

# Deep Learning

Vazgen Mikayelyan

November 7, 2020



1 Dropout

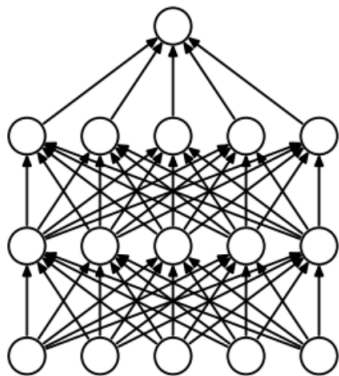
2 Moving Average

3 Batch Normalization

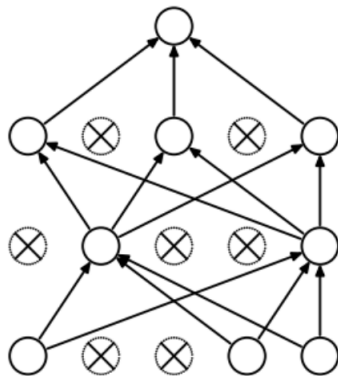
4 Other Optimizers

# Dropout

# Dropout

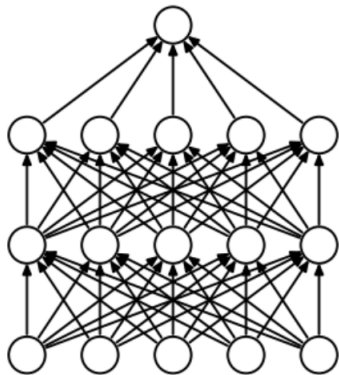


(a) Standard Neural Net

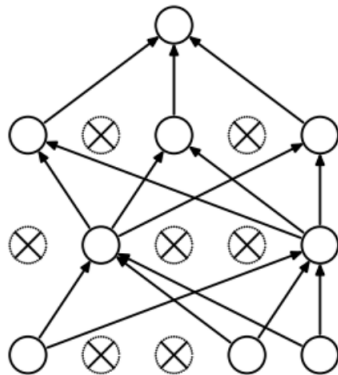


(b) After applying dropout.

# Dropout



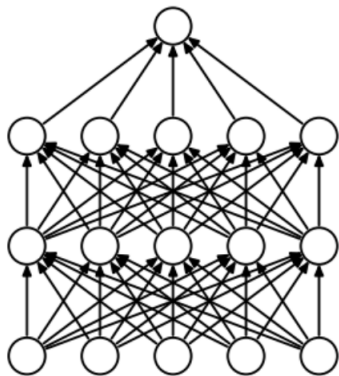
(a) Standard Neural Net



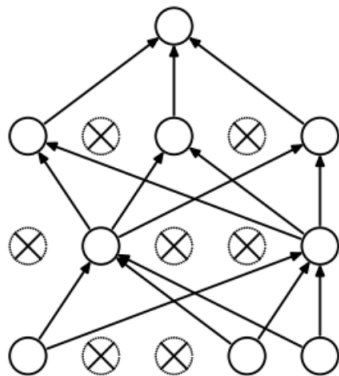
(b) After applying dropout.

What to do during the inference?

# Dropout



(a) Standard Neural Net



(b) After applying dropout.

What to do during the inference?

Answer: Scale units by  $\frac{1}{1 - rate}$  during the training and set  $rate=1$  during the inference.

# Outline

- 1 Dropout
- 2 Moving Average**
- 3 Batch Normalization
- 4 Other Optimizers

## Definition 1

*Simple moving average of the given data is the arithmetic mean of the previous  $k$  data.*



## Definition 1

*Simple moving average of the given data is the arithmetic mean of the previous  $k$  data.*

If you have the data  $x_1, x_2, \dots$ , then its simple moving average will be the following

$$\mu_n = \frac{x_{n-k+1} + \dots + x_n}{k}, n = k, k + 1, \dots$$

## Definition 2

*Cumulative moving average of the given data is the arithmetic mean of the all previous data up to the current time.*

# Cumulative Moving Average

## Definition 2

*Cumulative moving average of the given data is the arithmetic mean of the all previous data up to the current time.*

If you have the data  $x_1, x_2, \dots$ , then its cumulative moving average will be the following

$$\mu_n = \frac{x_1 + x_2 + \dots + x_n}{n}, n = 1, 2, \dots$$

# Cumulative Moving Average

## Definition 2

*Cumulative moving average of the given data is the arithmetic mean of the all previous data up to the current time.*

If you have the data  $x_1, x_2, \dots$ , then its cumulative moving average will be the following

$$\mu_n = \frac{x_1 + x_2 + \dots + x_n}{n}, n = 1, 2, \dots$$

Note that

$$\begin{aligned}\mu_n &= \frac{(x_1 + x_2 + \dots + x_{n-1}) + x_n}{n} = \frac{(n-1)\mu_{n-1} + x_n}{n} \\ &= \left(1 - \frac{1}{n}\right) \mu_{n-1} + \frac{1}{n} x_n.\end{aligned}$$

# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\mu_1 = x_1,$$

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2$$

# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\mu_1 = x_1,$$

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2$$

Note that

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n = \alpha(\alpha\mu_{n-2} + (1 - \alpha)x_{n-1}) + (1 - \alpha)x_n$$

# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\mu_1 = x_1,$$

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2$$

Note that

$$\begin{aligned}\mu_n &= \alpha\mu_{n-1} + (1 - \alpha)x_n = \alpha(\alpha\mu_{n-2} + (1 - \alpha)x_{n-1}) + (1 - \alpha)x_n \\ &= \alpha^2\mu_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n\end{aligned}$$

# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\mu_1 = x_1,$$

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2$$

Note that

$$\begin{aligned}\mu_n &= \alpha\mu_{n-1} + (1 - \alpha)x_n = \alpha(\alpha\mu_{n-2} + (1 - \alpha)x_{n-1}) + (1 - \alpha)x_n \\ &= \alpha^2\mu_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n \\ &= \alpha^3\mu_{n-3} + (1 - \alpha)\alpha^2 x_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n\end{aligned}$$



# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\mu_1 = x_1,$$

$$\mu_n = \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2$$

Note that

$$\begin{aligned}\mu_n &= \alpha\mu_{n-1} + (1 - \alpha)x_n = \alpha(\alpha\mu_{n-2} + (1 - \alpha)x_{n-1}) + (1 - \alpha)x_n \\ &= \alpha^2\mu_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n \\ &= \alpha^3\mu_{n-3} + (1 - \alpha)\alpha^2 x_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n \\ &= \alpha^{n-1}\mu_1 + (1 - \alpha)\alpha^{n-2}x_2 + \dots + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n.\end{aligned}$$

# Exponential Moving Average

If you have the data  $x_1, x_2, \dots$ , then its exponential moving average will be the following

$$\begin{aligned}\mu_1 &= x_1, \\ \mu_n &= \alpha\mu_{n-1} + (1 - \alpha)x_n, \quad n \geq 2\end{aligned}$$

Note that

$$\begin{aligned}\mu_n &= \alpha\mu_{n-1} + (1 - \alpha)x_n = \alpha(\alpha\mu_{n-2} + (1 - \alpha)x_{n-1}) + (1 - \alpha)x_n \\ &= \alpha^2\mu_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n \\ &= \alpha^3\mu_{n-3} + (1 - \alpha)\alpha^2 x_{n-2} + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n \\ &= \alpha^{n-1}\mu_1 + (1 - \alpha)\alpha^{n-2}x_2 + \dots + (1 - \alpha)\alpha x_{n-1} + (1 - \alpha)x_n.\end{aligned}$$

It is easy to see that the sum of the coefficients is equal to 1.

- 1 Dropout
- 2 Moving Average
- 3 Batch Normalization**
- 4 Other Optimizers

- Problem:
  - The distribution of each layer's input changes during training.

# Batch Normalization

- Problem:
  - The distribution of each layer's input changes during training.
- Solution:
  - Fix the distribution of inputs into subnetwork.

- Problem:
  - The distribution of each layer's input changes during training.
- Solution:
  - Fix the distribution of inputs into subnetwork.
- Effects:
  - Improve accuracy.
  - Faster learning.
  - Availability of high learning rates.

**Input:** Values of  $x$  over a mini-batch:  $\mathcal{B} = \{x_1 \dots x_m\}$ ;  
Parameters to be learned:  $\gamma, \beta$

**Output:**  $\{y_i = \text{BN}_{\gamma, \beta}(x_i)\}$

$$\mu_{\mathcal{B}} \leftarrow \frac{1}{m} \sum_{i=1}^m x_i \quad // \text{ mini-batch mean}$$

$$\sigma_{\mathcal{B}}^2 \leftarrow \frac{1}{m} \sum_{i=1}^m (x_i - \mu_{\mathcal{B}})^2 \quad // \text{ mini-batch variance}$$

$$\hat{x}_i \leftarrow \frac{x_i - \mu_{\mathcal{B}}}{\sqrt{\sigma_{\mathcal{B}}^2 + \epsilon}} \quad // \text{ normalize}$$

$$y_i \leftarrow \gamma \hat{x}_i + \beta \equiv \text{BN}_{\gamma, \beta}(x_i) \quad // \text{ scale and shift}$$

- 1 Can we do stochastic gradient descent in this case?



- 1 Can we do stochastic gradient descent in this case?
- 2 What to do during the test?

- 1 Can we do stochastic gradient descent in this case?
- 2 What to do during the test?
- 3 What about biases?

- 1 Dropout
- 2 Moving Average
- 3 Batch Normalization
- 4 Other Optimizers**

# Gradient Descent with Momentum

Let  $L(w)$  be a loss function that we want to minimize. The algorithm gradient descent with momentum is the following

$$v_0 = 0,$$

$$v_t = \beta v_{t-1} + (1 - \beta) \nabla L(w_t),$$

$$w_{t+1} = w_t - \alpha v_t,$$

where  $\alpha$  is the learning rate and  $\beta \in [0, 1)$  is the parameter of exponential moving average.

# Gradient Descent with Momentum



Image 2: SGD without momentum



Image 3: SGD with momentum

Let  $L(w)$  be a loss function that we want to minimize. The algorithm RMSProp is the following

$$\begin{aligned}v_0 &= 0, \\v_t &= \beta v_{t-1} + (1 - \beta) (\nabla L(w_t))^2, \\w_{t+1} &= w_t - \alpha \frac{\nabla L(w_t)}{\sqrt{v_t} + \epsilon},\end{aligned}$$

where  $\alpha$  is the learning rate and  $\beta \in [0, 1)$  is the parameter of exponential moving average.