

Deep Learning

Shuhei Watanabe

April 4, 2021

1 Introduction

Applications and research of deep learning (DL) are currently active thanks to **abundant libraries**, which allow non experts to implement algorithms without deep knowledge, and achieve **quick and strong performance** and **wide range of application** such as to image, video, graph data and so on. Additionally, DL can find good and abstract representation of given data ¹. Basically, DL learns such representations through the learning process. In the conventional ML, the domain experts extract such features or better representation of data. On the other hand, DL allows users to omit this process and enables the one pipeline of ML process, i.e. end-to-end learning. In the history of DL, there are two big waves. However, each of them declined for the following reasons:

1. Linear models cannot represent the XOR function and backlash against biologically inspired learning
2. Researchers could not realize big networks because of the lack of big data and adequate computational resources

At the first wave, logistic regression was invented, but this model ² cannot represent the XOR function. At the second wave, these kinds of problems were resolved, but there were not sufficient computational power. However, the large computational server is now accessible for anyone and that led to the current third wave.

2 KL divergence vs Maximum likelihood estimation

The goal of DL is to approximate the distribution of the output \mathbf{y} given a data point \mathbf{x} :

$$\min_{\mathbf{w}} \text{distance}(\hat{p}(\mathbf{y}|\mathbf{x}; \mathbf{w}), p(\mathbf{y}|\mathbf{x})) \quad (1)$$

It can be achieved by the minimization of **KL divergence** between the modeled distribution and the target distribution as follows:

$$D_{\text{KL}}(p||\hat{p}) = \int p(\mathbf{y}|\mathbf{x}) \log \frac{p(\mathbf{y}|\mathbf{x})}{\hat{p}(\mathbf{y}|\mathbf{x}; \mathbf{w})} d\mathbf{y} \quad (2)$$

Note that this metric is **not symmetric**, i.e. $D_{\text{KL}}(p||\hat{p}) \neq D_{\text{KL}}(\hat{p}||p)$ and the divergence is equal to zero if and only if $p = \hat{p}$ for all \mathbf{y} . In fact, the minimization of this metric is equivalent to the minimization of the cross entropy. To validate this statement, let's look at both formal and empirical form of the cross entropy:

$$\begin{aligned} H(p, \hat{p}) &= - \int p(\mathbf{y}|\mathbf{x}) \log \hat{p}(\mathbf{y}|\mathbf{x}; \mathbf{w}) d\mathbf{y} \\ &\simeq - \frac{1}{n} \sum_{i=1}^n \log \hat{p}(\mathbf{y}_i|\mathbf{x}_i; \mathbf{w}) \end{aligned} \quad (3)$$

¹For example, the roman numeral system is not a good representation for computation, but DL can extract better representation for this.

²Logistic regression can represent NAND, OR and AND. XOR is (NAND) AND (OR).

where $(\mathbf{x}_i, \mathbf{y}_i)$ for all i ($1 \leq i \leq n$) is the observation. Then, the conventional information entropy is the following:

$$\begin{aligned} H(p) &= - \int p(\mathbf{y}|\mathbf{x}) \log p(\mathbf{y}|\mathbf{x}) d\mathbf{y} \\ &\simeq -\frac{1}{n} \sum_{i=1}^n \log p(\mathbf{y}_i|\mathbf{x}_i) \end{aligned} \quad (4)$$

By combining both entropy, we will obtain the following:

$$D_{\text{KL}}(p||\hat{p}) = H(p, \hat{p}) - H(p) \quad (5)$$

The goal of DL is to identify the weight vector \mathbf{w} , which minimizes the cross entropy $H(p, \hat{p})$, and this equation can be reformulated as follows:

$$\mathbf{w}^* \in \operatorname{argmin}_{\mathbf{w}} H(p(\mathbf{y}|\mathbf{x}; \mathcal{D}), \hat{p}(\mathbf{y}|\mathbf{x}; \mathbf{w})) \quad (6)$$

where \mathcal{D} is a given dataset. Since the posterior distribution $p(\mathbf{y}|\mathbf{x}; \mathcal{D})$ given a dataset \mathcal{D} is independent from \mathbf{w} , the following holds:

$$\operatorname{argmin}_{\mathbf{w}} H(p, \hat{p}) = \operatorname{argmin}_{\mathbf{w}} (H(p, \hat{p}) - H(p)) = \operatorname{argmin}_{\mathbf{w}} D_{\text{KL}}(p||\hat{p}) \quad (7)$$

where $p = p(\mathbf{y}|\mathbf{x}; \mathcal{D})$, $\hat{p} = \hat{p}(\mathbf{y}|\mathbf{x}; \mathbf{w})$. By transforming the equation, we will obtain the following result:

$$\begin{aligned} \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} H(p, \hat{p}) = \operatorname{argmin}_{\mathbf{w}} -p \log \hat{p} = \operatorname{argmax}_{\mathbf{w}} p \log \hat{p} \\ &= \operatorname{argmax}_{\mathbf{w}} \sum_{i=1}^n \log \hat{p}(\mathbf{y}_i|\mathbf{x}_i; \mathbf{w}) = \operatorname{argmax}_{\mathbf{w}} \prod_{i=1}^n \hat{p}(\mathbf{y}_i|\mathbf{x}_i; \mathbf{w}) \end{aligned} \quad (8)$$

From this result, the maximum likelihood estimation (MLE) is proven to be equivalent to the minimization of KL divergence. KL divergence gives us zero if distribution is reproduced regardless of if the distribution has noise or not; therefore, KL divergence is the interest of observation in some contexts. Note that both formulations give negative infinity loss for wrong answers with high confidence and strictly punish them.

3 Neural networks

3.1 Non-linear regression vs neural networks

Neural networks learn the non-linear feature space through the learning process and **users do not have to define such non-linear mappings** so much. While each layer computes the weighted sum of the values from the previous layer, each node applies non-linear mapping and each layer stacks on such outputs. In the end, neural networks obtain the non-linear feature space suitable for a given task. On the other hand, non-linear regression takes the feature space by selecting the basis functions and these functions are basically not parameterized. While neural networks automatically learn the feature space, users have to design a suitable feature space by hand in non-linear regression.

3.2 Forward pass

Each layer of neural networks computes the weighted sum of the values from the previous layer and the activation function of the weighted sum. Let B be the batch size of input data, $m^{(l)}$ be the number of units in the l -th hidden layer and $a^{(l)}(\cdot)$ be the activation function at the $(l+1)$ -th layer. Additionally, let the output from the l -th hidden layer be $\mathbf{Z}^{(l)} \in \mathbb{R}^{B \times D}$, the weight vector from the l -th to $(l+1)$ -th layer ³ be $\mathbf{W}^{(l)} \in \mathbb{R}^{m^{(l)} \times m^{(l+1)}}$, the bias vector be $\mathbf{b}^{(l)} \in \mathbb{R}^{m^{(l+1)}}$. Then, the forward pass of neural

³the 0-th layer is the input layer and the $(L+1)$ -th layer is the output layer.

Table 1: The main usage of the activation function

Task	Activation functions for output layer
Regression	Linear
Binary classification	Logistic or tanh
Multi-class classification	softmax

networks are $\mathbf{Z}^{(l+1)} = \mathbf{a}^{(l)}(\mathbf{Z}^{(l)})\mathbf{W}^{(l)} + \mathbf{B}^{(l)} = \mathbf{A}^{(l)}\mathbf{W}^{(l)} + \mathbf{B}^{(l)}$ where $\mathbf{B}^{(l)}$ is the broadcasted $\mathbf{b}^{(l)}$ defined as follows:

$$\mathbf{B}^{(l)} = \begin{bmatrix} (\mathbf{b}^{(l)})^\top \\ (\mathbf{b}^{(l)})^\top \\ \vdots \\ (\mathbf{b}^{(l)})^\top \end{bmatrix} \quad (9)$$

The computational complexity and the memory complexity of the forward pass in the l -th layer is $O(Bm^{(l)}m^{(l+1)})$ and $O(m^{(l)}m^{(l+1)})$, respectively ⁴. If $m^{(l)}$ is identical for all l , then the total computational complexity and the total memory complexity of the forward pass is $O(Bm^2L)$ and $O(m^2L)$, respectively.

3.3 The activation functions

- **sigmoid** $f(x) = 1/(1 + e^{-x})$: Bounded in $[0, 1]$ and useful for the probabilistical representation
- **softmax** $f(y = i|\mathbf{x}) = e^{x_i} / \sum_{j=1}^K e^{x_j}$: Bounded in $[0, 1]$ and used for multi-class classification
- **tanh** $f(x) = (e^x + e^{-x}) / (e^x - e^{-x})$: Bounded in $[-1, 1]$ and symmetric and does not change the center of distribution, useful for **hidden layers**
- **linear** $f(x) = x$: Useful for regression output layer
- **ReLU** $f(x) = \max(0, x)$: Widely used the most, non-differentiable, but practically no problem. The advantages are the constant derivative values and the sparsity, but the problem is killed neurons are not updated at all.

