

# Stochastic Optimization and Training Skills in Deep Learning

Wenjun Zeng

Reading Group of Ye Lab, UM

Sept. 16, 2018

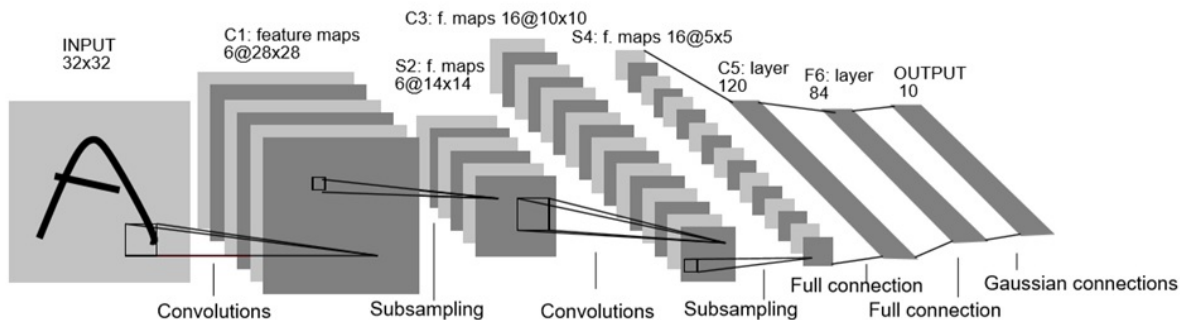
# Contents

- Architecture of Deep Neural Networks
- Back Propagation
- Stochastic Optimization Algorithms
  - ◆ Stochastic Gradient/Subgradient Descent
  - ◆ Accelerated Methods
  - ◆ Variance Reduction
- Training Skills
  - ◆ Batch Normalization
  - ◆ Dropout
  - ◆ Parameter Tuning (Alchemy?)...

# I. Architecture of Deep Neural Networks (DNN)

Take convolutional neural networks (CNN) as example.

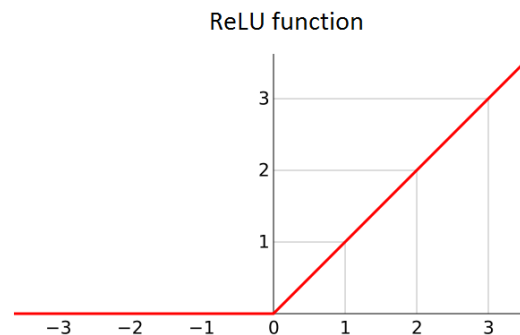
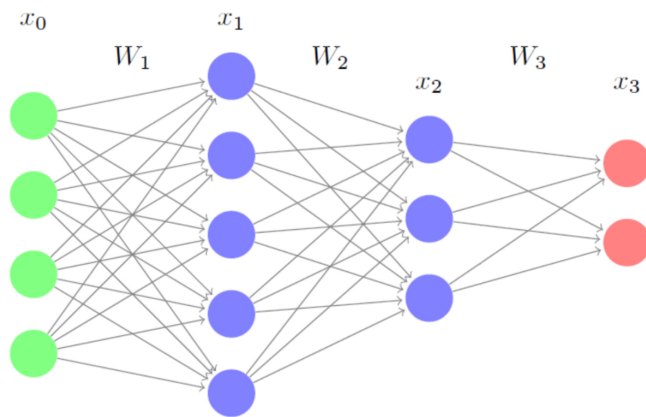
LeNet [LeCun *et al.* '1998]



- ☐ Convolutional layer
- ☐ Pooling layer (subsampling)
- ☐ Full connection layer
- ☐ Nonlinear mapping is sometimes also called as a layer (ReLU layer)

DNN: many (from several tens to nearly 200) layers.

## 1.1 Full Connection Layer



□  $\mathbf{x}_l \in \mathbb{R}^{n_l}$  ( $l = 0, 1, \dots, L$ ): vector of  $l$ -th layer;  $L$ : number of layers

□  $\mathbf{x}_0$  is input;  $\mathbf{x}_L$  is output; other  $\mathbf{x}_l$  are hidden layer variables.

