

Scribe Notes for Lecture 8

Lecturer: Navin Goyal*Scribe:* Kapil Pathak

1 Decomposition of test error into approximation, optimization, generalization errors

Consider a set up with supervised with feed-forward networks where \mathcal{X} is input distribution, \mathcal{Y} target and \mathcal{H} hypothesis class. Let \mathcal{A} be an algorithm for learning the network.

$h_{A,S}$ = Hypothesis output by \mathcal{A} on input \mathcal{S} .

$$h_S^* = \operatorname{argmin}_{h \in \mathcal{H}} L_S(h)$$

$$h^* = \operatorname{argmin}_{h \in \mathcal{H}} L_D(h)$$

For good generalization, a neural network needs to approximate classification function (assuming a classification task). For that, the training error should be as small as possible. But that doesn't imply the neural network will have a good generalization over the unseen data.

According to above notation, $L_D(h_{A,S})$ is a generalization error due to a hypothesis obtained from algorithm \mathcal{A} and training set \mathcal{S} . $L_D(h_{A,S})$ can be decomposed as,

$$L_D(h_{A,S}) = L_D(h_{A,S}) - L_S(h_{A,S}) + L_S(h_{A,S}) - L_S(h_S^*) + L_S(h_S^*) - L_D(h^*) + L_D(h^*)$$

where,

$$\text{Approximation error}(\epsilon_{app}) = L_D(h^*)$$

$$\text{Optimization error}(\epsilon_{opt}) = L_S(h_{A,S}) - L_S(h_S^*)$$

$$\text{Generalization error}(\epsilon_{gen}) = \sup_{h \in \mathcal{H}} |L_S(h) - L_D(h)|$$

$$L_S(h_S^*) - L_D(h^*) \leq L_S(h^*) - L_D(h^*) \leq \text{Generalization error}$$

From above equations, we get

$$L_D(h_{A,S}) \leq \epsilon_{app} + \epsilon_{opt} + 2\epsilon_{gen}$$

Denote the Training Set $\mathcal{S} = \left\{ (x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), \dots, (x^{(m)}, y^{(m)}) \right\}$

where $x^{(i)} \in \mathbb{R}^d$, denoted by $x^{(i)} = (x_1^{(i)}, x_2^{(i)}, \dots, x_d^{(i)})$, and $y^{(i)} \in \{0, 1\}$

2 Approximation Error (ϵ_{app})

2.1 Universal Approximation Theorem [4][5][3]

Let $P : \mathbb{R} \rightarrow \mathbb{R}$ be non constant, bounded and continuous. Then for any $\epsilon > 0$, positive integer m , $f : [0, 1]^m \rightarrow \mathbb{R}$, \exists integer N , real numbers a_i, b_i and vectors $w_i \in \mathcal{R}^m$ for $i \in \{1, 2, \dots, N\}$, such that, $F(x) = \sum_{i \in [N]} a_i P(\langle w_i, x \rangle + b_i)$ is an ϵ_{approx} of $f : |f(x) - F(x)| \leq \epsilon \forall x \in [0, 1]^m$.

But the question arises whether the depth increases expressive power of a neural network.

3 Generalization Error (ϵ_{gen})

Diagram pending: Size of one hidden layer network on small MNIST.

Here, the point to note is that no overfitting as we overparametrize the neural network. This contradicts the bias-variance trade-off.

CIFAR-10 dataset is an image dataset having training set size= 5×10^4 Images of size $32 \times 32 \times 3 = 3072$

Architecture	No of parameters
AlexNet	1.3×10^6
Inception	1.6×10^6
Imagenet	1.2×10^6
Inception V4	4.2×10^7
ResNet-152	6×10^7
VGG-19	1.4×10^8
AmaebaNet -B (Top 1 accuracy=84.3 % and Top 5 accuracy=97 %)	5.5×10^8

