

Deep Learning

Loo Tiang Kuan, Leonard

A0098795B

supervised by

Associate Professor Ji Hui

co-supervised by

Assistant Professor Alvina Goh Siew Wee

May 5, 2017

Today's Agenda

- 1 Initialization
- 2 Other Gradient Descents
- 3 Other Regulariser

Session 1

- Activations
- Cost
- Stochastic Gradient Descent
- Backpropagation
- Overfitting
- Regularization

Session 3

- Popular ConvNet Architectures
 - AlexNet
 - ZFNet
 - VGG
 - GoogLeNet
 - ResNet
 - Denoising

Session 2

- Initialisation
- Hyper-parameters
- Variants of SGD
- More Regularisers
- Intro to ConvNets

Session 4/5

- Autoencoders
- Transposed Convolution
- Generative Network
- Generative Adversarial Network

Recap: Stochastic Gradient Descent

Randomly sample into m mini-batches of size k :

$$\begin{aligned}\mathbf{x} &= \left\{ \{x^{(1)}, x^{(2)}, \dots, x^{(k)}\}, \{x^{(k+1)}, x^{(k+2)}, \dots, x^{(2k)}\}, \dots, \{x^{(n-k+1)}, \dots, x^{(n)}\} \right\} \\ &= \{M_1, M_2, \dots, M_m\}\end{aligned}$$

$$\left. \frac{\partial C}{\partial w} \right|_{\mathbf{x}} \approx \left. \frac{\partial C}{\partial w} \right|_{\mathbf{x}_{M_j}}, \quad j = 1, \dots, m$$

In one Epoch, we perform the following with $j = 1, \dots, m$ times:

$$\begin{aligned}w_{k+1} &\leftarrow w_k - \eta \left. \frac{\partial C}{\partial w_k} \right|_{\mathbf{x}_{M_j}} \\ b_{k+1} &\leftarrow b_k - \eta \left. \frac{\partial C}{\partial b_k} \right|_{\mathbf{x}_{M_j}}\end{aligned} \tag{1}$$

- 1 Initialization
 - Weights Initialisation
 - Hyperparameters
- 2 Other Gradient Descents
- 3 Other Regulariser

Weights Initialisation

Initially, we use

$$w \sim \mathcal{N}(0, 1)$$

But with large number of w , we have for $z = w^T x + b$

$$z \sim \mathcal{N}(0, \sqrt{N})$$

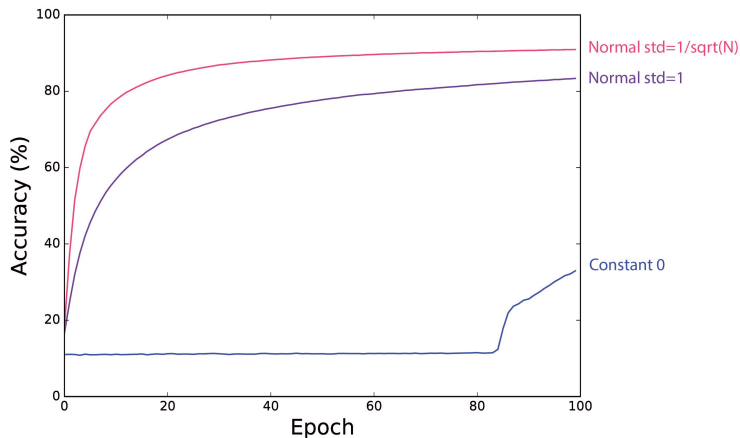
where z have higher chance of taking on large values in the first iteration.
Therefore we let

$$w \sim \mathcal{N}\left(0, \frac{1}{\sqrt{N}}\right)$$

so that

$$z \sim \mathcal{N}(0, 1)$$

Weights Initialisation



Glorot, X., & Bengio, Y. (2010, May). *"Understanding the difficulty of training deep feedforward neural networks."* In Aistats (Vol. 9, pp. 249-256).

1 Initialization

- Weights Initialisation
- Hyperparameters
 - Architecture
 - Learning Rate & Regularizer
 - Mini-batch Size

2 Other Gradient Descents

3 Other Regulariser

1 Initialization

- Weights Initialisation
- Hyperparameters
 - Architecture
 - Learning Rate & Regularizer
 - Mini-batch Size

2 Other Gradient Descents

3 Other Regulariser

Setting Up architecture

- Reduce problem to speed up learning
 - Classify from $\{0, 1\}$ instead of classifying $\{0, \dots, 9\}$
- Start with low number of hidden layers
- Monitor validation accuracy (start with 20% of training data)
- Overfit.