$$\mathbf{u}_l = \mathbf{W}_l \mathbf{x}_{l-1} + \mathbf{b}_l$$

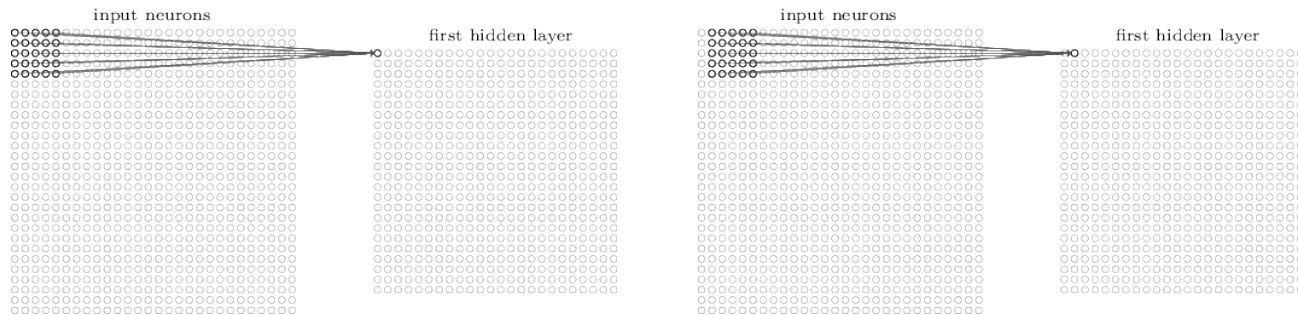
$$\mathbf{x}_l = \sigma(\mathbf{u}_l)$$

□  $\mathbf{W}_l \in \mathbb{R}^{n_l \times n_{l-1}}$  and  $\mathbf{b}_l \in \mathbb{R}^{n_l}$  are **weights** and **bias**.

□  $\sigma(\cdot)$  is nonlinear **activation function**, e.g., sigmoid or ReLU.

## 1.2 Convolutional Layer

Convolution layer in fact has similar structure as the full connection layer but a neuron of the next layer is merely connected with very few neurons of the last layer.



**Local receptive fields:** only make connections in small, localized regions of the input.

$$x_l(i, j) = \sigma \left( b + \sum_{p=0}^{FH} \sum_{q=0}^{FW} w_{p,q} x_{l-1}(i + p, j + q) \right)$$

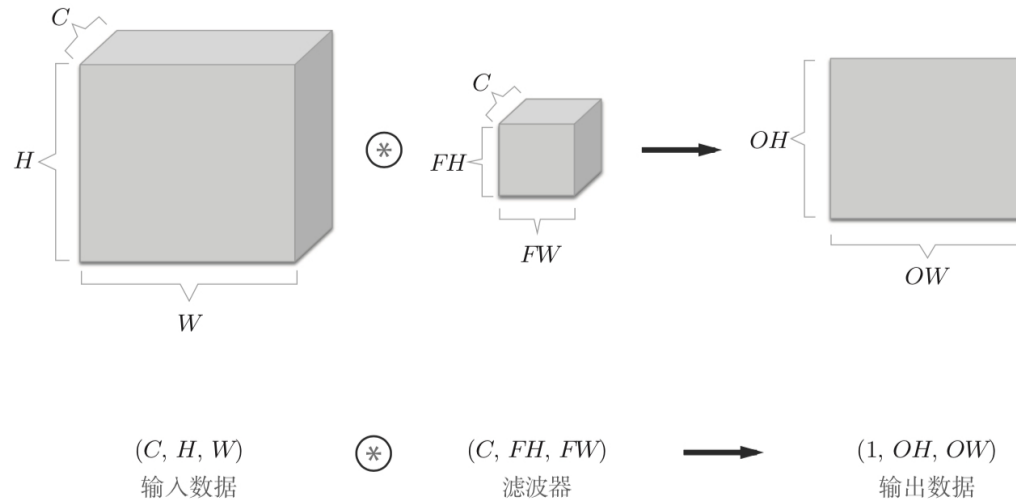
$\{w_{p,q}\}$  are coefficients of convolutional kernel (**filter**) with size  $FH \times FW$ .

**Shared weights and biases:** same convolution kernel applies to all pixels of the input.

Number of parameters significantly reduced compared with full connection case.