Note that ReLU is widely used because it can handle the **gradient vanishing problem** ⁵. Other than those, Parametric ReLU (PReLU), Exponential Linear Unit (ELU), Gaussian Error Linear Unit, SWISH are introduced. Logistic or softmax function are used for the output layer in classification tasks because of **the differentiability**. We can use hard-max for the forward pass, but we need the differentiability in the backward pass. Table 1 shows the main use of activation functions. When dealing with regression tasks, the linear function is preferred because it can represent the whole real numbers. On the other hand, for example, ReLU cannot represent the negative number, so it is not suitable for such regression tasks.

3.4 Representational capacity of deep learning

If Multi-layer perceptron (MLP) has more than two layers, the following theorem (**Universal Approximation Theorem**) holds:

⁴In the backward pass, the memory complexity is $O(m^{(l)}m^{(l+1)} + Bm^{(l)} + Bm^{(l+1)})$

⁵For example, the gradients of sigmoid and tanh are smaller than or equal to 0.25 and 1, respectively. As the architecture becomes deeper, it will be suffered from the gradient vanishing and slower convergence.

Theorem 1

1. Any Boolean function
2. Any bounded continuous function

can be realized or approximated with arbitrary precision by a MLP with one hidden layer.

However, we have to keep in mind the following notes:

1. The arbitrary function can be approximated if the **infinite** number of units are available
2. This theorem does not say that such MLP can be obtained, but it says there exists such MLPs

Therefore, we need more advanced models. In the recent research, it is showed that DL has more flexible and sophisticated representational capacity. The following theorem supports this proposition:

Theorem 2

A neural network with n_0 inputs and L layers of n units each, with ReLU activations can represent functions that have $\Omega((n/n_0)^{n_0(L-1)}n^{n_0})$ linear regions

Roughly speaking, current ML models with ReLU draw decision boundary or estimated functions curve using the joint of linear lines. If the length of each linear region goes to zero, such a curve or a decision boundary becomes smoother and more flexible. Therefore, the increase in the order of linear regions leads to more flexible and powerful models. Since the exponential part of Theorem 2 only has the number of layers, but not the number of units, we can tell DL has more representational capacity than shallower or wider models.

3.5 Backward pass

3.5.1 Naïve computation

The following chain rule is required to compute backward pass.

Theorem 3

Given $x \in \mathbb{R}^m, y \in \mathbb{R}^n, f : \mathbb{R}^m \rightarrow \mathbb{R}^n, g : \mathbb{R}^n \rightarrow \mathbb{R}$ and let z be $g(y)$ and y be $f(x)$, then the following equation holds:

$$\frac{\partial z}{\partial x_j} = \sum_{i=1}^n \frac{\partial z}{\partial y_i} \frac{\partial y_i}{\partial x_j} \quad (10)$$

Let $X \in \mathbb{R}^{B \times D}$ be the input data, $Y \in \mathbb{R}^{B \times C}$ be the output label data (each row has onehot vector) and the batch size is B as in the forward pass section. Then, the forward pass is $\mathbf{Z}^{(l+1)} = \mathbf{A}^{(l)}\mathbf{W}^{(l)} + \mathbf{B}^{(l)}$.

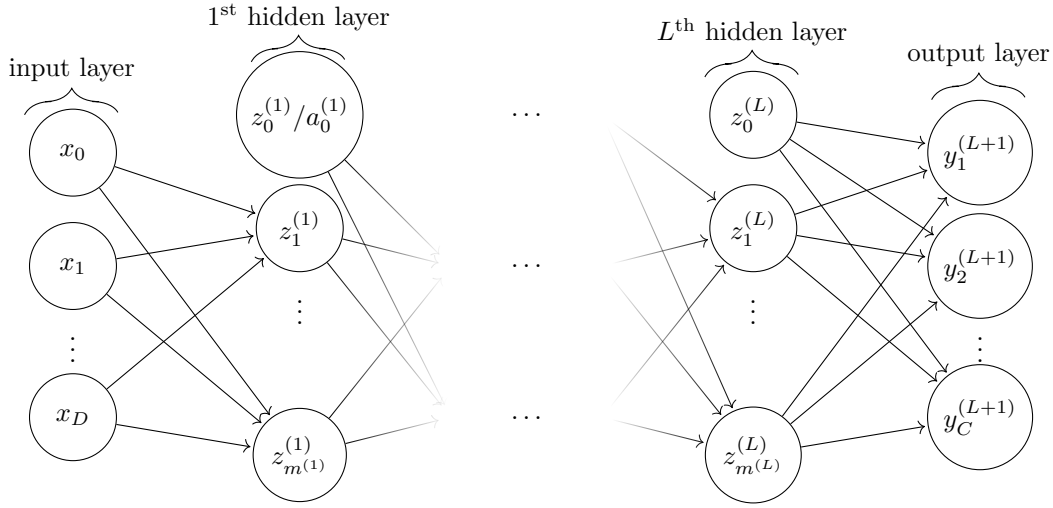


Figure 1: Network graph of a $(L + 1)$ -layer perceptron with D input units and C output units. The l -th hidden layer contains $m^{(l)}$ hidden units. $z_k^{(l)}$ is the weighted sum for the k -th node at the l -th hidden layer. The activation value for the corresponding node is $a_k^{(l)}$.

In the same vein, we introduce the backward pass as follows:

$$\begin{aligned}
 \frac{\partial L}{\partial \mathbf{A}^{(l)}} &= \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{A}^{(l)}} \\
 \frac{\partial L}{\partial \mathbf{W}^{(l)}} &= \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{W}^{(l)}} \\
 \frac{\partial L}{\partial \mathbf{b}^{(l)}} &= \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \frac{\partial \mathbf{Z}^{(l+1)}}{\partial \mathbf{b}^{(l)}} \\
 \frac{\partial L}{\partial \mathbf{Z}^{(l)}} &= \frac{\partial L}{\partial \mathbf{A}^{(l)}} \frac{\partial \mathbf{A}^{(l)}}{\partial \mathbf{Z}^{(l)}}
 \end{aligned} \tag{11}$$

where $L = L(\mathbf{y}, \hat{\mathbf{y}})$ is the loss function. In this formulation, we have to store the 4 dimensional tensor and this is often infeasible.

3.6 Matrix computation of back propagation

Since $\mathbf{Z}_{i,j}^{(l+1)}$ is the value of the i -th batch at the of j -th unit, the gradients with respect to i' -th batch $\mathbf{A}_{i',:}^{(l)}$ and the weight to j' -th unit $\mathbf{W}_{:,j'}^{(l)}$ is zero ($i' \neq i, j' \neq j$). Therefore, the computation of each element in $\mathbf{Z}^{(l+1)}$ is as follows:

$$\mathbf{Z}^{(l+1)} = \begin{bmatrix} z_{1,1}^{(l+1)} & \cdots & z_{1,m^{(l+1)}}^{(l+1)} \\ \vdots & \ddots & \vdots \\ z_{B,1}^{(l+1)} & \cdots & z_{B,m^{(l+1)}}^{(l+1)} \end{bmatrix} = \begin{bmatrix} b_1^{(l)} + \sum_{i=1}^{m^{(l)}} w_{i,1}^{(l)} a_{1,i}^{(l)} & \cdots & b_{m^{(l+1)}}^{(l)} + \sum_{i=1}^{m^{(l)}} w_{i,m^{(l+1)}}^{(l)} a_{1,i}^{(l)} \\ \vdots & \ddots & \vdots \\ b_1^{(l)} + \sum_{i=1}^{m^{(l)}} w_{i,1}^{(l)} a_{B,i}^{(l)} & \cdots & b_{m^{(l+1)}}^{(l)} + \sum_{i=1}^{m^{(l)}} w_{i,m^{(l+1)}}^{(l)} a_{B,i}^{(l)} \end{bmatrix} \tag{12}$$

