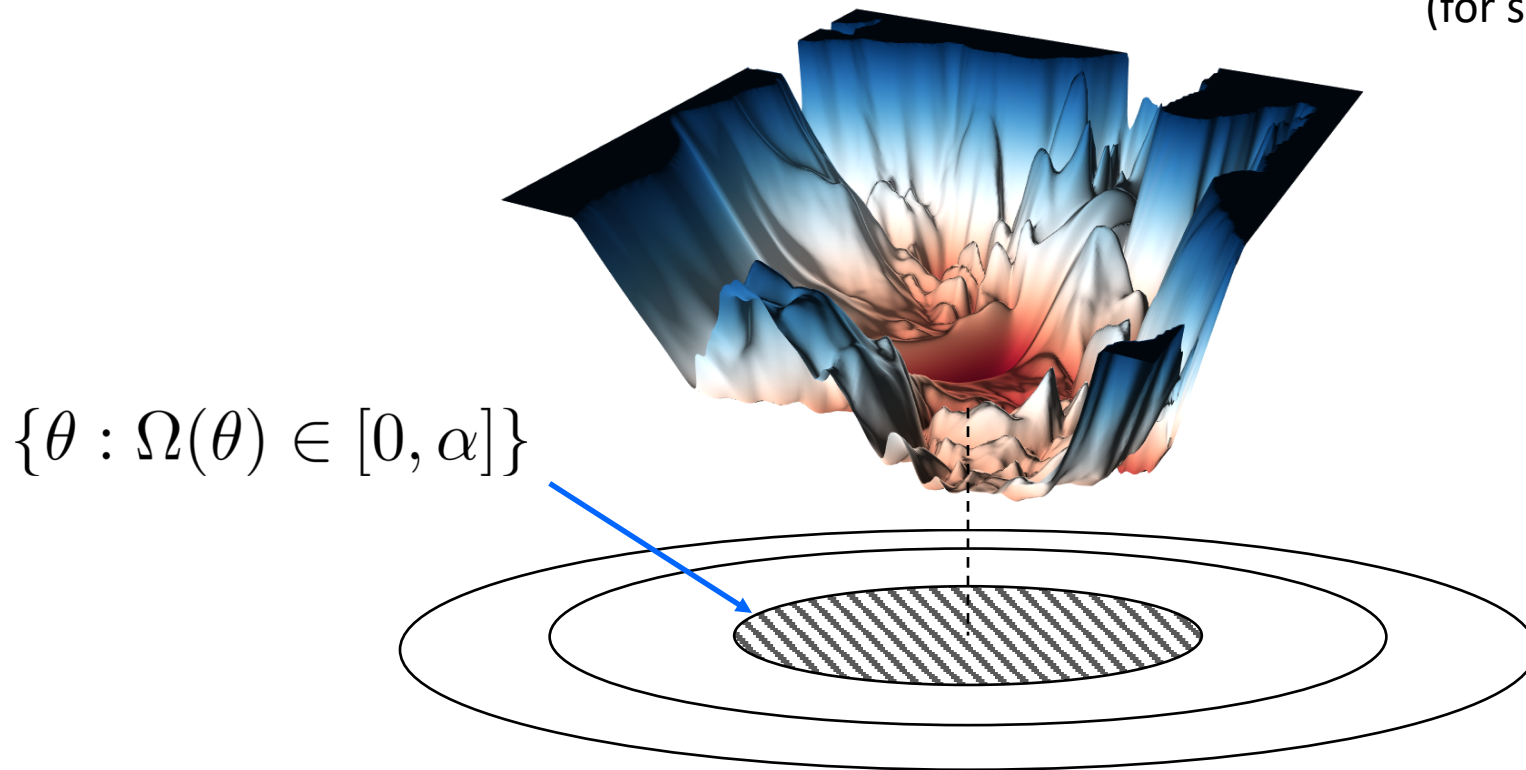
- $\lambda \in [0, \infty)$: a hyperparameter that controls the relative power of $\Omega(\boldsymbol{\theta})$
- Different penalty Ω results in a different solution being preferred



- **Parameter norm penalty:** constraint on the **search space** of parameters θ

$$\min_{\theta} L(\theta) + \lambda \Omega(\theta) \quad \Leftrightarrow \quad \min_{\theta} L(\theta) \quad \text{s.t.} \quad \Omega(\theta) \leq \alpha$$

(for some $\alpha > 0$)



Next, L^2 and L^1 regularization

- **Parameter norm penalty:** Penalizing on the **search space** of parameters θ
- The two most commonly used forms: **L^2 and L^1 penalty**

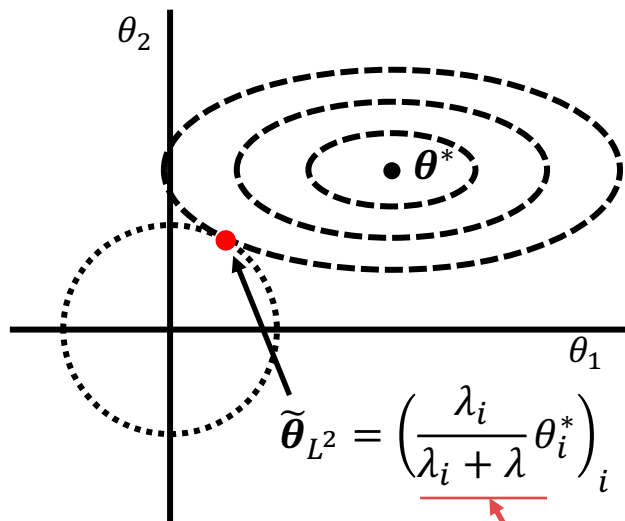
	L^2 ("weight decay")	L^1
$\Omega(\theta)$	$\frac{1}{2} \ \theta\ _2^2 := \frac{1}{2} \sum_i \theta_i^2$	$\ \theta\ _1 := \sum_i \theta_i $
Nickname	Ridge regression	LASSO
MAP Prior	$\mathcal{N}(\theta_i; 0, \frac{1}{\lambda})$	Laplace($\theta_i; 0, \frac{1}{\lambda}$)

- The solution maps to the *maximum a posteriori* (MAP) estimation under a certain **prior on weights**

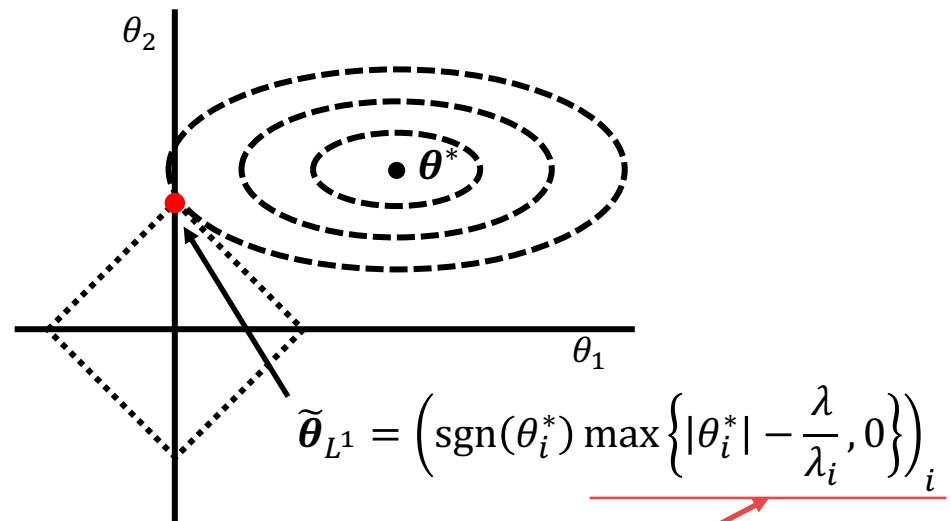
Parameter norm penalty

- If L is quadratic with diagonal Hessian $H = (\lambda_i)_{ii}$, we get the analytic solutions from each regularization [Goodfellow et al., 2016]:

$$\tilde{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \lambda \Omega(\boldsymbol{\theta})$$