Eldan, R. and Shamir, O., 2015. "*The Power of Depth for Feedforward Neural Networks.*" arXiv preprint [arXiv:1512.03965](https://arxiv.org/abs/1512.03965).

Hyperparameters

$w_{i,j}^l$: Weights

$b_{i,j}^l$: Biases

m : Mini-Batch Size

η : Learning Rate

λ : Regularizer

1 Initialization

- Weights Initialisation
- Hyperparameters
 - Architecture
 - Learning Rate & Regularizer
 - Mini-batch Size

2 Other Gradient Descents

3 Other Regulariser

Learning Rate & Regularizer

- Start with setting $\lambda = 0$
- Increase/Decrease η by a factor of 10 to find out a good order of magnitude
- Fine-tune both λ & η once a good order has been established

1 Initialization

- Weights Initialisation
- Hyperparameters
 - Architecture
 - Learning Rate & Regularizer
 - Mini-batch Size

2 Other Gradient Descents

3 Other Regulariser

Mini-batch Size

- If too small, GPU not sufficiently used
- If too large, weights not updated often enough

1 Initialization

2 Other Gradient Descents

- SGD
- Momentum
- Nesterov Momentum
- AdaGrad
- RMSProp
- Adam

3 Other Regulariser

Stochastic Gradient Descent

Randomly sample into m mini-batches of size k :

$$\begin{aligned}\mathbf{x} &= \left\{ \{x^{(1)}, x^{(2)}, \dots, x^{(k)}\}, \{x^{(k+1)}, x^{(k+2)}, \dots, x^{(2k)}\}, \dots, \{x^{(n-k+1)}, \dots, x^{(n)}\} \right\} \\ &= \{M_1, M_2, \dots, M_m\}\end{aligned}$$

$$\left. \frac{\partial C}{\partial w} \right|_{\mathbf{x}} \approx \left. \frac{\partial C}{\partial w} \right|_{\mathbf{x}_{M_j}}, \quad j = 1, \dots, m$$

In one Epoch, we perform the following with $j = 1, \dots, m$ times:

$$\begin{aligned}w_{k+1} &\leftarrow w_k - \eta \left. \frac{\partial C}{\partial w_k} \right|_{\mathbf{x}_{M_j}} \\ b_{k+1} &\leftarrow b_k - \eta \left. \frac{\partial C}{\partial b_k} \right|_{\mathbf{x}_{M_j}}\end{aligned} \tag{2}$$

1 Initialization

2 Other Gradient Descents

- SGD
- **Momentum**
- Nesterov Momentum
- AdaGrad
- RMSProp
- Adam

3 Other Regulariser

Momentum

Idea: Introduce friction/velocity component.

$$\begin{aligned}v_{k+1} &= \mu v_k - \eta \frac{\partial C}{\partial w} \\w_{k+1} &= w_k + v_{k+1} \\&= (w_k + \mu v_k) - \eta \frac{\partial C}{\partial w}\end{aligned}$$

where $0 < \mu < 1$ correspond to friction/velocity magnitude.

1 Initialization

2 Other Gradient Descents

- SGD
- Momentum
- **Nesterov Momentum**
- AdaGrad
- RMSProp
- Adam

3 Other Regulariser

Nesterov Momentum

Idea: Use gradient ONLY after travelling along μv_k ,

$$v_{k+1} = \mu v_k - \eta \left. \frac{\partial C}{\partial w} \right|_{w=w_k + \mu v_k}$$

$$\begin{aligned} w_{k+1} &= w_k + v_{k+1} \\ &= (w_k + \mu v_k) - \eta \left. \frac{\partial C}{\partial w} \right|_{w=w_k + \mu v_k} \end{aligned}$$

Let $\widehat{w}_k = w_k + \mu v_k$,

$$\widehat{w}_{k+1} = \widehat{w}_k - \mu v_k + (1 + \mu) v_{k+1}$$

1 Initialization

2 Other Gradient Descents

- SGD
- Momentum
- Nesterov Momentum
- **AdaGrad**
- RMSProp
- Adam

3 Other Regulariser

Idea: Store all past gradients,

$$c_{k+1} = c_k + \left(\frac{\partial C}{\partial w_k} \right)^2 = \sum_{i=0}^k \left(\frac{\partial C}{\partial w_k} \right)^2 \quad (\text{sum of squares: second moment})$$

$$w_{k+1} = w_k - \frac{\eta}{\sqrt{c_{k+1}}} \frac{\partial C}{\partial w_k}, \quad \text{Learning rate decay}$$
$$= w_k - \eta \frac{\partial C}{\partial w_k} \cdot \frac{1}{\sqrt{c_{k+1}}}, \quad \text{Entry-wise division}$$

- (+) Helps in learning rate decay.
- (+) Control amplitude of descent (not direction!).
- (-) Effective gradient converges to 0 too early.