## Tensor in Convolutional Layer

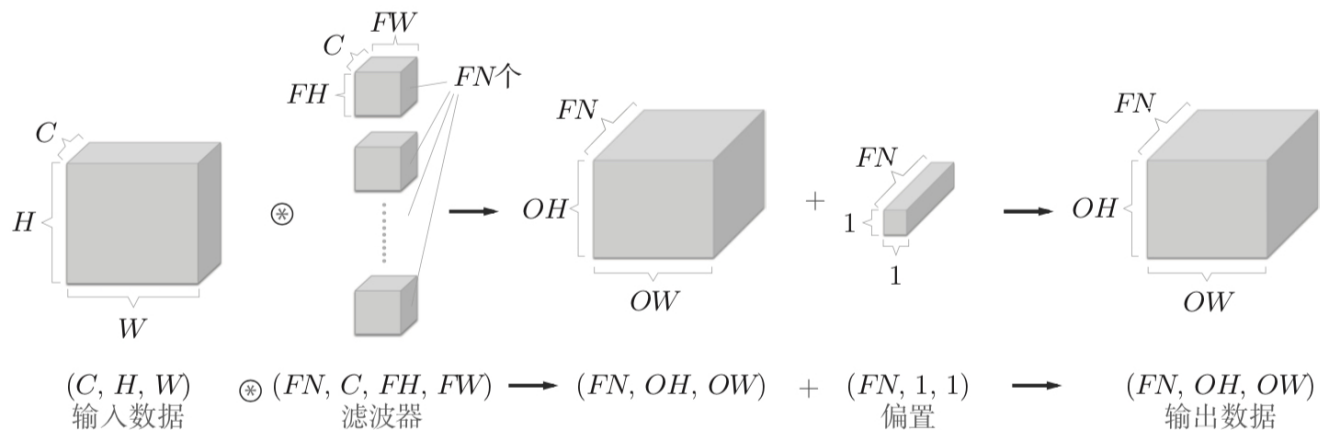
Using multiple filters, data of each convolutional layer are 3rd-order **tensor**.



$C$ : number of channels

- ❑ At the input layer,  $C = 3$  for colored images (RGB) and  $C = 1$  for grayed images.
- ❑ At other convolutional layers,  $C$  equals the number of convolutional kernels.
- ❑ The output corresponding to **one** filter is a matrix (**feature map**), not a tensor.

For  $FN$  convolutional kernels



Number of weights is  $FH \cdot FW \cdot C \cdot FN$  while number of bias is  $FN$ .

Total number of parameters  $(FH \cdot FW \cdot C + 1) \cdot FN \ll (n_l + 1)n_{l-1}$

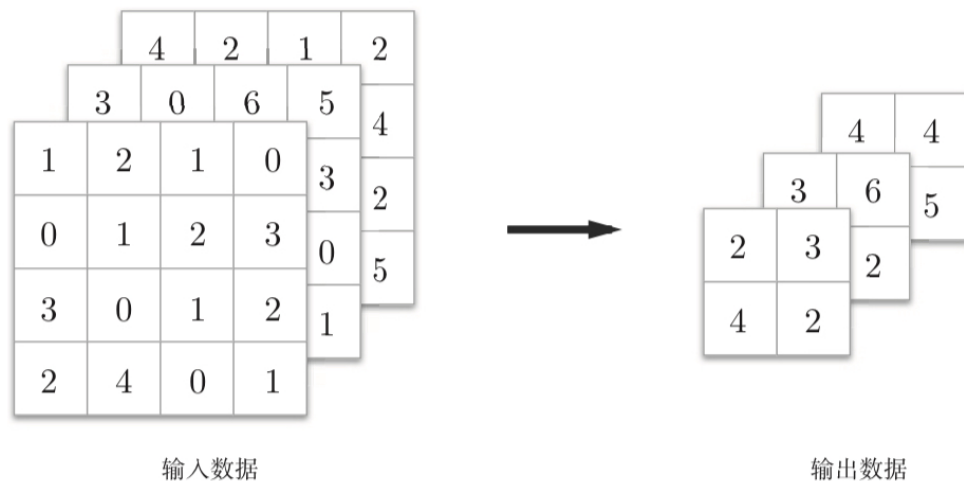
$$OH = \frac{H + 2P - FH}{S} + 1$$

$$OW = \frac{W + 2P - FW}{S} + 1$$

$P$  is number of padding and  $S$  is stride.

## 1.3 Pooling Layer

Pooling is nothing but subsampling. We prefer max pooling.



- ❑ No parameters need to learn in pooling layers (pay attention when back propagation).
- ❑ Number of channels remain unchanged.

Pooling is important because it makes DNN robust to slight variation of the data.



## II. Back Propagation (BP)

*D. E. Rumelhart, G. E. Hinton, R. J. Williams, “Learning representations by back-propagating errors,” **Nature**, 323, 533–536.*

- ❑ I think BP is the most important issue for NN and DNN. Without it, we cannot efficiently train an NN.
- ❑ It is based on the **chain rule** of multivariate derivative.
- ❑ Many courses, books, and papers present BP in a very long, tedious, puzzling, and even unclear way...