"rescaling" eigenvalues

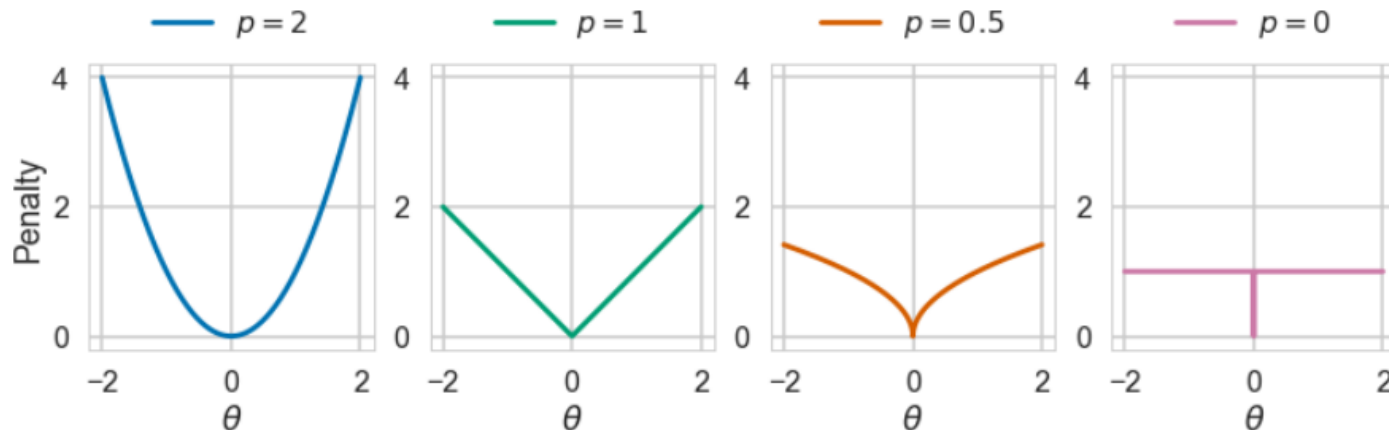


"sparse" solution

Next, L^0 -regularization

Parameter norm penalty: L^0 -regularization [Louizos et al., 2018]

- We typically use the popular L^1 -regularization to induce **sparsity**
 - Sparse models are advantageous on computational efficiency
 - Of course, it is a nice policy for regularization as well
- Why don't we use **L^0 -penalty**?
 - $\Omega(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_0 := |\{\theta_i : \theta_i \neq 0\}|$
 - A more direct measure of sparsity
 - It does not **shrink** the non-sparse weights



Parameter norm penalty: L^0 -regularization [Louizos et al., 2018]

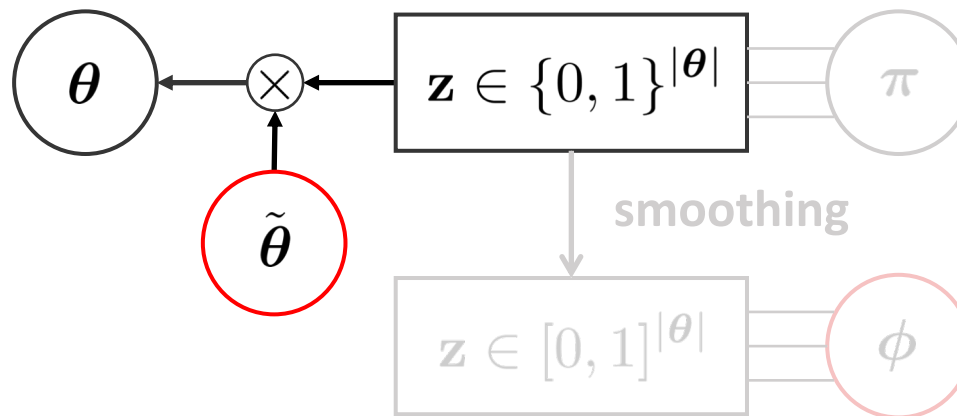
- We typically use the popular L^1 -regularization to induce **sparsity**
 - Sparse models are advantageous on computational efficiency
 - Of course, it is a nice policy for regularization as well
- Why don't we use **L^0 -penalty**?
 - $\Omega(\boldsymbol{\theta}) = \|\boldsymbol{\theta}\|_0 := |\{\theta_i : \theta_i \neq 0\}|$
 - A more direct measure of sparsity
 - It does not **shrink** the non-sparse weights
- **Problem:** Optimization with L^0 -penalty is **intractable** in general
 - Discrete optimization with $2^{|\boldsymbol{\theta}|}$ possible states
 - Standard gradient-based methods are not applicable
- Can we **relax this problem** so that to an efficient continuous optimization?

Parameter norm penalty: L^0 -regularization [Louizos et al., 2018]

- **Idea:** Regard θ as a random variable, where $\mathbb{E}[\|\theta\|_0]$ is **differentiable**
 1. Consider a simple **re-parametrization** of θ :

$$\theta_j = \tilde{\theta}_j z_j, \quad z_j \in \{0, 1\}, \quad \tilde{\theta}_j \neq 0$$

- Then, the L^0 -penalty becomes $\Omega(\theta) = \|\theta\|_0 = \sum_{j=0}^{|\theta|} z_j$



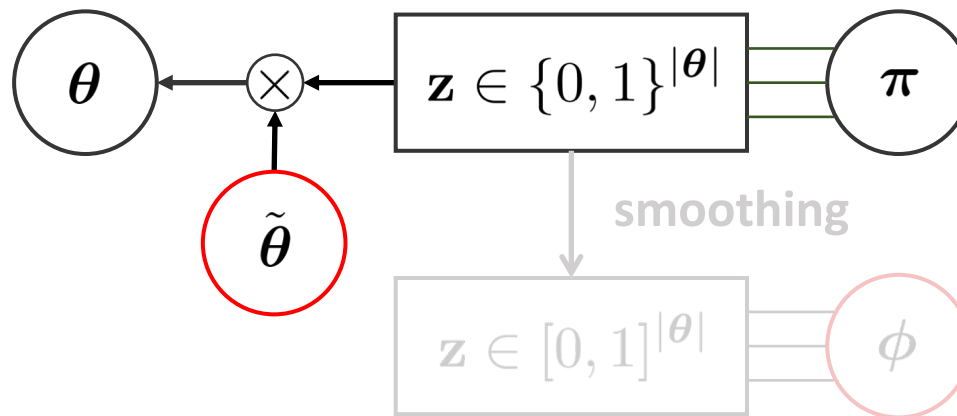
Parameter norm penalty: L^0 -regularization [Louizos et al., 2018]

- **Idea:** Regard θ as a random variable, where $\mathbb{E}[\|\theta\|_0]$ is **differentiable**

2. Letting $q(z_i|\pi_i) = \text{Bernoulli}(\pi_i)$, we define the **expected loss** \mathcal{R} :

$$\mathcal{R}(\tilde{\theta}, \pi) := \mathbb{E}_{q(\mathbf{z}|\pi)} \left[L(\tilde{\theta} \odot \mathbf{z}) \right] + \lambda \sum_{j=1}^{|\theta|} \pi_j$$

- However, optimizing $\mathcal{R}(\tilde{\theta}, \pi)$ is still hard
 - Estimating $\nabla \mathbb{E}_{q(\mathbf{z}|\pi)} \left[L(\tilde{\theta} \odot \mathbf{z}) \right]$ is not easy due to the discrete nature of \mathbf{z}

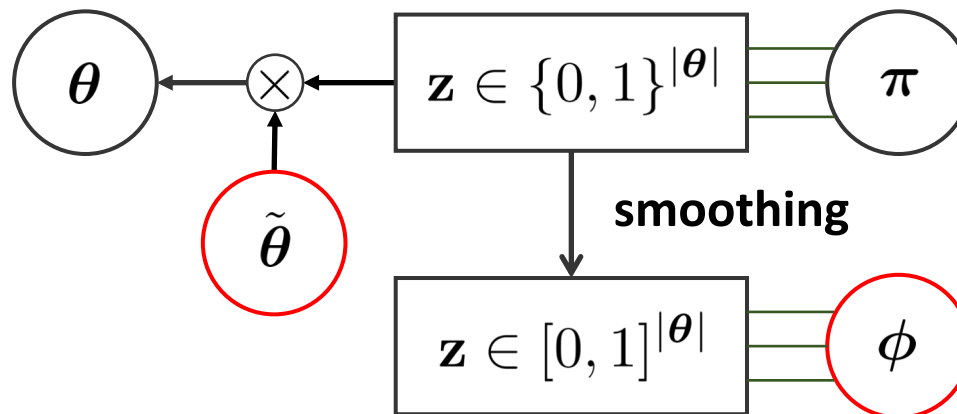


- **Idea:** Regard θ as a random variable, where $\mathbb{E}[\|\theta\|_0]$ is **differentiable**
- 3. Smoothing the discrete random variables \mathbf{z} via a continuous random variables \mathbf{s} :

$$\mathbf{z} = \min(1, \max(0, \mathbf{s})), \quad \mathbf{s} \sim q(\mathbf{s}|\phi)$$

- Since $q(\mathbf{z} \neq 0|\phi) = 1 - \mathbb{P}(\mathbf{s} \leq 0|\phi)$, we get:

$$\mathcal{R}(\tilde{\theta}, \phi) = \mathbb{E}_{q(\mathbf{s}|\phi)} \left[L(\tilde{\theta} \odot \min(1, \max(0, \mathbf{s}))) \right] + \lambda \sum_{j=1}^{|\theta|} (1 - \mathbb{P}(s_j \leq 0|\phi_j))$$