In the same vein, we do not have to compute most gradients and each element in matrices are the followings:

$$\begin{aligned}
\left(\frac{\partial L}{\partial \mathbf{A}^{(l)}} \right)_{i,j} &= \sum_{k=1}^{m^{(l+1)}} \frac{\partial L}{\partial z_{i,k}^{(l+1)}} \frac{\partial z_{i,k}^{(l+1)}}{\partial a_{i,j}^{(l)}} = \sum_{k=1}^{m^{(l+1)}} \frac{\partial J}{\partial z_{i,k}^{(l+1)}} w_{j,k}^{(l)} \\
\left(\frac{\partial L}{\partial \mathbf{W}^{(l)}} \right)_{i,j} &= \sum_{k=1}^B \frac{\partial L}{\partial z_{k,j}^{(l+1)}} \frac{\partial z_{k,j}^{(l+1)}}{\partial w_{i,j}^{(l)}} = \sum_{k=1}^B \frac{\partial L}{\partial z_{k,j}^{(l+1)}} a_{k,j}^{(l)} \\
\left(\frac{\partial L}{\partial \mathbf{b}^{(l)}} \right)_i &= \sum_{k=1}^B \frac{\partial L}{\partial z_{k,i}^{(l+1)}} \frac{\partial z_{k,i}^{(l+1)}}{\partial b_i^{(l)}} = \sum_{k=1}^B \frac{\partial L}{\partial z_{k,i}^{(l+1)}}
\end{aligned} \tag{13}$$

Since $a_{i,j}^{(l)}$ has connections to only the i -th data point at the following layer and $w_{i,j}^{(l)}$ has connections to only the j -th unit at the following layer, the chain rule has to be applied to only the aforementioned elements. Therefore, the computation can be written in the following matrix forms:

$$\begin{aligned}
\frac{\partial L}{\partial \mathbf{A}^{(l)}} &= \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} (\mathbf{W}^{(l)})^\top \\
\frac{\partial L}{\partial \mathbf{W}^{(l)}} &= (\mathbf{A}^{(l)})^\top \frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \\
\frac{\partial L}{\partial \mathbf{b}^{(l)}} &= \left(\frac{\partial L}{\partial \mathbf{Z}^{(l+1)}} \right)^\top \mathbf{1}^B
\end{aligned} \tag{14}$$

where $\mathbf{1}^B = [1, \dots, 1]^\top \in \mathbb{R}^B$. Finally, $\partial L / \partial \mathbf{Z}^{(l)}$ is obtained using the chain rule as follows:

$$\frac{\partial L}{\partial \mathbf{Z}^{(l)}} = \frac{\partial L}{\partial \mathbf{A}^{(l)}} \odot \mathbf{G}^{(l)} \tag{15}$$

where $G_{i,j}^{(l)} = \partial a_{i,j}^{(l)} / \partial z_{i,j}^{(l)} = \partial a^{(l)}(z_{i,j}^{(l)}) / \partial z_{i,j}^{(l)}$ and \odot is the Hadamard product⁶. The initial gradients are the gradients of the output $\hat{Y} = \mathbf{A}^{(L+1)}$ with respect to $\mathbf{Z}^{(L+1)}$ as follows:

$$\frac{\partial L}{\partial \mathbf{Z}^{(L+1)}} = \frac{\partial L}{\partial \mathbf{A}^{(L+1)}} \frac{\partial \mathbf{A}^{(L+1)}}{\partial \mathbf{Z}^{(L+1)}} = \frac{\partial L}{\partial \mathbf{A}^{(L+1)}} \odot \mathbf{G}^{(L+1)} \tag{16}$$

In the case where $g^{(L+1)}$ is the softmax, the gradients are the followings:

$$\frac{\partial L}{\partial \mathbf{Z}^{(L+1)}} = \frac{\hat{Y} - Y}{N} \tag{17}$$

By introducing the initial values to Eq (14), we will obtain the whole gradient values. The overall memory complexity of back propagation is $O(LM^2 + LMB)$ in the case of $M = m^{(i)}$ for all $i = 1, \dots, L$.

4 Optimization for deep learning

DL requires the optimization of the weight vector. The difficulty of optimization tasks depends on the properties of tasks. We start the discussion of convex optimization and turn to non-convex optimization. Finally, we introduce momentum, the adaptive learning rate and the initialization strategy of the weight vector.

⁶Element-wise multiplication

4.1 Learning vs pure optimization

Learning has the following four differences compared to pure optimization:

1. may have to use surrogate instead of the real objective such as softmax instead of hard-max
2. optimization on partial data distribution instead of the real distribution
3. often introduces regularization term to generalize
4. requires the summation over given data points, so computationally expensive for a large dataset

Note that surrogate may also enable better generalization.

4.2 Convex optimization

The definition of convex is as follows:

Definition 1

Given C^2 function $f(\mathbf{x}) : \mathbb{R}^D \rightarrow \mathbb{R}$ is convex iff the Hessian $H \in \mathbb{R}^{D \times D}$ where $H_{i,j} = \frac{\partial^2 f}{\partial x_i \partial x_j}$ is positive semidefinite.

Note that Hessian matrix is symmetric if and only if the given function has $\frac{\partial^2 f(\mathbf{x})}{\partial x_i \partial x_j}$ for all pairs (i, j) , $(i \neq j)$ and all of them are continuous. The definition of positive semidefinite is the following:

Definition 2

A matrix A is positive semidefinite if one of the followings is satisfied:

1. $\forall \mathbf{x} \in \mathbb{R}^D, \mathbf{x}^\top A \mathbf{x} \geq 0$
2. The all eigenvalues are non-negative
3. $\exists L \in \mathbb{R}^{D \times D}$ s.t. $A = LL^\top$

The convex optimization has a convergence guarantee. However, there is an ill-condition, which makes the convergence slower. The ill-condition is defined based on the condition number described as:

Definition 3

The condition number of the Hessian matrix H of a given convex function is $\kappa(H) = \frac{\lambda_{\max}}{\lambda_{\min}}$ where $\lambda_{\max}, \lambda_{\min}$ are the maximum and minimum eigenvalue, respectively.

When $\kappa(H)$ is large, the given function is called ill-conditioned and it causes the zig-zag behavior. There are two solutions for this.

1. **momentum**: By introducing the decayed step size from the previous iterations, it suppresses the zig-zag behavior
2. **preconditioning**: transforms the function into good condition and optimizes in the space. Then after optimizing in the space, it re-maps into the original space.

We will discuss the momentum ⁷ later and now discuss the preconditioning. To explain the method, let's take the quadratic example $f(\mathbf{x}) = \frac{1}{2} \mathbf{x}^\top H \mathbf{x}$ where H is positive definite. In the preconditioning,

⁷In gradient descent of the quadratic function, the convergence rate is $|\mathbf{x}_{t+1} - \mathbf{x}^*| \leq \frac{\kappa(H)-1}{\kappa(H)+1} |\mathbf{x}_t - \mathbf{x}^*|$. However, if it is with momentum, the convergence rate is $|\mathbf{x}_{t+1} - \mathbf{x}^*| \leq \frac{\sqrt{\kappa(H)}-1}{\sqrt{\kappa(H)}+1} |\mathbf{x}_t - \mathbf{x}^*|$.

we substitute this function with $g(\mathbf{y}) = f(P^{\frac{1}{2}}\mathbf{y})$ where P is an arbitrary positive definite symmetric matrix and $P^{\frac{1}{2}}$ is the square-root of P ⁸. In the preconditioning, we optimize $g(\mathbf{y})$ with respect to \mathbf{y} . Since the Hessian matrix of $g(\mathbf{y})$ is $P^{\frac{1}{2}}HP^{\frac{1}{2}}$, we can change the condition number of g . Ideally, $P^{\frac{1}{2}}HP^{\frac{1}{2}}$ should be an identity matrix. For this reason, we would like to take P to satisfy this condition and such P can be obtained as follows:

$$\begin{aligned} I &= P^{\frac{1}{2}}HP^{\frac{1}{2}} = P^{\frac{1}{2}}H^{\frac{1}{2}}H^{\frac{1}{2}}P^{\frac{1}{2}} = H^{-\frac{1}{2}}H^{\frac{1}{2}}H^{\frac{1}{2}}H^{-\frac{1}{2}} \\ (P^{\frac{1}{2}})^2 &= (H^{-\frac{1}{2}})^2 \\ P &= H^{-1} \end{aligned} \tag{18}$$

where we utilize the fact that H is positive definite and symmetric. In the case of non-convex optimization, there might not be H^{-1} and it does not mean the descent direction obtained from the Hessian matrix improves the function value even if we can compute H^{-1} .

4.3 Optimization of neural networks

There are several distinctive properties in the optimization of neural networks.

1. **Non-convex:** Hessian cannot tell the best direction and no guarantee of global solutions.
2. **Large datasets:** It requires the summation over data points, so the complexity goes up
3. **High dimensionality:** The computation of Hessian is computationally expensive
4. **Structure of neural networks:** The response surfaces are likely to be harder

The loss function of neural networks are typically non-convex function and there can be **many saddle points**⁹, which are the attractors for second-order methods. Therefore, we cannot guarantee the convergence in the optimization. Furthermore, the following properties in the structure of neural networks make optimizations harder.