One derivation by LeCun [LeCun’1988] using the **Lagrange multiplier** method may be better.

First, once  $\{\mathbf{W}_l, \mathbf{b}_l\}$  is initialized or updated, do **forward pass**

for  $l = 1, 2 \dots, L$

$$\mathbf{u}_l = \mathbf{W}_l \mathbf{x}_{l-1} + \mathbf{b}_l$$

$$\mathbf{x}_l = \sigma(\mathbf{u}_l)$$

end for

## 2.1 Training NN

□ One (or mini-batch) training sample  $\mathbf{x}_0$  and the desired target (e.g., label)  $\mathbf{y} \in \mathbb{R}_L^n$

□ Loss function, e.g., the  $\ell_2$ -loss:

$$f(\{\mathbf{W}_l, \mathbf{b}_l\}) = \frac{1}{2} \|\mathbf{x}_L - \mathbf{y}\|^2$$

We use stochastic gradient descent (**SGD**) (discuss soon)

$$\begin{aligned}\mathbf{W}_l^{k+1} &= \mathbf{W}_l^k - \alpha_k \frac{\partial f}{\partial \mathbf{W}_l^k} \\ \mathbf{b}_l^{k+1} &= \mathbf{b}_l^k - \eta_k \frac{\partial f}{\partial \mathbf{b}_l^k}\end{aligned}$$

where  $\alpha_k > 0$  and  $\eta_k > 0$  are **learning rates**.

Question: how to compute the gradients  $\frac{\partial f}{\partial \mathbf{W}_l}$  and  $\frac{\partial f}{\partial \mathbf{b}_l}$ ? BP solves it.

## 2.2 Summary of BP

1. For the last layer ( $l = L$ )

$$\frac{\partial f}{\partial \mathbf{u}_L} = (\mathbf{x}_L - \mathbf{y}) \odot \sigma'(\mathbf{u}_L)$$

2. For  $l = L - 1, L - 2, \dots, 1$  (**backward pass**)

$$\frac{\partial f}{\partial \mathbf{u}_l} = \mathbf{W}_l^T \frac{\partial f}{\partial \mathbf{u}_{l+1}} \odot \sigma'(\mathbf{u}_l)$$

where  $\odot$  is component-wise product and  $\sigma'(\cdot)$  is the derivative of  $\sigma(\cdot)$ .

Then, for all  $l = L, L - 1, \dots, 1$

$$\begin{aligned} \frac{\partial f}{\partial \mathbf{W}_l} &= \frac{\partial f}{\partial \mathbf{u}_l} \mathbf{x}_l^T \\ \frac{\partial f}{\partial \mathbf{b}_l} &= \frac{\partial f}{\partial \mathbf{u}_l} \end{aligned}$$

For convolutional layer, BP is slightly different from the full connection case. But the principle is the same.

### III. Stochastic Optimization Algorithms

Many (almost all?) deep learning libraries use **stochastic gradient** methods to train DNN.

- ❑ **Caffe**: SGD, Adadelata, Adagrad, Adam, Nesterov, Rmsprop
- ❑ **TensorFlow**: SGD, Adadelata, AdagradDA, Adagrad, ProximalAdagrad, Ftrl, Momentum, Adam, CenteredRMSProp.
- ❑ **Lasagne**: SGD, Momentum, Nesterov, Adagrad, Rmsprop, Adadelata, Adam, Adamax...

*Stochastic Optimization Algorithms is crucial for big DNN training.*

### 3.1 Big- $N$ Problem

Consider **finite sum** or empirical risk minimization (**ERM**)

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) \right\} \quad (1)$$

□  $\mathbf{x} \in \mathbb{R}^d$  denotes parameters of neural networks, i.e.,  $\{\mathbf{W}_l, \mathbf{b}_l\}$ ;  $d$  can be very large.

□  $N$ : number of training samples; usually very large, e.g., million.

□  $f_i(\mathbf{x})$ : loss function due to sample  $i$ .

Sometimes **regularization** is used to suppress over-fitting

$$\min_{\mathbf{x} \in \mathbb{R}^d} \left\{ f(\mathbf{x}) := \frac{1}{N} \sum_{i=1}^N f_i(\mathbf{x}) + \lambda r(\mathbf{x}) \right\} \quad (2)$$

Common regularizer used in Caffe and TensorFlow  $r(\mathbf{x}) = \|\mathbf{x}\|_2^2$  and  $r(\mathbf{x}) = \|\mathbf{x}\|_1$ .