1 Initialization

2 Other Gradient Descents

- SGD
- Momentum
- Nesterov Momentum
- AdaGrad
- **RMSPprop**
- Adam

3 Other Regulariser

Idea: Introduce decaying rate α in storing of past gradients,

$$\begin{aligned}c_{k+1} &= \alpha c_k + (1 - \alpha) \left(\frac{\partial C}{\partial w_k} \right)^2 \\&= (1 - \alpha) \sum_{i=0}^k \alpha^{k-i} \left(\frac{\partial C}{\partial w_k} \right)^2 \\w_{k+1} &= w_k - \eta \frac{\partial C}{\partial w_k} \cdot \frac{1}{\sqrt{c_{k+1}}}\end{aligned}$$

- Does not lead to zero gradient update in the long run.
- Default value: $\alpha = 0.9$.

1 Initialization

2 Other Gradient Descents

- SGD
- Momentum
- Nesterov Momentum
- AdaGrad
- RMSProp
- Adam

3 Other Regulariser

Idea: Incorporate first moments v_{k+1} and second moments c_{k+1} ,

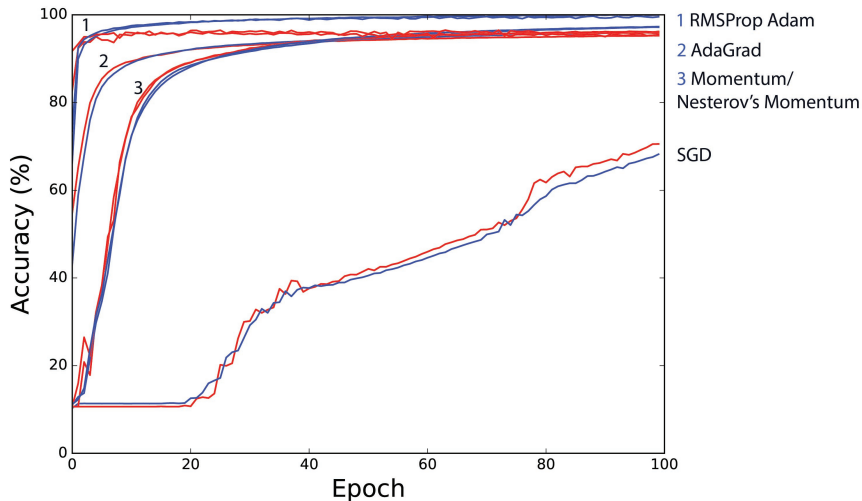
$$v_{k+1} = \beta_1 v_k + (1 - \beta_1) \frac{\partial C}{\partial w_k} = (1 - \beta_1) \sum_{i=0}^k \beta_1^{k-i} \frac{\partial C}{\partial w_k},$$

$$c_{k+1} = \beta_2 c_k + (1 - \beta_2) \left(\frac{\partial C}{\partial w_k} \right)^2 = (1 - \beta_2) \sum_{i=0}^k \beta_2^{k-i} \left(\frac{\partial C}{\partial w_k} \right)^2,$$

$$\begin{aligned} w_{k+1} &= w_k - \eta \frac{v_{k+1}}{\sqrt{c_{k+1}}} \\ &= w_k - \eta \frac{\beta_1 v_k + (1 - \beta_1) \frac{\partial C}{\partial w_k}}{\sqrt{\beta_2 c_k + (1 - \beta_2) \left(\frac{\partial C}{\partial w_k} \right)^2}}. \end{aligned}$$

Typically, $\beta_1 \approx 0.9$, $\beta_2 \approx 0.999$.

Comparison



References



Sutskever, I., Martens, J., Dahl, G. E., & Hinton, G. E. (2013). *“On the importance of initialization and momentum in deep learning.”* ICML (3), 28, 1139-1147.



Duchi, J., Hazan, E., & Singer, Y. (2011). *“Adaptive subgradient methods for online learning and stochastic optimization.”* Journal of Machine Learning Research, 12(Jul), 2121-2159.