1. **Clif-like surfaces in the searching space:** The gradients suddenly drop or rise significantly at some points especially in RNN.
2. **Gradient vanishing and exploding:** Since we have to multiply matrices again and again, each value is going to be very small or large by the power factor.

4.4 Stochastic gradient descent (SGD)

While the gradient descent handles **learning by epoch**, SGD computes the loss value using a part of the dataset. Note that learning by epoch is to learn all the training data at the same time. Typically, the optimization by gradient descent is stable, but it takes a lot more time. Therefore, we use so-called **mini-batch**, i.e. **a part of the dataset**, for one update step. SGD stochastically picks a batch of data points and learns the batch at each step. Note that both methods see all the data points in 1 iteration and this iteration is called **epoch**. Mini-batch learning is preferred due to the following reasons:

1. Much quicker training allows much more iteration
2. Since the standard error of empirical loss function against true loss function decreases by the order of $O(\frac{1}{\sqrt{B}})$ where B is the batch size, the larger batch size does not make a big difference.

⁸In the case of symmetric positive definite matrix, but not in the case of positive semidefinite matrix, Cholesky decomposition also works.

⁹When the loss goes down, it is more likely to have local minima and less likely to have saddle points.

3. The training dataset can include many similar data points, so mini-batch learning is likely to yield the similar result to that of the full-batch training.
4. Since mini-batch includes less data points and the standard error is larger, the neural networks are **not likely to overfit** the statistics of given data distribution.

The standard error is, roughly speaking, the expectable error of the distribution of a bootstrapped dataset from the true distribution and the empirical loss approaches the loss using full data as B goes to the number of data points. Note that when we use smaller batch size, we have to decrease the learning rate as well. This is because the smaller batch size has the same effect as the larger learning rate.

4.5 Learning rate and momentum

4.5.1 Learning rate

It is known that the convergence (including local minima) of SGD is guaranteed when we set the learning rate α_t scheduling as follows:

$$\begin{aligned} \sum_{t=1}^{\infty} \alpha_t &= \infty \\ \sum_{t=1}^{\infty} \alpha_t^2 &< \infty \end{aligned} \tag{19}$$

In this sense, we should schedule the learning rate using some methods and the followings achieves such scheduling:

1. **Linear decay:** $\alpha_t = (1 - \frac{t}{\tau})\alpha_0 + \frac{t}{\tau}\alpha_{\tau}$ (τ is the final epoch)
2. **Exponential decay:** $\alpha_t = c^t \alpha_0$
3. **Step decay:** decay the learning rate by a certain interval
4. **Cosine decay:** $\alpha_t = \frac{\alpha_0}{2}(1 + \cos \frac{t\pi}{\tau})$

Note that we conventionally adapt the learning rate by checking the validation and training learning curve of NN.

4.5.2 Momentum

The momentum has been introduced to suppress the zigzag behavior¹⁰ and it can be formulated as follows:

$$\begin{aligned} \mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha \hat{\mathbf{g}}_t + \beta \underbrace{(\mathbf{w}_t - \mathbf{w}_{t-1})}_{\text{Let it be } \Delta \mathbf{w}_{t-1}} \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha \hat{\mathbf{g}}_t + \beta(-\alpha \hat{\mathbf{g}}_{t-1} + \beta \Delta \mathbf{w}_{t-2}) \\ &\dots \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \alpha \sum_{i=0}^t \beta^{k-i} \hat{\mathbf{g}}_i \end{aligned} \tag{20}$$

By transforming the equation, we obtain the following:

$$\begin{aligned} \mathbf{G} &= \beta \mathbf{G} + \alpha \hat{\mathbf{g}} \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \mathbf{G} \end{aligned} \tag{21}$$

¹⁰Since the momentum takes over the dynamics from the previous iteration, it cancels out the force to go to the opposite direction.

The advantages of the momentum are **the smoothing effect of zig-zagging, the acceleration of learning at flat spots and the deacceleration of the learning when the sign of derivatives change**. On the other hand, the disadvantages of momentum are the additional parameter β and it potentially triggers additional zig-zagging when we use a strong momentum.

4.6 Gradient descent with adaptive learning rate

As mentioned previously, the learning rate has to be adapted for the learning convergence and different parameters (or dimensions) typically require a different learning rate. Therefore, **AdaGrad** has been invented. The update of AdaGrad is as follows:

$$\begin{aligned} \mathbf{G} &= \mathbf{G} + \hat{\mathbf{g}} \odot \hat{\mathbf{g}} \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{G}} + \epsilon} \odot \hat{\mathbf{g}} \end{aligned} \quad (22)$$

where $\epsilon \ll 1$ is a positive number. When the accumulated gradients \mathbf{G} is large, it implies the corresponding weight values have changed much compared to those have smaller \mathbf{G} . Therefore, AdaGrad slows down the learning by making the learning rate smaller. The problem of AdaGrad is the long term influences of the initial gradients and it potentially leads to premature learning.

The revised version of AdaGrad is **RMSProp**. The update of this equation is the following:

$$\begin{aligned} \mathbf{G} &= \rho \mathbf{G} + (1 - \rho) \hat{\mathbf{g}} \odot \hat{\mathbf{g}} \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{\alpha}{\sqrt{\mathbf{G}} + \epsilon} \odot \hat{\mathbf{g}} \end{aligned} \quad (23)$$

By introducing the decay rate ρ , RMSProp can reduce the influences from older epochs.

The last one is **Adam**.

$$\begin{aligned} \mathbf{G}_1 &= \rho_1 \mathbf{G}_1 + (1 - \rho_1) \hat{\mathbf{g}} \\ \mathbf{G}_2 &= \rho_2 \mathbf{G}_2 + (1 - \rho_2) \hat{\mathbf{g}} \odot \hat{\mathbf{g}} \\ \hat{\mathbf{G}}_1 &= \frac{\mathbf{G}_1}{1 - \rho_1^t} \\ \hat{\mathbf{G}}_2 &= \frac{\mathbf{G}_2}{1 - \rho_2^t} \\ \mathbf{w}_{t+1} &= \mathbf{w}_t - \frac{\alpha}{\sqrt{\hat{\mathbf{G}}_2} + \epsilon} \odot \hat{\mathbf{G}}_1 \end{aligned} \quad (24)$$

The correction is based on the equation $\sum_{i=1}^t \rho^{t-i} (1 - \rho) = (1 - \rho) \frac{1 - \rho^t}{1 - \rho} = 1 - \rho^t$. Note that the first moment \mathbf{G}_1 is equivalent to momentum. Although these methods adaptively change learning rate, it does not mean we do not have to schedule the learning rate. In fact, there is a paper stating that the learning rate scheduling can improve the performance of Adam.

4.7 The initialization of weights

Basically, weights are initialized using uniform or normal distribution¹¹ and bias is initialized using a constant number. In practice, the scale of the initialization is important. When the scale is larger, break symmetry¹² happens. On the other hand, if the scale is smaller, the training result would be similar to each other and thus it can be generalized. In the same way, when the scale is larger, we can make sure most activation functions do not output 0 in the forward pass. However, if we take smaller scale, we can avoid the gradient explosion.

¹¹A well-known example is Xavier. $w \sim \mathcal{U}\left(-\sqrt{\frac{6}{d_{\text{in}} + d_{\text{out}}}}, \sqrt{\frac{6}{d_{\text{in}} + d_{\text{out}}}}\right)$

¹²When we set all the weights as an identical number, all the weights will be updated in the same way. Break symmetry is to break this phenomenon.

The initialization has to be implemented in each layer separately and one effective strategy for this would be to check the scale of activation function values in each layer using mini-batch data and adapt the scale to make sure the even initialization for each layer. Since we use the same learning rate for each layer, we cannot expect the same learning for each layer when we do not have **the same scales**; therefore, the learning becomes difficult.

For the initialization of bias, there are some exceptional cases where we do not initialize by zero:

1. **The bias before ReLU unit:** We set a positive number (e.g. 0.1) to avoid the zero output
2. **Output units:** If we have prior information regarding the task, we put this information explicitly

5 Regularization

The main idea behind the regularization is to reduce the variance in the expected test loss ($= \text{bias}^2 + \text{variance} + \text{noise}$) in exchange of some bias. One of the most naïve methods for the regularization is **bagging**. Bagging is to bootstrap samples from a training dataset and ensemble the trained models on each bootstrapped dataset, but this is not used in DL due to the expensive training cost.