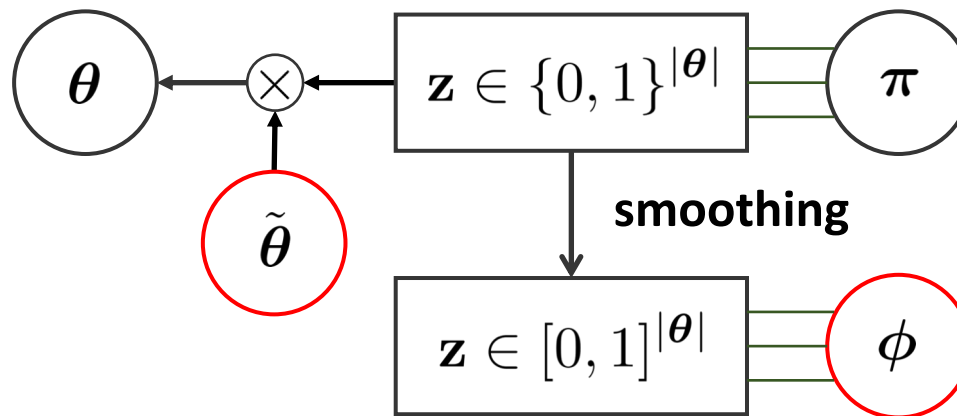


- **Idea:** Regard θ as a random variable, where $\mathbb{E}[\|\theta\|_0]$ is **differentiable**

- Finally, the original loss \tilde{L} is transformed by:

$$\mathcal{R}(\tilde{\theta}, \phi) = \mathbb{E}_{q(s|\phi)} \left[L(\tilde{\theta} \odot \min(1, \max(0, s))) \right] + \lambda \sum_{j=1}^{|\theta|} (1 - \mathbb{P}(s_j \leq 0 | \phi_j))$$

- We can **optimize** this via minibatch-based gradient estimation methods
 - For details, see [Kingma et al., 2013]



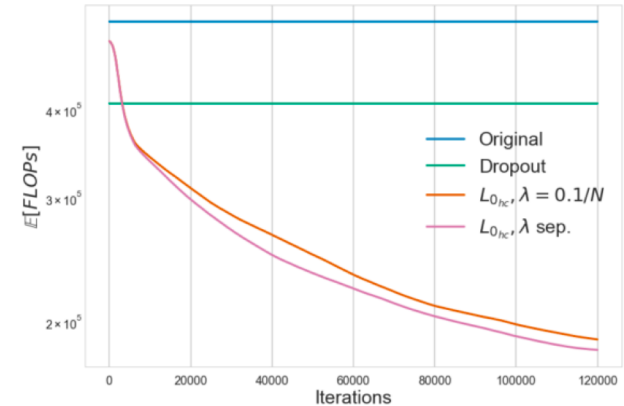
Parameter norm penalty: L^0 -regularization [Louizos et al., 2018]

- L^0 -regularization leads the networks to a sparse solution, with a good regularization as well on MNIST and CIFAR-10/100

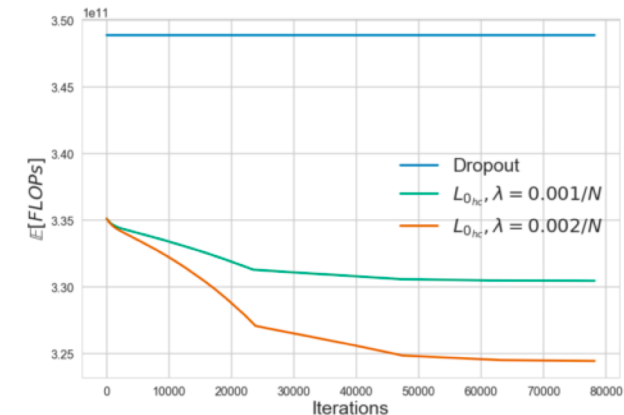
Network & size	Method	Pruned architecture	Error (%)
MLP 784-300-100	Sparse VD (Molchanov et al., 2017)	512-114-72	1.8
	BC-GNJ (Louizos et al., 2017)	278-98-13	1.8
	BC-GHS (Louizos et al., 2017)	311-86-14	1.8
	$L_{0_{hc}}, \lambda = 0.1/N$	219-214-100	1.4
	$L_{0_{hc}}, \lambda \text{ sep.}$	266-88-33	1.8
LeNet-5-Caffe 20-50-800-500	Sparse VD (Molchanov et al., 2017)	14-19-242-131	1.0
	GL (Wen et al., 2016)	3-12-192-500	1.0
	GD (Srinivas & Babu, 2016)	7-13-208-16	1.1
	SBP (Neklyudov et al., 2017)	3-18-284-283	0.9
	BC-GNJ (Louizos et al., 2017)	8-13-88-13	1.0
	BC-GHS (Louizos et al., 2017)	5-10-76-16	1.0
	$L_{0_{hc}}, \lambda = 0.1/N$	20-25-45-462	0.9
	$L_{0_{hc}}, \lambda \text{ sep.}$	9-18-65-25	1.0

Network	CIFAR-10	CIFAR-100
original-ResNet-110 (He et al., 2016a)	6.43	25.16
pre-act-ResNet-110 (He et al., 2016b)	6.37	-
WRN-28-10 (Zagoruyko & Komodakis, 2016)	4.00	21.18
WRN-28-10-dropout (Zagoruyko & Komodakis, 2016)	3.89	18.85
WRN-28-10- $L_{0_{hc}}, \lambda = 0.001/N$	3.83	18.75
WRN-28-10- $L_{0_{hc}}, \lambda = 0.002/N$	3.93	19.04

MLP



WRN-28-10



Next, complexity regularization

Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods


3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

- Reducing **complexity of a model** might be a direct way of regularization
 - But, how do we know whether a model is complex or not?
 - **Computational learning theory** provides a way for it
- Suppose we have a **model** F , i.e. **a set of hypothesis functions**
- **DARC** attempts to reduce the **Rademacher complexity** of F :


$$\text{Rad}_m(F) := \mathbb{E}_{\mathbf{x} \sim \mathcal{D}^m} \left[\frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in F} \sum_{i=1}^m \sigma_i f(x_i) \right] \right]$$


sample size

- $\sigma_1, \dots, \sigma_m$: *i.i.d.* random variables, $\mathbb{P}(\sigma_i = 1) = \mathbb{P}(\sigma_i = -1) = \frac{1}{2}$
- High $\text{Rad}_m(F) \Rightarrow F$ is more **expressive** on \mathcal{D}^m
- It can be used to give a bound of the generalization error in ERM
 - For details, see [Shalev-Shwartz et al., 2014]


- **DARC** attempts to reduce the **Rademacher complexity** of F :

$$\text{Rad}_m(F) := \mathbb{E}_{\mathbf{X} \sim \mathcal{D}^m} \left[\frac{1}{m} \mathbb{E}_{\boldsymbol{\sigma}} \left[\sup_{f \in F} \sum_{i=1}^m \sigma_i f(x_i) \right] \right]$$

 sample size

- Of course, computing $\text{Rad}_m(F)$ is **intractable** when F is a family of NNs
- Instead, DARC uses **a rough approximation** of it:

$$\tilde{L}(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) = L(\boldsymbol{\theta}; \mathbf{X}, \mathbf{y}) + \lambda \left(\frac{1}{m} \max_k \sum_{i=1}^m |f_k(x_i; \boldsymbol{\theta})| \right)$$

 mini-batch size

- $f = (f_1, \dots, f_d) \in \mathbb{R}^d$: the model to optimize (e.g. neural network)
- In other words, here F is approximated by $\{f_k : k = 1, \dots, d\}$

- Despite its simplicity, DARC improves state-of-the-art level models
 - Results on MNIST and CIFAR-10 are presented

Table 1: Test error (%). A standard variant of LeNet (LeCun et al., 1998) and ResNeXt-29(16 × 64d) (Xie et al., 2016) are used for MNIST and CIFAR-10, and compared with the addition of the studied regularizer.