T. Tieleman & G. E. Hinton. Lecture 6.5-rmsprop: *“Divide the gradient by running average of its recent magnitude.”* , 2012



Kingma, D., & Ba, J. (2014). *“Adam: A method for stochastic optimization.”* arXiv preprint arXiv:1412.6980.

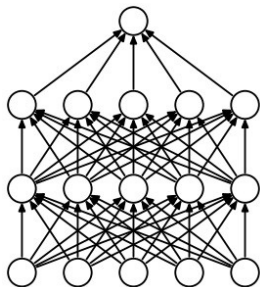


“CS231n Winter 2016: Lecture 6: Neural Networks Part 3/Intro to ConvNets”

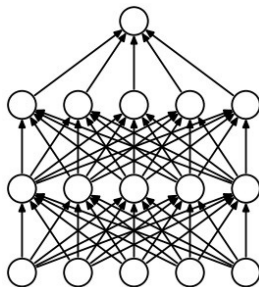
Outline

- 1 Initialization
- 2 Other Gradient Descents
- 3 Other Regulariser
 - Dropout
 - Batch Normalization

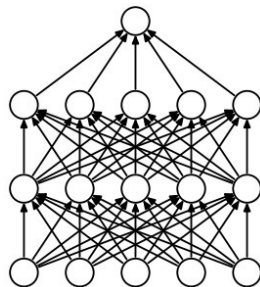
Ensemble



Network 1



Network 2



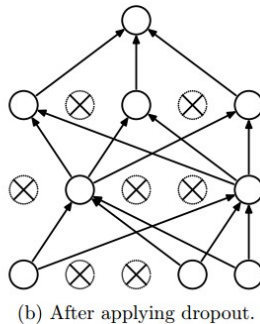
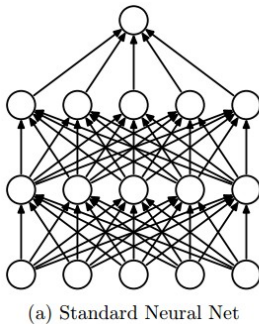
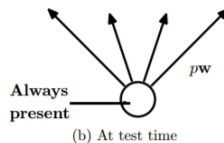
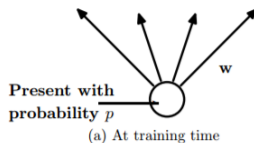
Network 3

Train each network with different weights and hyperparameter initialisations.

Compute probability of each class for each data.

Taking their average improves accuracy by about 2% (empirically).

Dropout



Method	Testing Error
mlpconv + Fully Connected	11.59%
mlpconv + Fully Connected + Dropout	10.88%
mlpconv + Global Average Pooling	10.41%

References



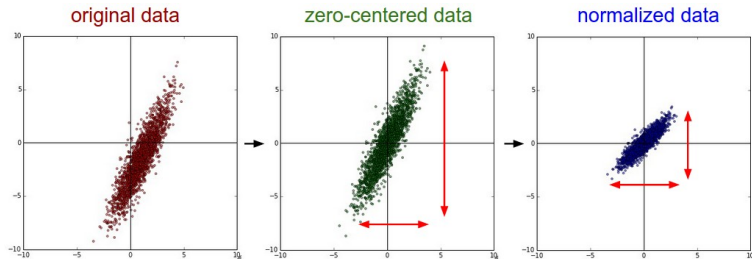
Srivastava, Nitish, et al. *"Dropout: a simple way to prevent neural networks from overfitting."* Journal of Machine Learning Research 15.1 (2014): 1929-1958.

- 1 Initialization
- 2 Other Gradient Descents
- 3 Other Regulariser
 - Dropout
 - Batch Normalization
 - Problem
 - Batch Normalization

Outline

- 1 Initialization
- 2 Other Gradient Descents
- 3 Other Regulariser
 - Dropout
 - Batch Normalization
 - Problem
 - Batch Normalization

Problem



- Internal Covariate Shift

- Consider layer a^l ,

$$a^l = g(W^l \cdot a^{l-1} + b^l)$$

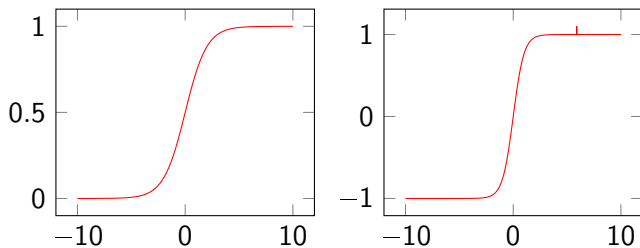
$$W_{k+1}^l = W_k^l - \eta \frac{\partial C}{\partial W_k^l}$$

- Distribution of the a^{l-1} changes every epoch, difficult to learn
- When there is fixed distribution in a^{l-1} , W^l does not have to readjust.