5.1 Parameter sharing

If we can share parameters in the network over multiple tasks, we implicitly have more data for training the model. By the parameter sharing, we can achieve the followings:

1. Since dealing with multiple tasks is harder, it can control the network's capacity
2. Since it shares parameters among similar tasks, it encourages the search for regular patterns

There are two major methods to achieve parameter sharing. One is to **share the shallower parts of parameters over multiple tasks** and since it implicitly increases the amount of training data, it can reduce both memory storage and computational time. The other is called **parameter tying** and it makes parameters of given two networks similar using a regularization term between their parameters as follows:

$$\Omega(\mathbf{w}^{(A)}, \mathbf{w}^{(B)}) = \|\mathbf{w}^{(A)} - \mathbf{w}^{(B)}\|_2 \quad (25)$$

where $\mathbf{w}^{(A)}, \mathbf{w}^{(B)}$ are the weight parameters for task A and B .

5.2 Dropout

Dropout is to drop some weights in a given neural network model with a certain probability p during training. The basic concept is to discourage the co-adaptation of weights. In other words, dropout prevents the neural network model from relying on some weights too much and encourages each weight to have some more contributions to solving given tasks. Additionally, since the dropped graphs are the subgraphs of a given model, dropout can implicitly realize an ensemble. The important point is that we have to **store the location of dropped weights in the forward pass to ignore the update in the backward pass**. Note that since dropout drops some weights in the neural network model, we have to scale the output by considering the dropped weights. The right solution for this is to compute by the following:

$$f(\mathbf{x}; \mathbf{w}) = \sum_{\mathbf{w}_s \in \mathbf{w}} p(\mathbf{w}_s) f(\mathbf{x}; \mathbf{w}_s) \quad (26)$$

However, this is not tractable, so it is performed by dividing each output by $\frac{1}{1-p}$ empirically.

5.3 Data augmentation / Noise robustness

Data augmentation is to increase the data amount ¹³ by applying some simple transformation to a training dataset and add this transformed data to the training dataset. Since deep learning typically requires so much data to yield high performance and it is cumbersome to collect more data, data augmentation is the effective solution. The followings are the concrete examples:

1. **translation**: To move data along each axis
2. **scaling**: To change the size scale of data
3. **reflection**: To invert the left and right side of data ¹⁴
4. **rotation**: To tilt data
5. **stretching**: To change the aspect ratio of data

Note that since **CNNs have invariance with respect to the movement along each axis, translation does not make sense**. The drawback of data augmentation is to require more computational time proportional to the data size. Additionally, the improvement diminishes as the augmented amount increases. For this reason, we do not increase more than 100x in the applications.

Another option is to add random noise to various components of neural networks as follows:

1. **Inputs**: data augmentation
2. **Hidden units**: data augmentation of extracted features
3. **Weights**: Bayesian neural networks $p(f(\mathbf{x})|\mathbf{x}, \mathcal{D}) = \int p(f(\mathbf{x})|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w}$
4. **Outputs**: label smoothing

Label smoothing is to introduce some noise to label. For example, suppose we would like to classify given images into 3 classes, the labels for this task looks like one-hot vector $\mathbf{y}_i = [0, 1, 0]$. However, by introducing noise, we can make it as $\mathbf{y}_i = [0.1, 0.8, 0.1]$. This method stops the neural networks from overfitting to given specific classes. The drawback is to require the tuning of the noise control parameter.

5.4 Adversarial training

Adversarial data is the data which deceive the trained model. This method is called fast gradient sign method and formulated as follows:

$$\mathbf{x}' = \mathbf{x} + \epsilon \operatorname{sign}\left(\frac{\partial L(y, \hat{y})}{\partial \mathbf{x}}\right) \quad (27)$$

where ϵ is a sufficiently small positive number. Basically, this new input \mathbf{x}' increases the loss function. By adding adversarial data to training dataset, we can encourage local constancy¹⁵ and mitigates sensitivity to small changes in inputs.

¹³We usually augment only training dataset, but people sometimes augment test or validation dataset as well for the sake of measuring robustness.

¹⁴It does not make sense to augment by reflection in the case of letter recognitions.

¹⁵The function is said to have local constancy if the output does not change much within a small region of space.

5.5 Weight preference

First, it is ideal to stop the training when the validation loss reaches its minimum and starts to go up and **early-stopping** is the strategy stopping the training earlier. The second method is called **shrinkage method** or weight decay. This term is defined as the L^p -norm of the weight vector. $L1$ - and $L2$ -norm are well known and $L1$ -norm promotes the sparsity and $L2$ -norm solves the colinearity. The update of the weight vector can be described as follows:

$$\begin{aligned}
 L &= L_{\mathcal{D}} + \frac{\lambda}{2} \|\mathbf{w}\|^2 \\
 \mathbf{g} &= \frac{\partial L}{\partial \mathbf{w}} = \mathbf{g}_{\mathcal{D}} + \lambda \mathbf{w} \\
 \mathbf{w} &= \mathbf{w} - \alpha(\mathbf{g}_{\mathcal{D}} + \lambda \mathbf{w}) \\
 \mathbf{w} &= \underbrace{(1 - \alpha\lambda)}_{\text{Approaches 1}} \mathbf{w} - \alpha \mathbf{g}_{\mathcal{D}}
 \end{aligned} \tag{28}$$

where α is a learning rate. In the case of the optimization methods with adaptive learning rate such as Adam and RMSProp, since the learning rate α for some weights goes down quickly, the weight decay effect is small for the corresponding weights. In this sense, we have to adjust the regularization term by decoupling as follows:

$$\mathbf{w} = (1 - \beta)\mathbf{w} - \alpha \mathbf{g}_{\mathcal{D}} \tag{29}$$

where β is a positive number. By doing this, even when the learning rate goes down, the regularization works properly. Note that while SGD with the weight decoupling is equivalent to SGD with L2 regularization, Adam with the weight decoupling is not equivalent to Adam with L2 regularization. This is because the learning rate for each weight is different in Adam.

6 Convolutional neural networks (CNN)

6.1 Forward pass

CNNs are one type of neural networks and it performs well on data with grid-like structure such as image and time series. As real-world applications, it is applied to object detection, segmentation¹⁶, image captioning, skeletal pose estimation. The forward pass of CNNs uses the operation called **cross-correlation**. Before introducing cross-correlation, we have to define the **convolution** operation as follows:

Definition 4

Given two matrix $\mathbf{A} \in \mathbb{R}^{H \times W}$, $\mathbf{K} \in \mathbb{R}^{H' \times W'}$, the convolution $\mathbb{R}^{H' \times W'} \times \mathbb{R}^{H \times W} \rightarrow \mathbb{R}^{(H-H'+1) \times (W-W'+1)}$ is the following operation:

$$(\mathbf{K} * \mathbf{A})_{h,w} = \sum_{x=1}^{W'} \sum_{y=1}^{H'} K_{y,x} A_{h-y,w-x} \tag{30}$$

where $W' + 1 \leq w \leq W + 1$, $H' + 1 \leq h \leq H + 1$.

¹⁶It detects objects in a pixel-wise way.

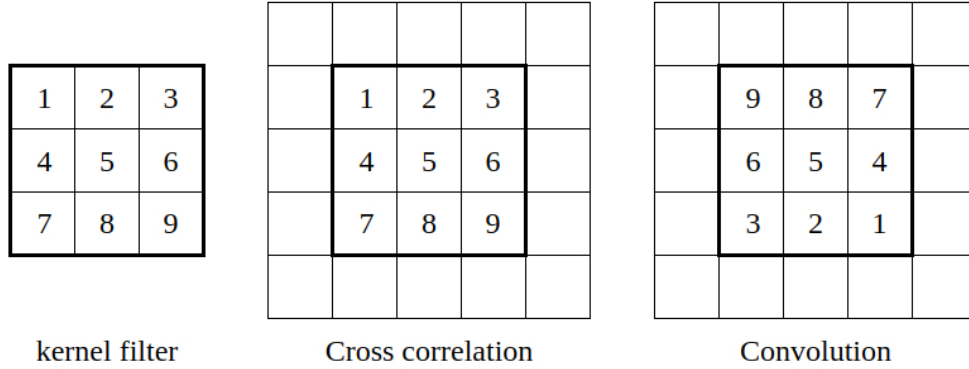


Figure 2: The conceptual visualization of cross correlation.