Method	MNIST	CIFAR-10
Baseline	0.26	3.52
DARC1	<u>0.20</u>	<u>3.43</u>

- Comparisons in the values of DARC penalty
 - Data augmentation by itself implicitly regularize the DARC penalty

Table 3: Values of $\frac{1}{m} \left(\max_k \sum_{i=1}^m |h_k^{(H+1)}(x_i)| \right)$

Method	MNIST (ND)		MNIST		CIFAR-10	
	mean	stdv	mean	stdv	mean	stdv
Base	<u>17.2</u>	2.40	<u>8.85</u>	0.60	12.2	0.32
DARC1	1.30	0.07	1.35	0.02	0.96	0.01

(ND) = no data augmentation

Next, Noise robustness

Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

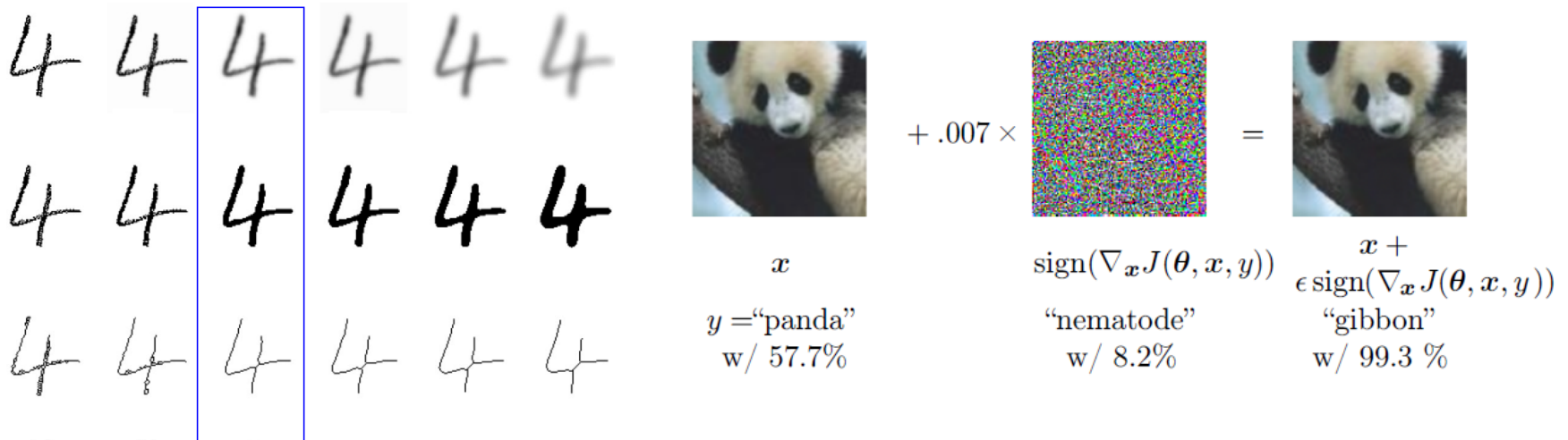
- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

- **Prior:** Most AI tasks have certain levels of **resilience on noise**
- One can incorporate such prior by **injecting noises** to the network



- Noise robustness is also related to **adversarial examples**
 - We will discuss this topic more in detail later

*sources :

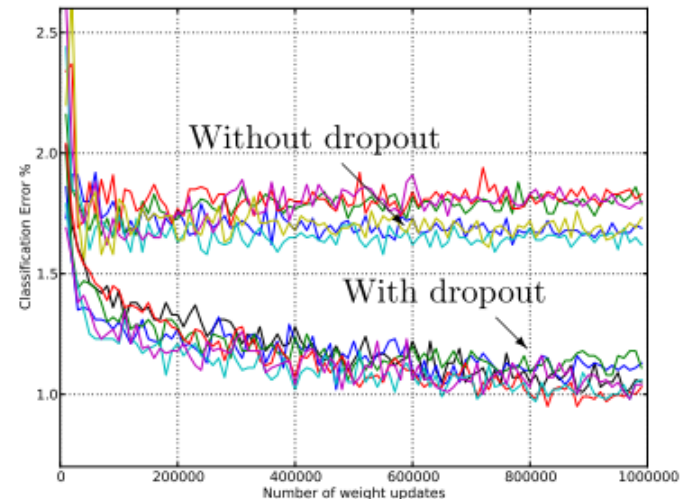
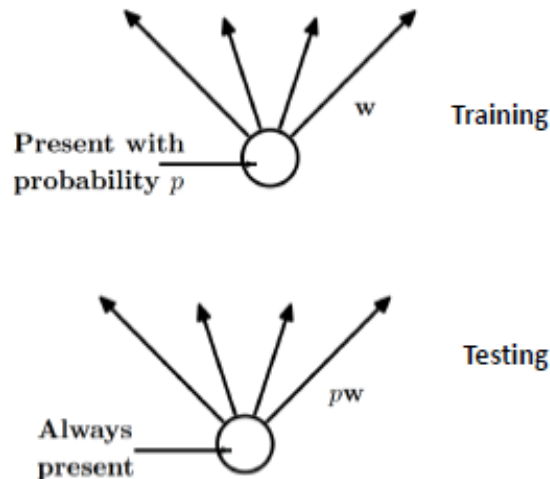
- Chatbri, Houssem et al. "Using scale space filtering to make thinning algorithms robust against noise in sketch images." Pattern Recognition Letters 42 (2014): 1-10.
- <https://www.deeplearningbook.org/contents/ml.html>

- **Prior:** Most AI tasks have certain levels of **resilience on noise**
- One can incorporate such prior by **injecting noises** to the network
- There can be many ways to impose noises:
 1. On **inputs** or **hidden units** (e.g. *Dropout*)
 - Noise with infinitesimal variance at the input is equivalent to imposing a penalty on the norm of the weights for some models [Bishop, 1995a,b]
 2. On **model parameters** (e.g. *Variational dropout*)
 - A stochastic implementation of a Bayesian inference over the weights
 3. On **gradients** during optimization (e.g. *Shake-shake regularization*)
 - In practice, SGD can generalize better than full GD in training DNNs [Keskar et al., 2016]

Next, Dropout

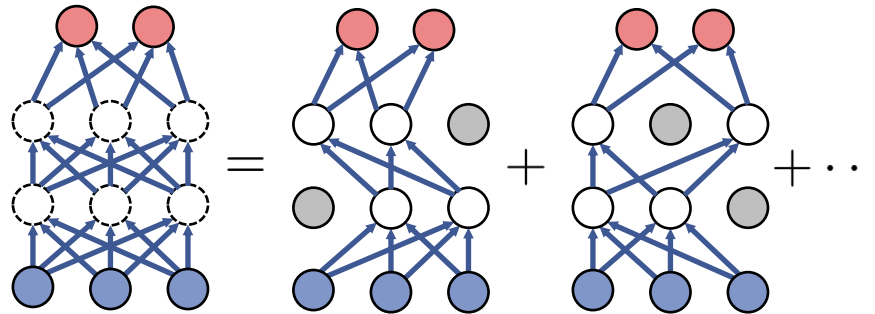
Noises on inputs or hidden units: Dropout [Srivastava et al., 2014]

- **Dropout** [Srivastava et al., 2014] randomly drops a neuron with probability p during training
 - Same as **multiplying a noise** $\mu \sim \text{Bernulli}(p)$ to each neuron
- At testing, each weights are scaled by p
- Dropout is applied to **hidden units** typically
 - Destruction of **high-level information** e.g. edges, nose, ...

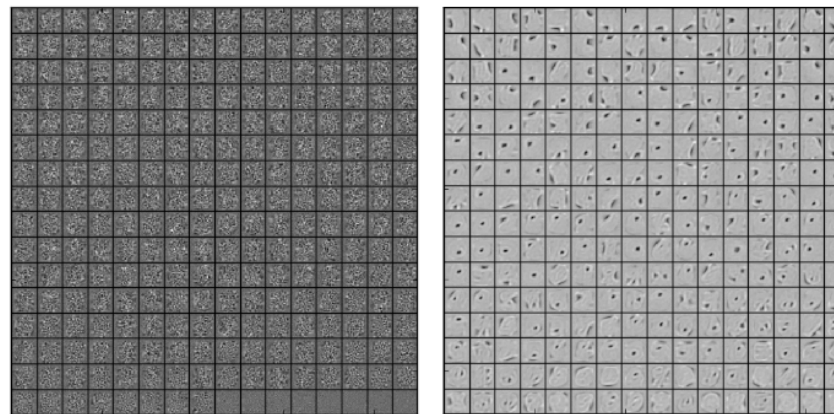


Why dropout generalizes well?

1. It can be thought of as **ensemble of 2^n subnets** with parameter sharing



2. Dropout prevents **co-adaptation of neurons**
 - Noisy neurons are **less reliable**
 - Each neuron must be prepared on which other neurons are dropped



(a) Without dropout

(b) Dropout with $p = 0.5$.