$\Rightarrow$  alternative to **weight decay**.

## 3.2 Gradient Method

Deterministic gradient descent (GD) [Cauchy' 1847]:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f(\mathbf{x}^k) \quad \text{with } \nabla f(\mathbf{x}^k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}^k) \quad (3)$$

Well-established rules to determine **learning rate**  $\alpha_k$ :

- $0 < \alpha_k < \frac{2}{L}$ , e.g.,  $\alpha_k = \frac{1}{L}$ , with  $L$  the **Lipschitz constant** of  $\nabla f(\mathbf{x})$ .
- $\alpha_k$  is determined by **backtracking line search**, e.g., Armijo's rule [Bertsekas' 1999].

**Convergence Guarantees** [Nestrov' 2004]

- Linear convergence rate of  $\mathcal{O}((1 - \mu/L)^k)$  for  $\mu$ -strongly convex objective  $f(\mathbf{x})$ .
- Sublinear rate of  $\mathcal{O}(1/k)$  without strong convexity.

**Drawback:** Too expensive to compute  $\nabla f(\mathbf{x}^k) = \frac{1}{N} \sum_{i=1}^N \nabla f_i(\mathbf{x}^k)$  for big  $N$ .

### 3.3 Stochastic Gradient Method

Stochastic gradient descent (SGD) [Robbins-Monro' 1951]:

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \alpha_k \nabla f_{i_k}(\mathbf{x}^k) \quad (4)$$

where  $i_k$  is selected **uniformly at random** from  $\{1, 2, \dots, N\}$ .

**Motivation:**

$$\mathbb{E} [\nabla f_{i_k}(\mathbf{x}^k)] = \sum_{i=1}^N \frac{1}{N} \nabla f_i(\mathbf{x}^k) = \nabla f(\mathbf{x}^k)$$

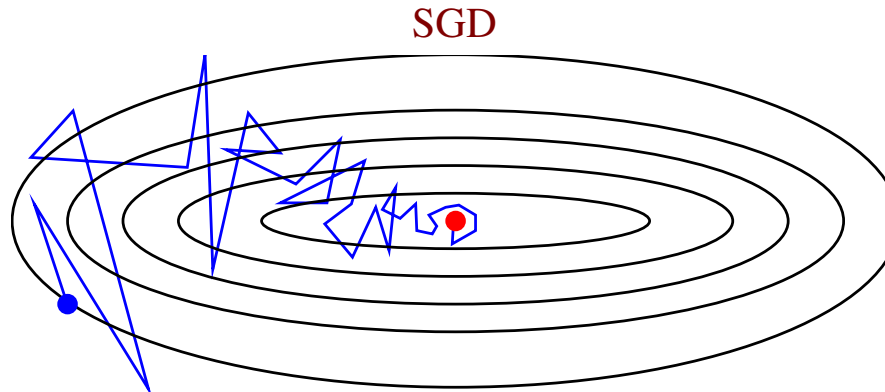
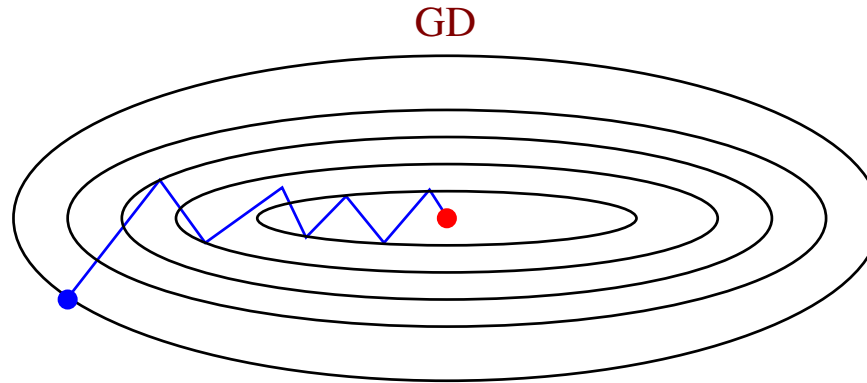
$\Rightarrow \nabla f_{i_k}(\mathbf{x}^k)$  is a cheap unbiased estimate of true gradient

□ Per-iteration cost of SGD is  $N$  times lower than GD.