The cross-correlation is defined using convolution and flipped kernel \mathbf{K} as follows:

$$\begin{aligned}
 (\mathbf{K} \otimes \mathbf{A})_{h,w} &= (\mathbf{K}' * \mathbf{A})_{h+H'+1,w+W'+1} = \sum_{x=1}^{W'} \sum_{y=1}^{H'} K'_{y,x} A_{h+H'-y+1,w+W'-x+1} \\
 &= \sum_{x=1}^{W'} \sum_{y=1}^{H'} K_{H'-y+1,W'-x+1} A_{h+H'-y+1,w+W'-x+1} \quad (31) \\
 &= \sum_{x=1}^{W'} \sum_{y=1}^{H'} K_{y,x} A_{h+y,w+x}
 \end{aligned}$$

where \mathbf{K}' is 180-degree flipped kernel filter, i.e. $K'_{y,x} = K_{H'-y+1,W'-x+1}$ and $0 \leq h \leq H - H', 0 \leq w \leq W - W'$. From this equation, it can be said that **cross-correlation is equivalent to the combination of flipped kernel and convolution**. The example of flipped kernel is illustrated in Figure 2. The forward pass of CNNs are computed as follows:

$$Z_{B,c} = \sum_{j=1}^{C_{\text{in}}} K_{c,j} \otimes A_{B,j} + b_c \mathbf{1}_{H_{\text{out}} \times W_{\text{out}}} \quad (32)$$

where B, c are the index of batch and the output channel, $C_{\text{in}}, C_{\text{out}}$ are the number of channels in the current and the following layers and $H_{\text{out}}, W_{\text{out}}$ are the height and the width of the output image and $\mathbf{K} \in \mathbb{R}^{C_{\text{out}} \times C_{\text{in}} \times H'_{\text{in}} \times W'_{\text{in}}}$ is a kernel tensor and $\mathbf{A} \in \mathbb{R}^{B \times C_{\text{in}} \times H'_{\text{in}} \times W'_{\text{in}}}$ is a input value from the previous layer and $\mathbf{b} \in \mathbb{R}^{C_{\text{out}}}$ is a bias vector and $\mathbf{Z} \in \mathbb{R}^{B \times C_{\text{out}} \times H_{\text{out}} \times W_{\text{out}}}$ is the output of the layer. There are several components of CNNs.

1. **stride**: It changes the step size of a convolution. For example, if a stride is set to 2, the filter will move two pixel at a time.
2. **padding**: To pad the perimeter of image by some numbers, e.g. zero or the average of each pixel.
3. **valid convolution**: To not perform the convolution if a filter is not within the target image.
4. **same convolution**: To pad to keep the image size constant after the convolution.

Note that when a stride is S and a padding size P and the image size along one axis is X_{in} and the kernel size is K , then the output image size X_{out} along the axis is the following:

$$X_{\text{out}} = \left\lfloor \frac{X_{\text{in}} + 2P - K}{S} \right\rfloor + 1 \quad (33)$$

Additionally, there is a technique called **pooling**. Pooling is used to down-size an input image and it is invariant to small translations. Basically, **max**- or **average**-pooling is often used. For example, when we apply pooling to a given input, the formulation is the following:

$$\begin{aligned} \text{max pooling : } A_{h,w}^{\text{out}} &= \max_{1 \leq x \leq W', 1 \leq y \leq H'} A_{h+y, w+x}^{\text{in}} \\ \text{average pooling : } A_{h,w}^{\text{out}} &= \frac{1}{H'W'} \sum_{x=1}^{W'} \sum_{y=1}^{H'} A_{h+y, w+x}^{\text{in}} \end{aligned} \quad (34)$$

where $0 \leq h \leq H - H', 0 \leq w \leq W - W'$. Note that it does not have any weight parameters and it is not common to use zero-padding for pooling layers.

6.2 Properties of CNNs

The CNNs extract abstract features from an input and the following convolution layer builds on top of the extracted features, so CNNs can learn more complex features. Additionally, there are the following three major properties in CNN:

1. **Sparse interactions**: Each element in a current layer interacts with limited number of units in the following layer and this leads to **the sparsity** and **less memory complexity**¹⁷. However, as the layers go deeper, each unit receives much wider **receptive field**, because each element takes over multiple features from the previous layer and the number of interacted elements increases exponentially with respect to the depth.
2. **Parameter sharing**: Since CNNs reuse the same weights multiple times in each layer at different regions and each parameter tries to capture regular patterns, it **reduces the memory complexity** and leads to **the generalization**, i.e. regularization.
3. **Equivariant representation**¹⁸: For example, when we define $f(x)$ as CNN and $g(x)$ as location shift. Since $f(g(x)) = g(f(x))$ holds, CNNs can yield the same result in both the prediction on shifted x , i.e. $f(g(x))$, and the shifted prediction on non-shifted x , i.e. $g(f(x))$. In other words, it has **translation invariant**. This property comes from weight sharing. Note that CNNs are not naturally equivariant to some other transformations such as changes in the scale or rotation of the image.

6.3 Miscellaneous convolutions

1. **3D convolution**: Normal convolution is 2D, but this convolution is expanded to 3D. It is useful when we handle data such as videos (height, width, time), MRI or CT images (height, width, depth)
2. **1×1 convolution**: It is used for the purpose of reducing channels and thus we can compute forward pass efficiently from the following layers. Since it extracts features and down-sizes data, it is called **feature pooling** and enables more complicated representation using non-linear activation after the procedure.

¹⁷CNNs can reduce the memory complexity from $O(C_{\text{in}}C_{\text{out}}H_{\text{in}}W_{\text{in}}H_{\text{out}}W_{\text{out}})$ to $O(C_{\text{in}}C_{\text{out}}H'_{\text{in}}W'_{\text{in}})$.

¹⁸The definition: $f(g(x)) = g(f(x))$

3. **Transposed convolution:** This convolution up-samples data by combining the padding and convolution.
4. **Dilated convolution:** It inflates the kernel by inserting spaces between kernel elements while keeping the kernel size; therefore no additional cost is required. Since it convolutes elements from a larger range, we can often get much larger receptive field especially when we use multi-dilated convolutions that are stuck together.
5. **Grouped convolution:** Filters at each layer are separated into certain numbers of groups and each group is responsible for convolutions of the corresponding group¹⁹. The advantages of this convolution are efficient training (each group can be allocated to different GPUs) and efficient model in terms of the number of model parameters²⁰.
6. **Spatially separable convolutions:** It first convolutes only along one axis and then it convolutes the convoluted data along the other axis. The advantages of this convolution are less parameters and less matrix multiplications. On the other hand, training results can be sub-optimal.
7. **Depthwise separable convolutions:** It first convolutes each channel separately²¹, and then it applies 1×1 convolution to the data. This convolution also yields the same benefits and problems as spatially separable convolutions.

7 Recurrent neural networks (RNN)

7.1 General settings

RNNs have cyclic graph structures and typically skilled at dealing with sequential data, because it provides a time context and the concept of memory. RNNs implement dynamical systems rather than function mappings. Additionally, it is known as turing-complete²². There are several ways to use RNNs:

1. **one to many:** such as extracting the description (a sequence of words) of a given picture
2. **many to one:** such as giving a score (one real value) to a given review (a sequence of words)
3. **many to many:** such as translation from English (a sequence of words) to German (a sequence of words)

The forward pass of RNNs is the following:

$$\begin{aligned} \mathbf{h}_t &= \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_i \mathbf{x}_t + \mathbf{b}_h) \\ \hat{\mathbf{y}}_t &= g(\mathbf{o}_t) = g(\mathbf{b}_o + \mathbf{W}_o \mathbf{h}_t) \end{aligned} \tag{35}$$

where $\mathbf{h} \in \mathbb{R}^H$ is a hidden state vector, $\mathbf{x} \in \mathbb{R}^D$ is an input vector, $\mathbf{W}_h \in \mathbb{R}^{H \times H}$, $\mathbf{W}_i \in \mathbb{R}^{H \times D}$, $\mathbf{W}_o \in \mathbb{R}^{H \times K}$ are shared weight vectors for the hidden states, inputs and outputs, respectively, $\mathbf{b}_h \in \mathbb{R}^H$, $\mathbf{b}_o \in \mathbb{R}^K$ are shared bias vectors for the hidden states and outputs respectively, a is an activation function and K is the number of outputs for each time step. Note that the number of parameters does not depend on the sequence length and the memory complexity of the weight matrix is only $O(H(D+H))$.

¹⁹Conventional CNNs convolute all the input channels.

²⁰In the normal setting, the memory complexity of the kernel is $O(C_{\text{out}} C_{\text{in}} H'_{\text{in}} W'_{\text{in}})$. On the other hand, the kernel is $O(G C_{\text{out}} / G C_{\text{in}} / G H'_{\text{in}} W'_{\text{in}}) = O(C_{\text{out}} C_{\text{in}} H'_{\text{in}} W'_{\text{in}} / G)$ in the grouped convolution where G is the number of groups.

²¹More precisely, $Z_{B,c} = K_{c,1} \otimes A_{B,c}$ for all $c = 1, \dots, C_{\text{in}}$. We do not take the summation over input channels.