*source : Srivastava et al. "Dropout: A Simple Way to Prevent Neural Networks from Overfitting". JMLR 2014 46

The fully understanding on why dropout works is still an open question

- **Stochasticity might not be necessary**
 - **Fast dropout** [Wang et al., 2013]: A deterministic version of dropout with analytic marginalization
- **Dropout as an ensemble is not enough**
 - Dropout offers **additional improvements** to generalization error beyond those obtained by ensembles of independent models [Warde-Farley et al., 2013]
- **Dropping neurons are not necessary**
 - In principle, **any kind** of random modification is admissible
 - **Gaussian dropout**, i.e. $\mu \sim \mathcal{N}(1, \frac{1-p}{p})$, can work as well as the original dropout with probability p , or even work better

Next, Variational dropout

- In dropout, one have to find the best p manually
 - What if we want different rates for each of neurons?
- **Variational dropout** (VD) allows to **learn** the dropout rates separately
- Unlike Dropout, VD imposes noises on **model parameters** θ :

$$w_i := \theta_i \cdot \xi_i, \quad \text{where } p_{\alpha_i}(\xi_i) = \mathcal{N}(1, \alpha_i)$$

- A Bayesian generalization of **Gaussian dropout** [Srivastava et al., 2014]
- The random vector $\mathbf{w} = (w_i)_i$ is **adapted to data** in Bayesian sense by updating α and θ
- **Re-parametrization trick** allows \mathbf{w} to be learned via minibatch-based gradient estimation methods [Kingma et al., 2013]
 - α and θ can be “optimized” separated from noises

$$w_i = \theta_i + (\theta_i \sqrt{\alpha_i}) \cdot \varepsilon_i, \quad \text{where } \varepsilon_i \sim \mathcal{N}(0, 1)$$

Noises on model parameters: Variational dropout [Kingma et al., 2015]

- VD lead to a better model than dropout
- VD could also improve CNN as well, while dropout could not^(1b)

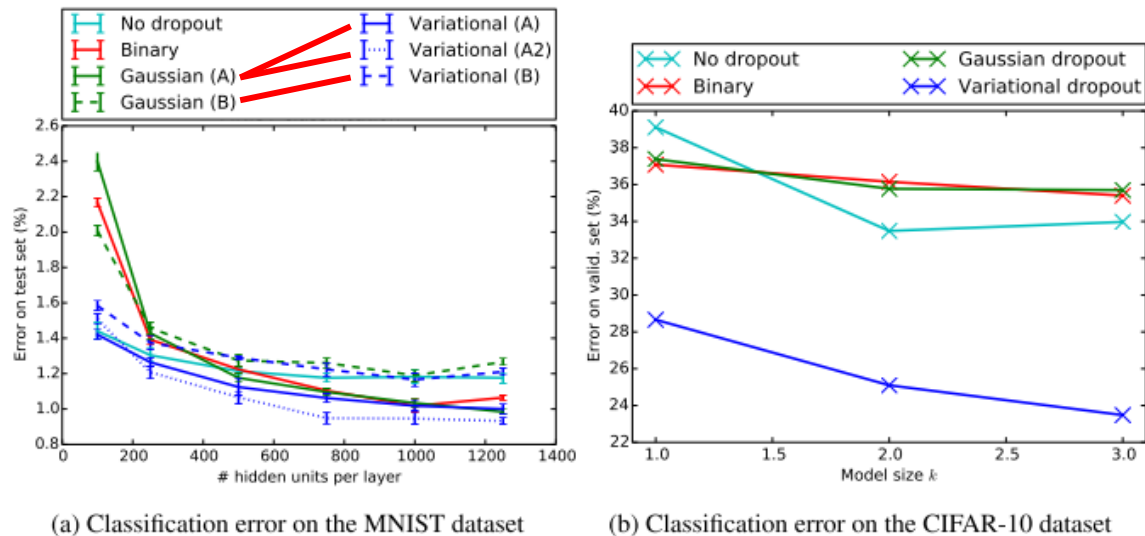
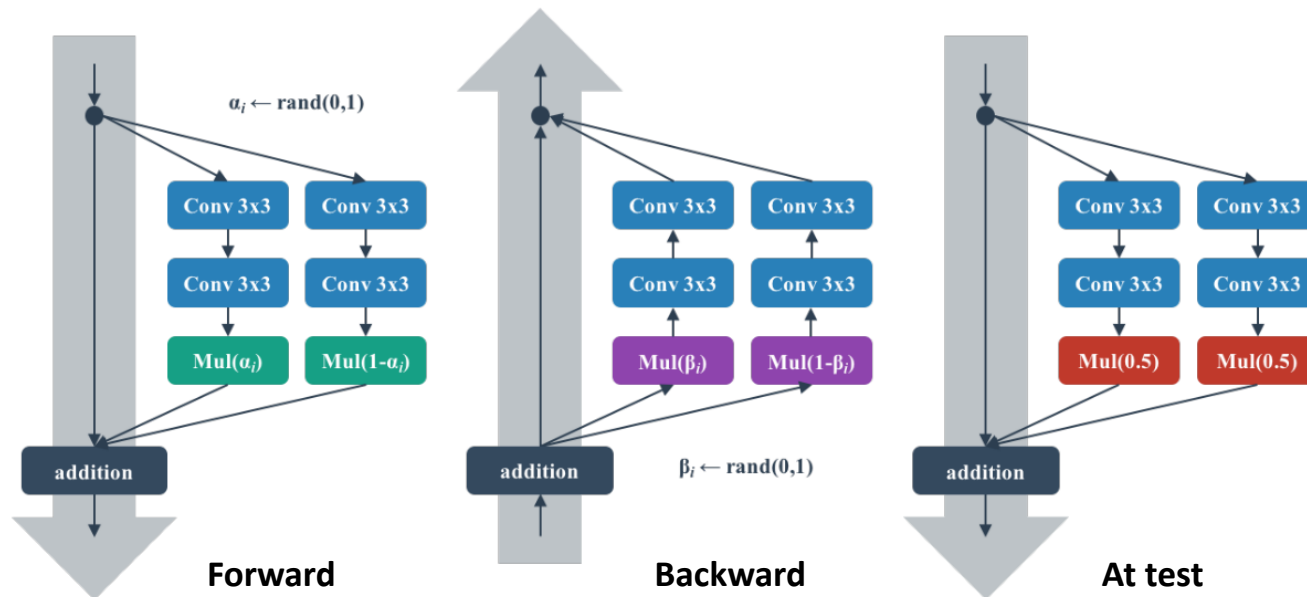


Figure 1: Best viewed in color. **(a)** Comparison of various dropout methods, when applied to fully-connected neural networks for classification on the MNIST dataset. Shown is the classification error of networks with 3 hidden layers, averaged over 5 runs. The variational versions of Gaussian dropout perform equal or better than their non-adaptive counterparts; the difference is especially large with smaller models, where regular dropout often results in severe underfitting. **(b)** Comparison of dropout methods when applied to convolutional net trained on the CIFAR-10 dataset, for different settings of network size k . The network has two convolutional layers with each $32k$ and $64k$ feature maps, respectively, each with stride 2 and followed by a softplus nonlinearity. This is followed by two fully connected layers with each $128k$ hidden units.

Next, Shake-shake regularization

- Noises can be injected even in **gradients** during back-propagation
- Shake-shake regularization** considers a 3-branch ResNeXt [Xie et al., 2017]



- Here, notice that α_i and β_i are independent random variables
 - α_i stochastically **blends** the outputs from two branches
 - β_i randomly **re-distributes** the returning gradient between two branches
- Those re-scaling are done in channel-wise

Noises on gradients: Shake-shake regularization [Gastaldi, 2017]

- Shake-shake shows one of the current state-of-the-art result on CIFAR-10/100

Method	Depth	Params	C10	C100
Wide ResNet	28	36.5M	3.8	18.3
ResNeXt-29, 16x64d	29	68.1M	3.58	17.31
DenseNet-BC (k=40)	190	25.6M	3.46	17.18
C10 Model S-S-I	26	26.2M	2.86	-
C100 Model S-E-I	29	34.4M	-	15.85

- Shake-shake reduces layer-wise correlations between two branches

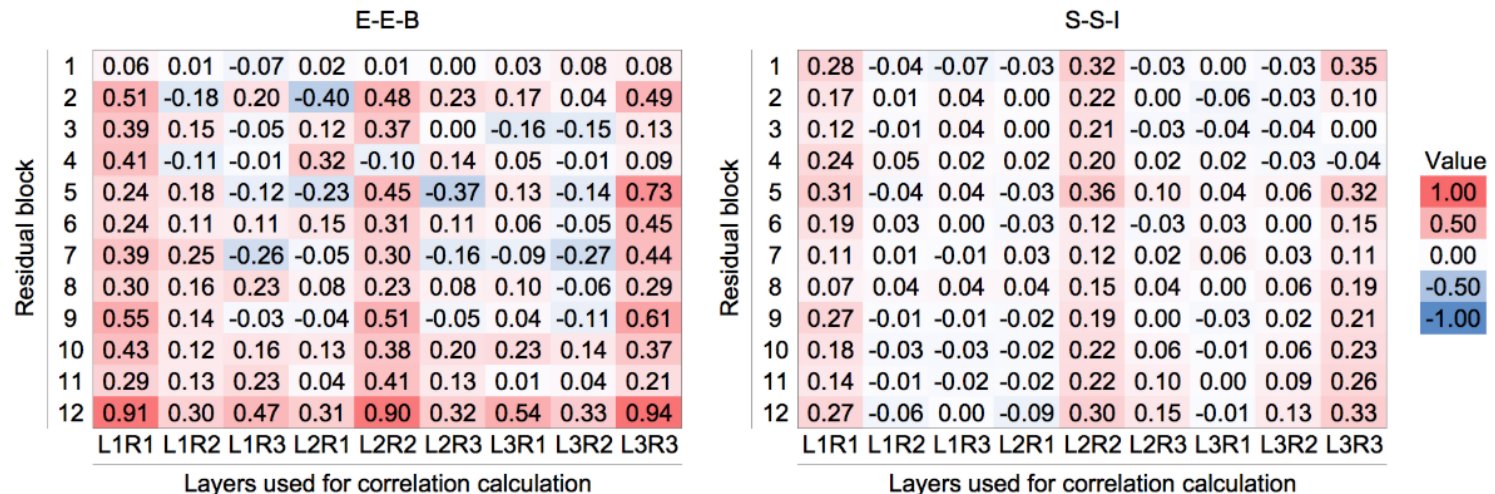


Table of Contents

1. Introduction

- Empirical risk minimization (ERM)