□ Convergence requires  $\alpha_k \rightarrow 0 \Rightarrow$  *line search not applicable—learning rate needs to tune!*

Typical settings:  $\alpha_k \sim \mathcal{O}(1/k)$  or  $\alpha_k \sim \mathcal{O}(1/\sqrt{k})$ .

## Convergence Behavior



The objective of SGD does not monotonically decrease.



## Theoretical Convergence of SGD

Assumptions:

- $\nabla f(\mathbf{x})$  is Lipschitz continuous with constant  $L$  (DNN loss function locally satisfies)
- Gradient is bounded:  $\mathbb{E}_i [\|\nabla f_i(\mathbf{x})\|^2] \leq M^2$ .

**Theorem** Define  $\Delta_k = \|\mathbf{x} - \mathbf{x}^*\|$ . For  $\mu$ -strongly convex objective and a diminishing stepsize  $\alpha_k = \frac{\beta}{k+\gamma}$  for some  $\beta > \frac{1}{2\mu}$  and  $\gamma > 0$  such that  $\alpha_1 \leq \frac{1}{2\mu}$ . Then, for any  $k \geq 1$

$$\mathbb{E}[f(k) - f(\mathbf{x}^*)] \leq \frac{L}{2} \mathbb{E}[\Delta_k^2] \leq \frac{L}{2} \frac{\nu}{k + \gamma}$$

where  $\nu = \max\left(\frac{\beta^2 M^2}{2\beta\mu - 1}, (\gamma + 1)\Delta_1^2\right)$ .

- For fixed stepsize, we don't have convergence.
- For diminishing stepsize, convergence rate is  $\mathcal{O}(1/k)$ .
- For nonsmooth case (e.g., ReLU), convergence rate of **stochastic subgradient** is  $\mathcal{O}(1/\sqrt{k})$ .

## 3.4 Accelerated Methods

### SGD with Momentum

$$\begin{aligned}\mathbf{v}^{k+1} &= \mu_k \mathbf{v}^k - \alpha_k \nabla f_{i_k}(\mathbf{x}^k) \\ \mathbf{x}^{k+1} &= \mathbf{x}^k + \mathbf{v}^{k+1}\end{aligned}$$

### Nesterov Accelerated Gradient

Yurii Nesterov



Y. Nesterov, “A method for solving the convex programming problem with convergence rate  $O(1/k^2)$ ,” *Dokl. Acad. Nauk SSSR*, pp. 543–547, 1983.

## Original Version of Nesterov

$$\mathbf{v}^{k+1} = (1 + \mu_k)\mathbf{x}^k - \mu_k\mathbf{x}^{k-1}$$

$$\mathbf{x}^{k+1} = \mathbf{v}^{k+1} - \alpha_k \nabla f_{i_k}(\mathbf{v}^{k+1})$$

where  $\mu_k = \frac{k+2}{k+5}$ .

## Momentum Version of Nesterov

$$\mathbf{v}^{k+1} = \mu_k\mathbf{v}^k - \alpha_k \nabla f_{i_k}(\mathbf{x}^k + \mu_k\mathbf{v}^k)$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \mathbf{v}^{k+1}$$

- For GD, convergence rate of Nesterov is improved from  $\mathcal{O}(1/k)$  to  $\mathcal{O}(1/k^2)$ .  $\mathcal{O}(1/k^2)$  is the **optimal bound** for non-strongly convex objective of all first-order methods! Recently still many researches on it. [Su-Boyd-Candès'2016], [Bubeck-Lee-Singh'2015]
- For stochastic optimization, is there any theoretical analysis on the improvement of convergence rate of Nesterov's acceleration? I don't find any at present...

## Adaptive Subgradient (Adagrad) [Duchi-Hazan-Singer'2011]

Let  $\mathbf{g}^k = \nabla f_{i_k}(\mathbf{x}^k)$

$$[\mathbf{x}^{k+1}]_n = [\mathbf{x}^k]_n - \frac{\alpha_k}{\sqrt{\sum_{j=1}^k [\mathbf{g}^j]_n^2 + \epsilon}} [\mathbf{g}^k]_n, \quad n = 1, 2, \dots, d$$

where  $\epsilon > 0$  is small to avoid being divided by zero.

Drawback: denominator becomes larger and larger  $\Rightarrow$  learning rate approaches 0 quickly.