²²Roughly speaking, if a given system can solve all the problems, which can be resolved by current computer systems, it is said to be turing-complete.

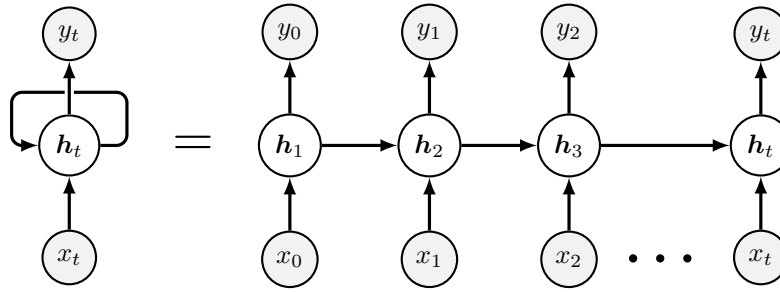


Figure 3: The visualization of RNN. The left figure is the short version and the right figure is the unfolded version.

7.2 Output-to-hidden connections and teacher forcing

When thinking about the back propagation through time (BPTT)²³ for Eq (35), we need to compute all the gradients with respect to h_t using the information from the successive time steps. Therefore, it cannot be computed in parallel and it is more complicated process. One solution is to connect between the output o_t and the following hidden state h_{t+1} instead of h_t and h_{t+1} . Since it lacks hidden-to-hidden recurrent connections, they are **likely to miss the necessary information about the past history** of the input. On the other hand, the advantage is that all the gradients computation for each time step can be decoupled and that is why they can be computed in parallel. One variant is **teacher forcing**. Teacher forcing has connections from labels y_t to the successive hidden state h_{t+1} during training and from predictions o_t to h_{t+1} during test. Since this method uses the labels directly, the **training is easier**. However, it approximates based on the predicted values during test and it is **likely to accumulate errors over time**.

7.3 Miscellaneous architectures

By changing Eq (35), we can easily obtain the formulations for one-to-many, many-to-many, many-to-one.

1. **many to one**: Just limit the approximation to only at the last step to summarize the information and train by BPTT
2. **one to many**: Use same x for each time step $t = 1, 2, \dots, T$
3. **many to many**: Can be Combined with BPTT and teacher forcing

As other variants, there are **bidirectional RNNs** and **encoder-decoder sequence-to-sequence architectures**. Bidirectional RNNs has literally two directional hidden states. One is the forward direction and the other is the backward direction. The advantage is to be able to use information from both past and future. On the other hand, since it requires future information, we cannot use it for online tasks. Encoder-decoder sequence-to-sequence architectures are used when the length of inputs and outputs are different such as translation.

²³BPTT is also weight sharing.

7.4 Back propagation

7.4.1 Computation

Suppose we would like to solve many-to-many tasks and we classify into K class at each time step. Then the forward pass at time step t is the following:

$$\begin{aligned}
 \mathbf{h}_t &= \tanh(\mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_i \mathbf{x}_t + \mathbf{b}_h) \\
 \mathbf{o}_t &= \mathbf{b}_o + \mathbf{W}_o \mathbf{h}_t \\
 \hat{\mathbf{y}}_t &= \text{softmax}(\mathbf{o}_t) \\
 L_t &= L(\hat{\mathbf{y}}_t, \mathbf{y}_t) \\
 L &= \sum_{t=1}^T L_t
 \end{aligned} \tag{36}$$

Therefore, the backward pass at time step t , which does not require the whole pass, are the followings:

$$\begin{aligned}
 \frac{\partial L}{\partial L_t} &= 1 \\
 \frac{\partial L_t}{\partial \mathbf{o}_t} &= \hat{\mathbf{y}}_t - \mathbf{y}_t \\
 \frac{\partial L_t}{\partial \mathbf{W}_o} &= \frac{\partial L_t}{\partial \mathbf{o}_t} \frac{\partial \mathbf{o}_t}{\partial \mathbf{W}_o} = \frac{\partial L_t}{\partial \mathbf{o}_t} \mathbf{h}_t^\top \left(W_{o \ i, j} = \frac{\partial L_t}{\partial o_{t,i}} h_{t,j} \right) \\
 \frac{\partial L_t}{\partial \mathbf{b}_o} &= \frac{\partial L_t}{\partial \mathbf{o}_t} \frac{\partial \mathbf{o}_t}{\partial \mathbf{b}_o} = \frac{\partial L_t}{\partial \mathbf{o}_t} \left(b_{o \ i} = \frac{\partial L_t}{\partial o_{t,i}} \right)
 \end{aligned} \tag{37}$$

The gradients with respect to $\mathbf{h}_t, \mathbf{h}_{t-1}$ are the following:

$$\begin{aligned}
 \frac{\partial L_t}{\partial \mathbf{h}_t} &= \frac{\partial L_t}{\partial \mathbf{o}_t} \frac{\partial \mathbf{o}_t}{\partial \mathbf{h}_t} = \mathbf{W}_o^\top \frac{\partial L_t}{\partial \mathbf{o}_t} \left(\because \frac{\partial L_t}{\partial h_{t,i}} = \sum_{j=1}^K W_{o,j,i} o_{t,j} \right) \\
 \frac{\partial L_t}{\partial \mathbf{h}_{t-1}} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{h}_{t-1}} = \mathbf{W}_h^\top \left(\frac{\partial L_t}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2) \right) \left(\because \frac{\partial L_t}{\partial h_{t-1}} = \sum_{j=1}^H W_{h \ j,i} \frac{\partial L_t}{\partial h_{t,j}} (1 - h_{t,j}^2) \right)
 \end{aligned} \tag{38}$$

where $\frac{\partial y}{\partial x} = (1 - y^2)w$ when $y = \tanh(wx + b)$. Using the gradients above, other gradients can be computed as follows:

$$\begin{aligned}
 \frac{\partial L_t}{\partial \mathbf{W}_h} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_h} = \frac{\partial L_t}{\partial \mathbf{h}_t} (\mathbf{1} - \mathbf{h}_t^2) \mathbf{h}_{t-1}^\top = \left(\frac{\partial L_t}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2) \right) \mathbf{h}_{t-1}^\top \left(\because \frac{\partial L_t}{\partial W_{h \ i, j}} = \frac{\partial L_t}{\partial h_{t,i}} (1 - h_{t,i}^2) h_{t-1,j} \right) \\
 \frac{\partial L_t}{\partial \mathbf{W}_i} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{W}_i} = \frac{\partial L_t}{\partial \mathbf{h}_t} (\mathbf{1} - \mathbf{h}_t^2) \mathbf{x}_t^\top = \left(\frac{\partial L_t}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2) \right) \mathbf{x}_t^\top \left(\because \frac{\partial L_t}{\partial W_{i \ i, j}} = \frac{\partial L_t}{\partial h_{t,i}} (1 - h_{t,i}^2) x_{t,j} \right) \\
 \frac{\partial L_t}{\partial \mathbf{b}_h} &= \frac{\partial L_t}{\partial \mathbf{h}_t} \frac{\partial \mathbf{h}_t}{\partial \mathbf{b}_h} = \frac{\partial L_t}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2)^\top \left(\because \frac{\partial L_t}{\partial b_{h,i}} = \frac{\partial L_t}{\partial h_{t,i}} (1 - h_{t,i}^2) \right)
 \end{aligned} \tag{39}$$

In the whole pass, we use $\frac{\partial L}{\partial \cdot} = \sum_{t=1}^T \frac{\partial L}{\partial L_t} \frac{\partial L_t}{\partial \cdot}$ and the gradients of hidden states are the followings:

$$\frac{\partial L}{\partial \mathbf{h}_t} = \frac{\partial L_t}{\partial \mathbf{h}_t} + \frac{\partial L_{t+1}}{\partial \mathbf{h}_t} = \mathbf{W}_o^\top \frac{\partial L}{\partial \mathbf{o}_t} + \mathbf{W}_h^\top \left(\frac{\partial L}{\partial \mathbf{h}_{t+1}} \odot (\mathbf{1} - \mathbf{h}_{t+1}^2) \right) \tag{40}$$

After the whole computation of the gradients of loss with respect to hidden states, we can compute the gradients of each matrix as follows:

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{W}_h} &= \sum_{t=2}^T \left(\frac{\partial L}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2) \right) \mathbf{h}_{t-1}^\top \\ \frac{\partial L}{\partial \mathbf{W}_i} &= \sum_{t=1}^T \left(\frac{\partial L}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2) \right) \mathbf{x}_t^\top \\ \frac{\partial L}{\partial \mathbf{b}_h} &= \sum_{t=1}^T \frac{\partial L}{\partial \mathbf{h}_t} \odot (\mathbf{1} - \mathbf{h}_t^2)^\top\end{aligned}\tag{41}$$