2. Stochastic Gradient Descent

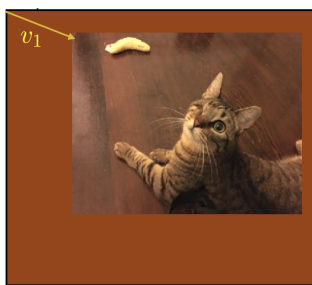
- Gradient descent (GD) and stochastic gradient descent (SGD)
- Momentum and adaptive learning rate methods

3. Regularization

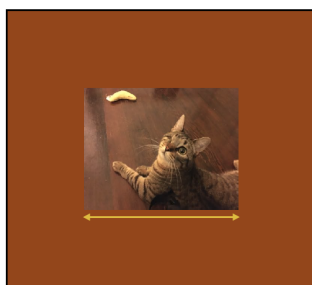
- Loss penalty with L2/L1/L0 norm
- Directly approximately regularizing complexity
- Noises on hidden units/gradients
- Data augmentations

4. Summary

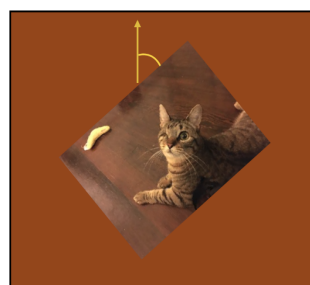
- **Prior:** The best way to generalize better is to gain more data
- Create fake data and add it to the training set
 - Requires some knowledge on **making good “fakes”**
- Particularly effective for classification tasks
 - Some tasks may not be readily applicable, e.g. density estimation
- Example: Rigid transformation symmetries
 - Translation, dilation, rotation, mirror symmetry, ...
 - Forms an affine group on pixels:
$$\begin{bmatrix} u_1 \\ u_2 \end{bmatrix} \mapsto \begin{bmatrix} v_1 \\ v_2 \end{bmatrix} + \begin{bmatrix} a_1 & a_2 \\ a_3 & a_4 \end{bmatrix} \begin{bmatrix} u_1 \\ u_2 \end{bmatrix}$$



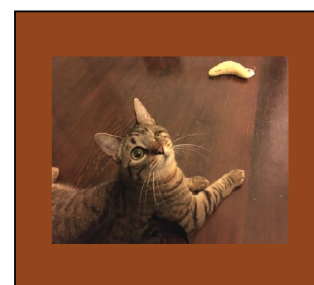
Translation



Dilation



Rotation



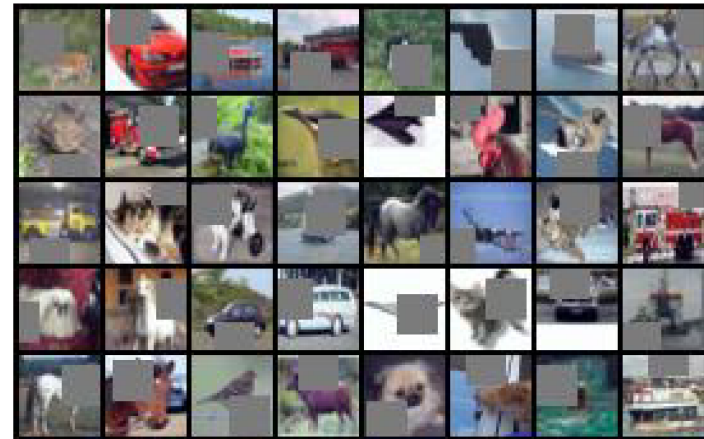
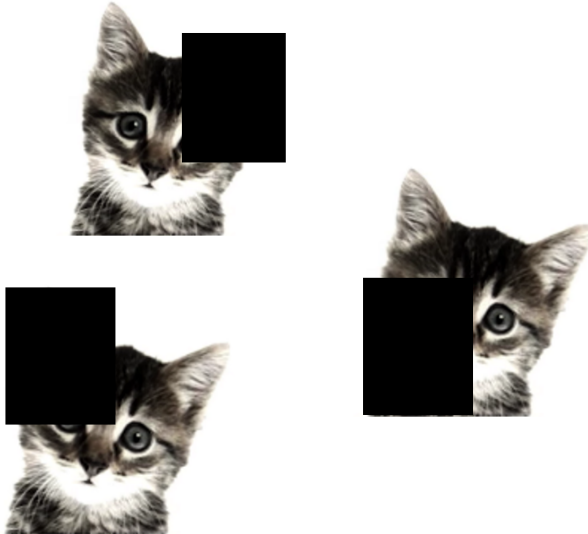
Mirror symmetry

Next, CutOut

- Dropout appears to be less powerful when used with **convolutional layers**
 - Dropping pixels randomly **may disturb gradients** due to parameter sharing
 - **Neighboring pixels in CNNs** would contains much of the dropped information
- Channel-wise dropout [Tompson et al., 2015] may alleviate these issues
 - However, the network capacity may be considerably reduced
- What do we expect by performing dropout on images?
 - Preventing **co-adaptation on high-level objects** (nose, eyes, ...)
 - For images, this can be also done by just using **local masking**



- What do we expect by performing dropout on images?
 - Preventing **co-adaptation on high-level objects** (nose, eyes, ...)
 - For images, this can be also done by just using **local masking**
- **CutOut** directly brings this into data augmentation
 - Data augmentation via **square-masking** randomly on images

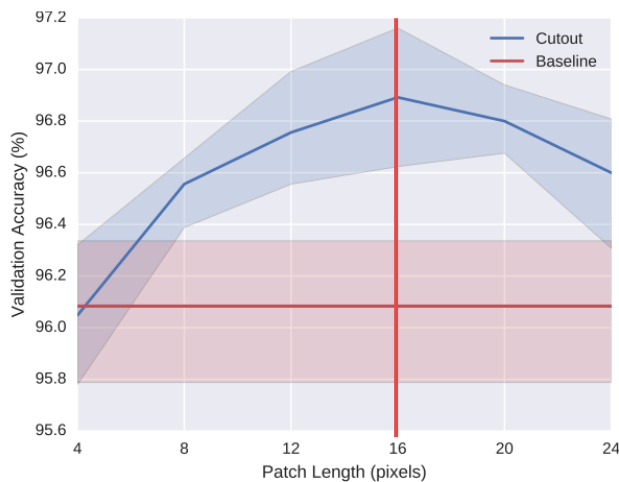


Making new data by local masking: CutOut [Devries et al., 2017]

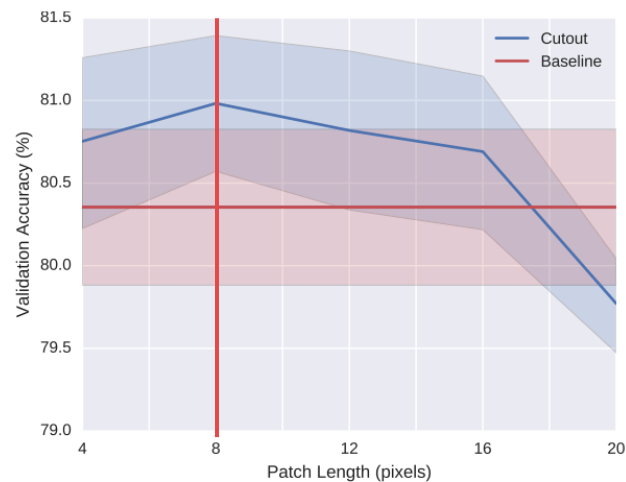
- CutOut further improved Shake-shake regularization [Gastaldi, 2017] achieving the state-of-the-art result on CIFAR-10/100

Method	C10	C10+	C100	C100+	SVHN
ResNet18 [5]	10.63 ± 0.26	4.72 ± 0.21	36.68 ± 0.57	22.46 ± 0.31	-
ResNet18 + cutout	9.31 ± 0.18	3.99 ± 0.13	34.98 ± 0.29	21.96 ± 0.24	-
WideResNet [22]	6.97 ± 0.22	3.87 ± 0.08	26.06 ± 0.22	18.8 ± 0.08	1.60 ± 0.05
WideResNet + cutout	5.54 ± 0.08	3.08 ± 0.16	23.94 ± 0.15	18.41 ± 0.27	1.30 ± 0.03
Shake-shake regularization [4]	-	2.86	-	15.85	-
Shake-shake regularization + cutout	-	2.56 ± 0.07	-	15.20 ± 0.21	-