*Adagrad greatly improved the robustness of SGD and used it for training large-scale neural nets at Google and for recognizing cats in Youtube videos.* [Dean, Ng, et al. 2012]

## RMSProp [Hinton]

Initialize  $\mathbb{E}[(\mathbf{g}^0)^2] = \mathbf{0}$ . At step  $k = 1, 2, \dots$

$$\begin{aligned} \mathbb{E}[(\mathbf{g}^k)^2] &= \rho \mathbb{E}[(\mathbf{g}^{k-1})^2] + (1 - \rho)(\mathbf{g}^k)^2 \\ \mathbf{x}^{k+1} &= \mathbf{x}^k - \frac{\alpha_k}{\sqrt{\mathbb{E}[(\mathbf{g}^k)^2] + \epsilon \mathbf{1}_d}} \nabla f_{i_k}(\mathbf{x}^k) \end{aligned}$$

where  $\rho = 0.9$  (like **forgetting factor** in adaptive filtering) and  $\alpha_k = 10^{-3}$  is a good choice.

The root mean square (**RMS**) of  $\mathbf{g}^k$  is denoted as  $\text{RMS}[\mathbf{g}^k] = \sqrt{\mathbb{E}[(\mathbf{g}^k)^2] + \epsilon \mathbf{1}_d}$ .

## Adadelta [Zeiler'2012]

The “delta” in Adadelta means the quantity  $\Delta \mathbf{x}^k = \mathbf{x}^{k+1} - \mathbf{x}^k$ .

Denote  $\text{RMS}[\mathbf{g}^k] = \sqrt{\mathbb{E}[(\mathbf{g}^k)^2] + \epsilon \mathbf{1}_d}$  and  $\text{RMS}[\Delta \mathbf{x}^k] = \sqrt{\mathbb{E}[(\Delta \mathbf{x}^k)^2] + \epsilon \mathbf{1}_d}$ .

Initialize  $\mathbb{E}[(\mathbf{g}^0)^2] = \mathbf{0}$  and  $\mathbb{E}[(\Delta \mathbf{x}^0)^2] = \mathbf{0}$ .

At step  $k = 1, 2, \dots$

$$\mathbb{E}[(\mathbf{g}^k)^2] = \rho \mathbb{E}[(\mathbf{g}^{k-1})^2] + (1 - \rho)(\mathbf{g}^k)^2$$

$$\Delta \mathbf{x}^k = -\frac{\text{RMS}[\Delta \mathbf{x}^{k-1}]}{\text{RMS}[\mathbf{g}^k]} \mathbf{g}^k$$

$$\mathbb{E}[(\Delta \mathbf{x}^k)^2] = \rho \mathbb{E}[(\Delta \mathbf{x}^{k-1})^2] + (1 - \rho)(\Delta \mathbf{x}^k)^2$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k + \Delta \mathbf{x}^k$$

where  $\rho$  is around 0.9.

Advantage: Adadelta does not need to set a learning rate.

## Adaptive Moment Estimation (Adam) [Kingma-Ba'2015]

Idea: Combine Adadelta/RMSProp and momentum.

**Moment:** First and second moments of (component-wise) gradient  $\mathbb{E}[\mathbf{g}^k]$  and  $\mathbb{E}[(\mathbf{g}^k)^2]$ .

Initialize  $\mathbb{E}[\mathbf{g}^0] = \mathbf{0}$  and  $\mathbb{E}[(\mathbf{g}^0)^2] = \mathbf{0}$ .

At step  $k = 1, 2, \dots$

$$\mathbb{E}[\mathbf{g}^k] = \rho_1 \mathbb{E}[\mathbf{g}^{k-1}] + (1 - \rho_1) \mathbf{g}^k \quad (\text{momentum term})$$

$$\mathbb{E}[(\mathbf{g}^k)^2] = \rho_2 \mathbb{E}[(\mathbf{g}^{k-1})^2] + (1 - \rho_2) (\mathbf{g}^k)^2$$

$$\widehat{\mathbb{E}}[\mathbf{g}^k] = \frac{1}{1 - \rho_1^k} \mathbb{E}[\mathbf{g}^k] \quad (\text{bias correct})$$

$$\widehat{\mathbb{E}}[(\mathbf{g}^k)^2] = \frac{1}{1 - \rho_2^k} \mathbb{E}[(\mathbf{g}^k)^2]$$

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \frac{\alpha}{\sqrt{\mathbb{E}[(\mathbf{g}^k)^2] + \epsilon} \mathbf{1}_d} \widehat{\mathbb{E}}[\mathbf{g}^k]$$

where  $\rho_1$  (0.9) and  $\rho_2$  (0.999) are decay rates while  $\alpha$  is learning rate.