Outline

- 1 Initialization
- 2 Other Gradient Descents
- 3 Other Regulariser
 - Dropout
 - Batch Normalization
 - Problem
 - Batch Normalization

Batch Normalization



- Batch Normalisation

- Whitening z^l before committing $a^l = g(z^l)$.
- $a^l = g(\hat{z}^l)$ where $\hat{z}^l = \text{norm}(z^l)$.
- $y^l = \gamma^l \hat{z}^l + \beta^l$
to prevent $\sigma(\hat{z}^l)$ or $\tanh(\hat{z}^l)$ being stuck at the linear region
- Extra terms to learn in each batch norm layer γ^l, β^l .

Backpropagation:

$$\frac{\partial C}{\partial \hat{z}_i^l} = \frac{\partial C}{\partial y_i^l} \gamma^l$$

$$\frac{\partial C}{\partial \sigma_B^2} = \sum_{i=1}^m \frac{\partial C}{\partial \hat{z}_i^l} \cdot (z_i^l - \mu_B) \cdot \frac{-1}{2} (\sigma_B^2)^{-\frac{3}{2}}$$

$$\frac{\partial C}{\partial \mu_B} = \left(\sum_{i=1}^m \frac{\partial C}{\partial \hat{z}_i^l} \cdot \frac{-1}{\sqrt{\sigma_B^2}} \right) + \frac{\partial C}{\partial \sigma_B^2} \cdot \frac{\sum_{i=1}^m 2(\mu_B - z_i^l)}{m}$$

$$\Rightarrow \frac{\partial C}{\partial y_i^l} = \left[\frac{\partial C}{\partial \hat{z}_i^l} \cdot \frac{-1}{\sqrt{\sigma_B^2}} \right] + \left[\frac{\partial C}{\partial \sigma_B^2} \cdot \frac{2(z_i^l - \mu_B)}{m} \right] + \left[\frac{\partial C}{\partial \mu_B} \cdot \frac{1}{m} \right]$$

Backpropagation:

$$\frac{\partial C}{\partial \gamma^l} = \sum_{i=1}^m \frac{\partial C}{\partial y_i^l} \hat{z}_i^l$$

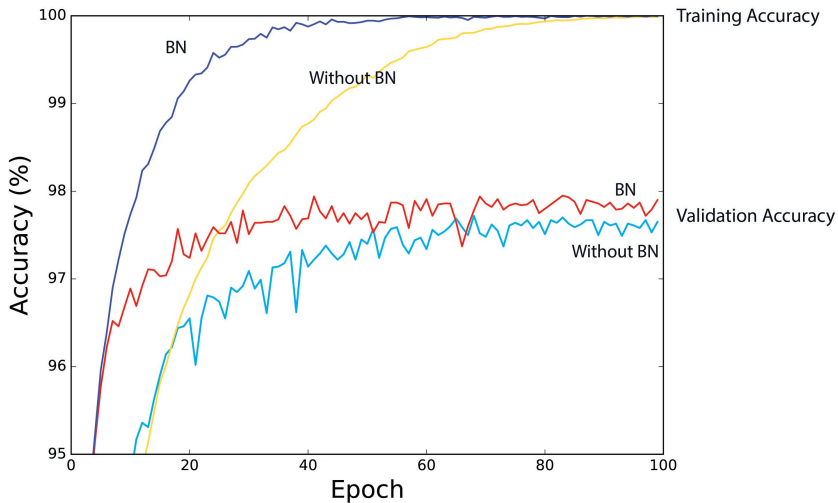
$$\frac{\partial C}{\partial \beta^l} = \sum_{i=1}^m \frac{\partial C}{\partial y_i^l}$$

$$\frac{\partial C}{\partial w^l} = \frac{\partial C}{\partial z_i^l} a^{l-1}$$



Ioffe, S., & Szegedy, C. (2015). *"Batch normalization: Accelerating deep network training by reducing internal covariate shift."* arXiv preprint arXiv:1502.03167.

Batch Normalization



- Benefits

- Larger learning rate η can be used.
- Accelerate η decay due to greater training speed.
- Acts as a regulariser.
 - Can reduce Dropout strength.
 - Can reduce λ in L1 or L2 weights penalty.
- Reduces saturation.

Thank You