- The size of the square should be set as a hyperparameter



(a) CIFAR-10



(b) CIFAR-100

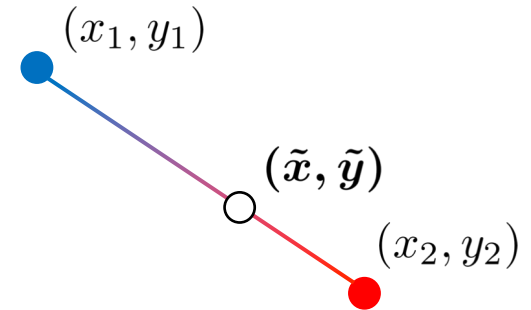
Next, Mixup

- In *mixup*, a new training example is constructed by:

$$\tilde{x} = \lambda x_1 + (1 - \lambda)x_2$$

$$\tilde{y} = \lambda y_1 + (1 - \lambda)y_2$$

- $\lambda \sim \text{Beta}(\alpha, \alpha) \in [0, 1]$, where α : hyperparameter
- (x_i, y_i) 's are uniformly sampled from the training data



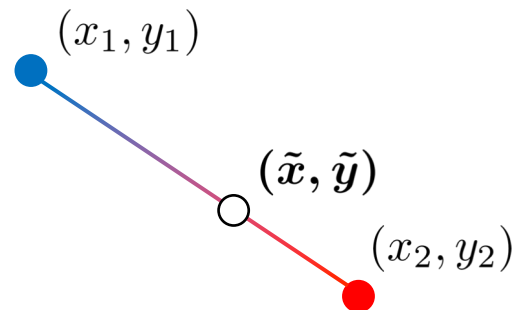
- Surprisingly, this simple scheme outperforms empirical risk minimization (ERM)
 - A new **state-of-art performance** on CIFAR-10/100 and ImageNet
 - Robustness when **learning from corrupt labels**
 - Handling adversarial examples
 - Stabilizing GANs
 - ...

- In *mixup*, a new training example is constructed by:

$$\tilde{x} = \lambda x_1 + (1 - \lambda)x_2$$

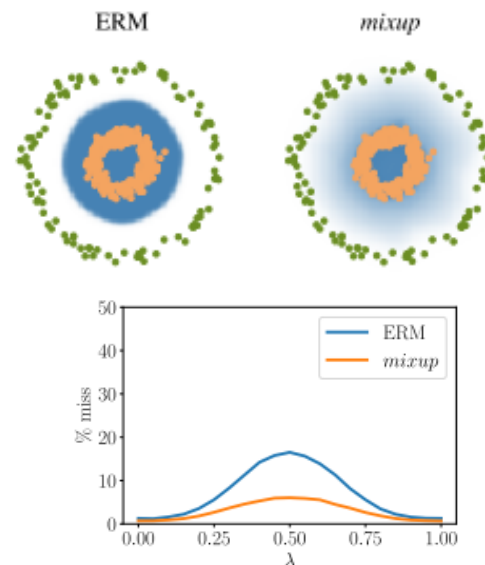
$$\tilde{y} = \lambda y_1 + (1 - \lambda)y_2$$

- $\lambda \sim \text{Beta}(\alpha, \alpha) \in [0, 1]$, where α : hyperparameter
- (x_i, y_i) 's are uniformly sampled from the training data



- What is *mixup* doing?**

- Incorporating **prior knowledge**: the model should behave **linearly** in-between training examples
- It reduces the amount of undesirable oscillations when predicting outside the training examples



(a) Prediction errors in-between training data. Evaluated at $x = \lambda x_i + (1 - \lambda)x_j$, a prediction is counted as a “miss” if it does not belong to $\{y_i, y_j\}$. The model trained with *mixup* has fewer misses.

Mixing two samples in dataset: *mixup* [Zhang et al., 2018]

- mixup* significantly improves generalization in CIFAR-10/100 and ImageNet

Dataset	Model	ERM	<i>mixup</i>
CIFAR-10	PreAct ResNet-18	5.6	3.9
	WideResNet-28-10	3.8	2.7
	DenseNet-BC-190	3.7	2.7
CIFAR-100	PreAct ResNet-18	25.6	21.1
	WideResNet-28-10	19.4	17.5
	DenseNet-BC-190	19.0	16.8

(a) Test errors for the CIFAR experiments.

Model	Method	Epochs	Top-1 Error	Top-5 Error
ResNet-50	ERM	200	23.6	7.0
	<i>mixup</i> $\alpha = 0.2$	200	22.1	6.1
ResNet-101	ERM	200	22.0	6.1
	<i>mixup</i> $\alpha = 0.2$	200	20.8	5.4
ResNeXt-101 32*4d	ERM	200	21.3	5.9
	<i>mixup</i> $\alpha = 0.4$	200	20.1	5.0

Table 1: Validation errors for ERM and *mixup* on the development set of ImageNet-2012.

- mixup* also shows robustness on corrupted labels while improving memorization [Zhang et al., 2016]

Label corruption	Method	Test error		Training error	
		Best	Last	Real	Corrupted
20%	ERM	12.7	16.6	0.05	0.28
	ERM + dropout ($p = 0.7$)	8.8	10.4	5.26	83.55
	<i>mixup</i> ($\alpha = 8$)	5.9	6.4	2.27	86.32
50%	ERM	18.8	44.6	0.26	0.64
	ERM + dropout ($p = 0.8$)	14.1	15.5	12.71	86.98
	<i>mixup</i> ($\alpha = 32$)	11.3	12.7	5.84	85.71
80%	ERM	36.5	73.9	0.62	0.83
	ERM + dropout ($p = 0.8$)	30.9	35.1	29.84	86.37
	<i>mixup</i> ($\alpha = 32$)	25.3	30.9	18.92	85.44

Summary

- SGD have been used as essential algorithms to deep learning as back-propagation.
 - Momentum methods improve the performance of gradient descend algorithms.
 - Annealing learning rates are critical for training loss functions
- In practice, **SGD + momentum** shows successful results, outperforming Adam!
 - For example, NLP (Huang et al., 2017) or machine translation (Wu et al., 2016)
- Reducing the test error, possibly at the expense of increased training error
 - No free lunch theorem says that there is **no best form of regularization**
- Developing effective regularizations is one of the major research in the field
- Nevertheless, as we are focusing on AI tasks, there could be some general strategies for a wide range of our problems
 - Loss penalty
 - Parameter sharing
 - Noise robustness
 - Dataset augmentation
 - ... **there can be many other ways!**

References

- [Nesterov' 1983] Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. 1983
link: <http://mpawankumar.info/teaching/cdt-big-data/nesterov83.pdf>
- [Duchi et al 2011], "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011
link : <http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>
- [Tieleman' 2012] Geoff Hinton's Lecture 6e of Coursera Class
link : http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- [Zeiler' 2012] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method
link : <https://arxiv.org/pdf/1212.5701.pdf>
- [Smith' 2015] Smith, Leslie N. "Cyclical learning rates for training neural networks."
link : <https://arxiv.org/pdf/1506.01186.pdf>
- [Kingma and Ba., 2015] Kingma and Ba. Adam: A method for stochastic optimization. ICLR 2015
link : <https://arxiv.org/pdf/1412.6980.pdf>
- [Dozat' 2016] Dozat, T. (2016). Incorporating Nesterov Momentum into Adam. ICLR Workshop,
link : http://cs229.stanford.edu/proj2015/054_report.pdf
- [Smith et al., 2017] Smith, Samuel L., Pieter-Jan Kindermans and Quoc V. Le. Don't Decay the Learning Rate, Increase the Batch Size. ICLR 2017.
link : <https://openreview.net/pdf?id=B1Yy1BxCZ>
- [Loshchilov et al., 2017] Loshchilov, I., & Hutter, F. (2017). SGDR: Stochastic Gradient Descent with Warm Restarts. ICLR 2017.
link : <https://arxiv.org/pdf/1608.03983.pdf>
- [Loshchilov et al., 2019] Loshchilov, I., & Hutter, F. (2019). Decoupled Weight Decay Regularization. ICLR 2019.
link : <https://arxiv.org/pdf/1711.05101.pdf>