3.1 A Convergence Theory for Deep Learning via Overparametrization [2]

3.2 Learning and Generalization in Overparametrized Neural Networks [1]

3.2.1 Informal Statement

- Training data = $\{(x_i, y_i^*)\}_{i \in [n]}$ where $x \in R^\partial, y_i^* \in R^d$
- Separability Assumption: \forall distinct $i, j \in [n]$ we have $\|x_i - x_j\| \geq \delta$
- L-hidden layer ReLU fully connected network with each other of size of m .
- He initialization
- l_2 - loss for regression, cross-entropy for classification

- For regression, $m \geq \text{poly}(n, L, \frac{1}{\delta})$ then gradient descent or stochastic gradient descent, find ϵ - error global minimum for l_2 - loss in $\text{poly}(n, L, \frac{1}{\delta} \log(\frac{1}{\epsilon}))$ iterations.
- For classification, gradient descent or stochastic gradient descent find a classifier with 0 classification error in $\text{poly}(n, L, \frac{1}{\delta})$ iterations with number of samples $m \geq \text{poly}(n, L, \frac{1}{\delta})$
- Extensions to other Lip-smooth losses and architectures.

3.2.2 Remarks:

- It makes ReLU crucial.
- Independent of dimension d
- Overparametrization in theorem implies in practice
- Dependency on L (Not clear if it's correct)

Main insight in [1] and [2]

$$\text{ReLU}(\langle w, x \rangle + b) = \langle w, x \rangle + b \text{ if } \langle w, x \rangle + b \geq 0$$

For highly overparametrized settings, activation pattern essentially remain constant, weights don't change much.

For $w^0, w^1, \dots, w^t, \|w^0 - w^t\|$ is small. Still small change is sufficient as the network is over-parametrized.

Notation: Consider l_2 -regression

$$\phi(x) = \max\{0, x\} = \text{ReLU}(x) \text{ For } v \in R^m, \phi(v) = (\phi(v_1), \phi(v_2), \dots, \phi(v_m)).$$

Data is normalized such that $\|x_i\| = 1$ and $(x_i)_\partial = \frac{1}{\sqrt{2}}$

We assume that there no biases added before applying activation functions.

Output from a neural network can be given as,

$$y_i = N(x_i) = B\phi(W_L\phi(\dots W_2\phi(W_1\phi(Ax)))) \text{ where } B \text{ denotes output layer,}$$

W_L denotes Hidden layer L, $A \in R^{m \times d}, W_l \in R^{m \times m}, B \in R^{d \times m}$

Initialize $A_{ij} \sim N(0, \frac{2}{m})$ for $i \in [m], j \in [\partial]$

$[W_l^{(0)}] \sim N(0, \frac{2}{m})$ for $i, j \in [m], l \in [L]$

$B_{ij} \sim N(0, \frac{1}{d})$ for $i \in [d], j \in [m]$

$$\vec{W} = (W_1, \dots, W_L), \vec{W}^{(0)} = (W_1^{(0)}, \dots, W_L^{(0)})$$

$$F(\vec{W}) = \sum_{i \in [n]} F_i(\vec{W})$$

$$F_i(\vec{W}) = \frac{1}{2} \|N(x_i) - y_i^*\|_2^2 \text{ for } i \in [n]$$

$$\nabla F(\vec{W}) = (\nabla_{W_1} F(\vec{W}), \dots, \nabla_{W_L} F(\vec{W}))$$

Here, we are training with respect to \vec{W} keeping A and B to initialize values.

References

- [1] Zeyuan Allen-Zhu, Yuanzhi Li, and Yingyu Liang. Learning and generalization in over-parameterized neural networks, going beyond two layers. *CoRR*, abs/1811.04918, 2018.
- [2] Zeyuan Allen-Zhu, Yuanzhi Li, and Zhao Song. A convergence theory for deep learning via over-parameterization. *CoRR*, abs/1811.03962, 2018.
- [3] A. R. Barron. Universal approximation bounds for superpositions of a sigmoidal function. *IEEE Trans. Inf. Theor.*, 39(3):930–945, May 1993.
- [4] G. Cybenko. Approximation by superpositions of a sigmoidal function. *Mathematics of Control, Signals, and Systems (MCSS)*, 2(4):303–314, December 1989.
- [5] Kurt Hornik. Approximation capabilities of multilayer feedforward networks. *Neural Networks*, 4(2):251–257, 1991.