Adam is widely used in training DNN. More other methods such as Nestrov Adam (**Nadam**).

## Mini-Batch Stochastic Gradient Descent

Just replace the stochastic gradient  $\nabla f_{i_k}(\mathbf{x}^k)$  with  $\frac{1}{m} \sum_{i_k=i_1}^{i_m} \nabla f_{i_k}(\mathbf{x}^k)$  where  $\{i_1, \dots, i_m\} \in \{1, \dots, N\}$  is the randomly selected mini-batch index.

Typical mini-batch size  $m = 100$ .

## IV. Training Skills

### Batch Normalization

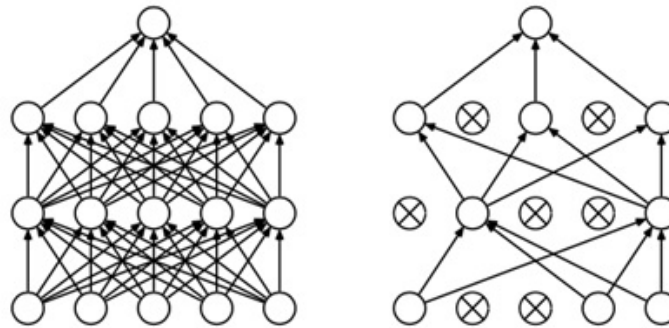
At each layer, shift min-batch data to zero-mean and unit variance.

$$\hat{x}_i = \frac{x_i - \mathbb{E}[x_i]}{\sqrt{\text{var}[x_i]}}$$

$$z_i = \gamma_i \hat{x}_i + \beta_i$$

New technique: Group normalization [\[He'2018\]](#).

### Dropout

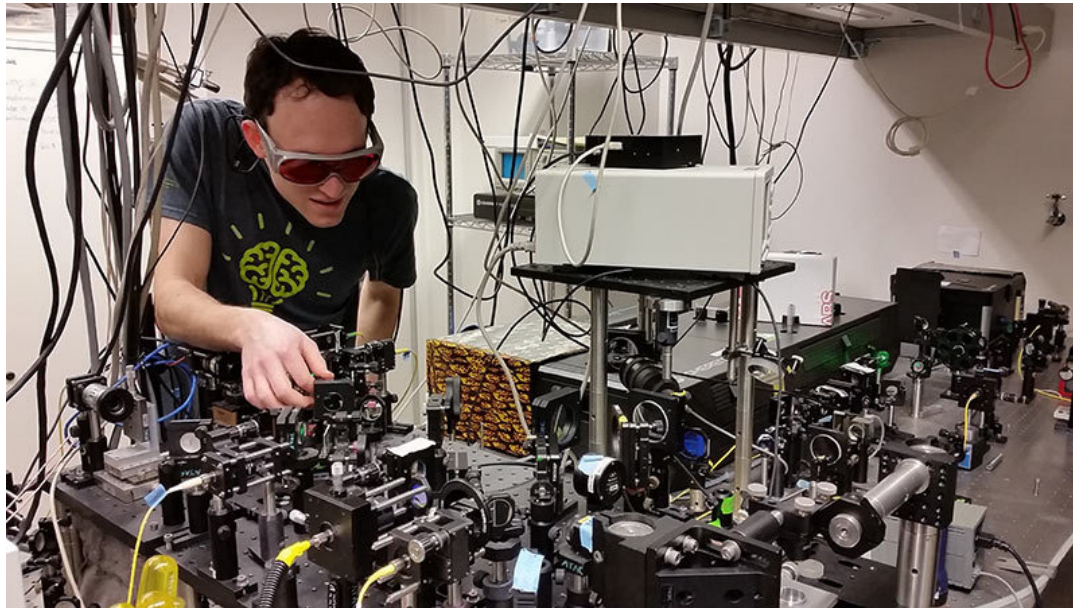




## Parameter Tuning (Alchemy?)

Parameter tuning is very important in deep learning.

Parameter tuning is somewhat like:



This picture has no relation with the context of deep learning.

**Hyper-parameters** need to tune in DNN:

- ☐ Number of layers (depth);
- ☐ Number of neurons of each layer (width)
- ☐ What is the structure of each layer? Linear layer or convolution layer?
- ☐ What is the activation function?
- ☐ Size of mini-batch
- ☐ Which optimization algorithm?
- ☐ Learning rate
- ☐ Any bias is needed for convolution layers?...

Thank you for your attention!