References

- [Bishop, 1995a] Bishop, C. (1995). Regularization and Complexity Control in Feed-forward Networks. In *Proceedings International Conference on Artificial Neural Networks* (pp. 141–148).
link : <https://www.microsoft.com/en-us/research/publication/regularization-and-complexity-control-in-feed-forward-networks/>
- [Bishop, 1995b] Bishop, C. (1995). Training with Noise is Equivalent to Tikhonov Regularization. *Neural Computation*, 7, 108–116.
link : <https://ieeexplore.ieee.org/document/6796505/>
- [Wolpert et al., 1997] Wolpert, D. H., & Macready, W. G. (1997). No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation*, 1(1), 67–82.
link : <https://ieeexplore.ieee.org/document/585893/>
- [Hinton et al., 2012] Hinton, G. E., Srivastava, N., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. R. (2012). Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*.
link : <https://arxiv.org/abs/1207.0580>
- [Kingma et al., 2013] Kingma, D. P., & Welling, M. (2013). Auto-encoding variational bayes. *arXiv preprint arXiv:1312.6114*.
link : <https://arxiv.org/abs/1312.6114>
- [Wang et al., 2013] Wang, S., & Manning, C. (2013). Fast dropout training. In *Proceedings of the 30th International Conference on Machine Learning* (Vol. 28, pp. 118–126). Atlanta, Georgia, USA: PMLR.
link : <http://proceedings.mlr.press/v28/wang13a.html>
- [Warde-Farley et al., 2013] Warde-Farley, D., Goodfellow, I. J., Courville, A., & Bengio, Y. (2013). An empirical analysis of dropout in piecewise linear networks. *ArXiv Preprint ArXiv:1312.6197*.
link : <https://arxiv.org/abs/1312.6197>

References

- [Dauphin et al., 2014] Dauphin, Y. N., Pascanu, R., Gulcehre, C., Cho, K., Ganguli, S., & Bengio, Y. (2014). Identifying and attacking the saddle point problem in high-dimensional non-convex optimization. *Advances in Neural Information Processing Systems* 27 (pp. 2933–2941).
link : <https://papers.nips.cc/paper/5486-identifying-and-attacking-the-saddle-point-problem-in-high-dimensional-non-convex-optimization>
- [Goodfellow et al., 2014] Goodfellow, I. J., & Vinyals, O. (2014). Qualitatively characterizing neural network optimization problems. *CoRR*, *abs/1412.6544*.
link : <https://arxiv.org/abs/1412.6544>
- [Shalev-Shwartz et al., 2014] Shalev-Shwartz, S., & Ben-David, S. (2014). Understanding Machine Learning: From Theory to Algorithms. Cambridge: Cambridge University Press. doi:10.1017/CBO9781107298019
link : <https://www.cambridge.org/core/books/understanding-machine-learning/3059695661405D25673058E43C8BE2A6>
- [Srivastava et al., 2014] Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., & Salakhutdinov, R. (2014). Dropout: A Simple Way to Prevent Neural Networks from Overfitting. *Journal of Machine Learning Research*, 15, 1929–1958.
link : <http://jmlr.org/papers/v15/srivastava14a.html>
- [Tompson et al., 2015] Tompson, J., Goroshin, R., Jain, A., LeCun, Y., & Bregler, C. (2015). Efficient object localization using convolutional networks. In *Computer Vision and Pattern Recognition* (pp. 648–656).
link : <https://arxiv.org/abs/1411.4280>
- [Goodfellow et al., 2016] Goodfellow, I., Bengio, Y., & Courville, A. (2016). *Deep Learning*. MIT Press, pp.221-265.
link : <https://www.deeplearningbook.org/>
- [Kingma et al., 2015] Kingma, D. P., Salimans, T., & Welling, M. (2015). Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems* (pp. 2575-2583).
link : <https://papers.nips.cc/paper/5666-variational-dropout-and-the-local-reparameterization-trick>

References

- [Maddison et al. 2016] Maddison, C. J., Mnih, A., & Teh, Y. W. (2016). The concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations*.
link : <https://openreview.net/forum?id=S1jE5L5gl>
- [Keskar et al., 2016] Keskar, N. S., Mudigere, D., Nocedal, J., Smelyanskiy, M., & Tang, P. T. P. (2016). On large-batch training for deep learning: Generalization gap and sharp minima. In *International Conference on Learning Representations*.
link : <https://openreview.net/forum?id=H1oyRIYgg>
- [Zhang et al., 2016] Zhang, C., Bengio, S., Hardt, M., Recht, B., & Vinyals, O. (2016). Understanding deep learning requires rethinking generalization. *CoRR*, *abs/1611.03530*.
link : <https://arxiv.org/abs/1611.03530>
- [Ravanbakhsh et al., 2017] Ravanbakhsh, S., Schneider, J., & Póczos, B. (2017). Equivariance Through Parameter-Sharing. In *Proceedings of the 34th International Conference on Machine Learning* (Vol. 70, pp. 2892–2901). International Convention Centre, Sydney, Australia: PMLR.
link : <http://proceedings.mlr.press/v70/ravanbakhsh17a.html>
- [Devries et al., 2017] Devries, T., & Taylor, G. W. (2017). Improved Regularization of Convolutional Neural Networks with Cutout. *CoRR*, *abs/1708.04552*. Retrieved from <http://arxiv.org/abs/1708.04552>
link : <https://arxiv.org/abs/1708.04552>
- [Gastaldi, 2017] Gastaldi, X. (2017). Shake-Shake regularization. *CoRR*, *abs/1705.07485*.
link : <http://arxiv.org/abs/1705.07485>
- [Kawaguchi et al., 2017] Kawaguchi, K., Kaelbling, L. P., & Bengio, Y. (2017). Generalization in deep learning. arXiv preprint arXiv:1710.05468.
link : <https://arxiv.org/abs/1710.05468>

References

- [Pereyra et al., 2017] Pereyra, G., Tucker, G., Chorowski, J., Kaiser, Ł., & Hinton, G. (2017). Regularizing neural networks by penalizing confident output distributions. arXiv preprint arXiv:1701.06548.
link : <https://arxiv.org/abs/1701.06548>
- [Xie et al., 2017] Xie, S., Girshick, R., Dollár, P., Tu, Z., & He, K. (2017). Aggregated residual transformations for deep neural networks. In *Computer Vision and Pattern Recognition* (pp. 5987–5995).
link : http://openaccess.thecvf.com/content_cvpr_2017/papers/Xie_Aggregated_Residual_Transformations_CVPR_2017_paper.pdf
- [Arora et al., 2018] Arora, S., Cohen, N., & Hazan, E. (2018). On the Optimization of Deep Networks: Implicit Acceleration by Overparameterization. In *Proceedings of the 35th International Conference on Machine Learning* (Vol. 80, pp. 244–253). Stockholmsmässan, Stockholm Sweden: PMLR.
link : <http://proceedings.mlr.press/v80/arora18a.html>
- [Hartford et al., 2018] Hartford, J., Graham, D., Leyton-Brown, K., & Ravanbakhsh, S. (2018). Deep Models of Interactions Across Sets. In *Proceedings of the 35th International Conference on Machine Learning* (Vol. 80, pp. 1909–1918). Stockholmsmässan, Stockholm Sweden: PMLR.
Link : <http://proceedings.mlr.press/v80/hartford18a.html>
- [Louizos et al., 2018] Louizos, C., Welling, M., & Kingma, D. P. (2018). Learning Sparse Neural Networks through L₀ Regularization. In *International Conference on Learning Representations*.
link : <https://openreview.net/forum?id=H1Y8hhg0b>
- [Zhang et al., 2018] Zhang, H., Cisse, M., Dauphin, Y. N., & Lopez-Paz, D. (2018). *mixup*: Beyond Empirical Risk Minimization. In *International Conference on Learning Representations*.
link : <https://openreview.net/forum?id=r1Ddp1-Rb>

References

- [Nesterov' 1983] Nesterov. A method of solving a convex programming problem with convergence rate $O(1/k^2)$. 1983
link: <http://mpawankumar.info/teaching/cdt-big-data/nesterov83.pdf>
- [Duchi et al 2011], "Adaptive subgradient methods for online learning and stochastic optimization", JMLR 2011
link : <http://www.jmlr.org/papers/volume12/duchi11a/duchi11a.pdf>
- [Tieleman' 2012] Geoff Hinton's Lecture 6e of Coursera Class
link : http://www.cs.toronto.edu/~tijmen/csc321/slides/lecture_slides_lec6.pdf
- [Zeiler' 2012] Zeiler, M. D. (2012). ADADELTA: An Adaptive Learning Rate Method
link : <https://arxiv.org/pdf/1212.5701.pdf>
- [Smith' 2015] Smith, Leslie N. "Cyclical learning rates for training neural networks."
link : <https://arxiv.org/pdf/1506.01186.pdf>
- [Kingma and Ba., 2015] Kingma and Ba. Adam: A method for stochastic optimization. ICLR 2015
link : <https://arxiv.org/pdf/1412.6980.pdf>
- [Dozat' 2016] Dozat, T. (2016). Incorporating Nesterov Momentum into Adam. ICLR Workshop,
link : http://cs229.stanford.edu/proj2015/054_report.pdf
- [Smith et al., 2017] Smith, Samuel L., Pieter-Jan Kindermans and Quoc V. Le. Don't Decay the Learning Rate, Increase the Batch Size. ICLR 2017.
link : <https://openreview.net/pdf?id=B1Yy1BxCZ>
- [Loshchilov et al., 2017] Loshchilov, I., & Hutter, F. (2017). SGDR: Stochastic Gradient Descent with Warm Restarts. ICLR 2017.
link : <https://arxiv.org/pdf/1608.03983.pdf>