7.4.2 Gradient vanishing and explosion

Although RNNs have powerful non-linear processing due to the repeating application of the same function, it also induces some problems. Typical problems are the gradient vanishing and explosion. This is due to the formulation of RNNs.

$$\begin{aligned}\frac{\partial L}{\partial \mathbf{h}_t} &\propto \mathbf{W}_h^\top \left(\frac{\partial L}{\partial \mathbf{h}_{t+1}} \odot (\mathbf{1} - \mathbf{h}_{t+1}^2) \right) \\ &\propto \mathbf{W}_h^\top \left(\mathbf{W}_h^\top \left(\frac{\partial L}{\partial \mathbf{h}_{t+2}} \odot (\mathbf{1} - \mathbf{h}_{t+2}^2) \right) \odot (\mathbf{1} - \mathbf{h}_{t+1}^2) \right) \\ &\dots \\ &\propto \left((\mathbf{W}_h^\top)^{T-t} \bigodot_{\tau=t} (\mathbf{1} - \mathbf{h}_\tau^2) \right) \frac{\partial L}{\partial \mathbf{h}_T}\end{aligned}\tag{42}$$

Let $\mathbf{W}_h = \mathbf{V}\mathbf{\Lambda}\mathbf{V}^{-1}$ be the eigendecomposition of the weight matrix. In this case, the following holds:

$$\frac{\partial L}{\partial \mathbf{h}_t} \propto (\mathbf{1} - \mathbf{h}^2)^{T-t} \mathbf{V} \mathbf{\Lambda}^{T-t} \mathbf{V}^{-1} \frac{\partial L}{\partial \mathbf{h}_T}\tag{43}$$

Therefore, if the maximum or minimum of absolute eigenvalue is larger or smaller than 1, it leads to explosion or vanishing. The solution for the explosion is **gradient clipping** and that for the vanishing is **LSTM**.

7.5 Long Short-Term Memory (LSTM)

LSTM is invented as a solution for the vanishing gradient. This model is composed of three types of gates, i.e. **forget gate**, **input gate**, **output gate**, as visualized in Figure 4. Since this model has a self-loop of inner state \mathbf{c} without any activations, the gradient with respect to \mathbf{c} does not vanish nearly as RNNs do.

7.5.1 Structure

The formulation is as follows:

$$\begin{aligned}\begin{bmatrix} \mathbf{f}_t^\top & \mathbf{i}_t^\top & \tilde{\mathbf{c}}_t^\top & \mathbf{o}_t^\top \end{bmatrix}^\top &= \mathbf{W}_h \mathbf{h}_{t-1} + \mathbf{W}_i \mathbf{x}_t + \mathbf{b}_h \\ \mathbf{c}_t &= \sigma(\mathbf{f}_t) \odot \mathbf{c}_{t-1} + \sigma(\mathbf{i}_t) \odot \tanh(\tilde{\mathbf{c}}_t) \\ \mathbf{h}_t &= \tanh(\mathbf{c}_t) \odot \sigma(\mathbf{o}_t) \\ \hat{\mathbf{y}}_t &= a(\mathbf{W}_o \mathbf{h}_t + \mathbf{b}_o)\end{aligned}\tag{44}$$

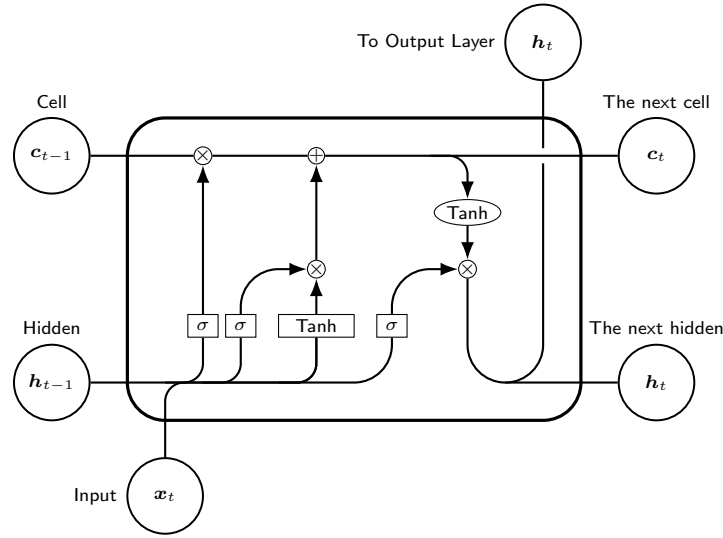


Figure 4: The visualization of LSTM.

where $\mathbf{f}, \mathbf{i}, \tilde{\mathbf{c}}, \mathbf{o} \in \mathbb{R}^H$. Since the derivatives move from \mathbf{c}_t to $\mathbf{f}_t, \mathbf{i}_t, \tilde{\mathbf{c}}_t$ and their derivatives move to \mathbf{h}_{t-1} and $\frac{\partial \mathbf{c}_t}{\partial \mathbf{c}_{t-1}} = \sigma(\mathbf{f}_t)$, the decaying rate is far better than that of Eq (43). The forget gate decays less important features from the past and maintain important features. The input gate takes two vectors \mathbf{i}_t and $\tilde{\mathbf{c}}_t$ and they control how much new information we should store and in which direction we should change the state. The output gate controls which features to use. The advantages of LSTM are that it can **selectively forget and remember features** and **each parameter of LSTM already have their meanings**. For this reason, the training result can be easily interpreted and LSTM can yield good performance with less parameters and less training.

7.6 Alternatives

1. **Gated recurrent units (GRU)**: Since it combines the forgetting factor and the new contribution to the state update in one gate, it has slightly simpler and fewer parameters compared to LSTM.
2. **Echo state networks**: It uses a large fixed hidden part called reservoir with random weights and only trains output layer with linear regression. The idea is to use the rich connections and rich dynamics in reservoir. Even though it is simple and quick, it shows strong performance. Note that we need to pay attention to the spectral radius ²⁴,
3. **Neural turing machines**: This is an extension of RNN using the concept called memory cells. This method reads and writes the features of tasks in memory cells and this mechanism is called **attention**. It is skilled at dealing with tasks requiring long-term memories.

8 Practical methodology

8.1 Attention mechanism

Attention mechanism is to pay attention to more important features in given data rather than dealing with all the inputs equally. It is especially helpful to process big images and it can reduce the computational complexity. For example, when we have images of 256×256 pixel size, but if they

²⁴The largest absolute eigenvalues of connectivity matrix should be smaller than 1.

only have information in the center of the size 32×32 , it can reduce the complexity significantly. Such attention mechanism is realized by weighting each feature differently. **Soft attention** and **hard attention** are one of the attention techniques. Soft attention first extracts feature using CNN and RNN learns weights for each feature at a given time step. Therefore, soft attention can be used for image captioning. Since soft attention requires all the pixels and puts weights on each pixel, it does not reduce the computational time.

For this reason, the hard attention has been introduced. Hard attention pays attention to only one location and since the sample of locations, i.e. the positions of cropping, cannot be learned using back propagation, such a location sampling is learned via **reinforcement learning**²⁵. Hard attention can be used for letter recognition. On the other hand, soft attention uses weighted sum of each feature from each location and **learns such weights using back propagation**. Note that there have been differentiable hard attention methods and one of them is **spatial transformer**. This method trains parameterized function that transforms inputs into cropped (, so smaller than the original inputs) inputs using back propagation and interpolates the transformed inputs using the bilinear interpolation if required.

8.2 Detection and segmentation

Semantic segmentation is a pixel-wise labeling and CNNs for each pixel are trained in the most naïve setting. However, the pixel-wise training is computationally expensive, so recent research suggests to extract feature first while maintaining the image size and give the likelihoods in the final layer. Although it can reduce the computational time, we need to maintain the resolution of the original images. For this reason, we combine the feature extraction of CNNs and **transpose convolution**, i.e. upsampling.

Object detection is a combination of classification and localization. Since it has to optimize two tasks, it will optimize the multitask loss of softmax loss from a classification task and L2 loss from a localization task. Note that localization task is to minimize the error between prediction of the vertices of bounding box and the central position and the actual location of the bounding box. In the case of multiple objects, we apply the task to each sub-image chosen by the sliding window manner. However, this computation is expensive due to the huge candidate pool²⁶. For this reason, Region CNN (RCNN) and You Only Look Once (YOLO) have been invented. RCNN learns region proposal network and searches from each region. YOLO is an end-to-end learning and first divides images into grid cells, then it approximates the shape of bounding box and the class label for each cell. Since grid cells are coarse, the accuracy is lower, but still achieves good performance.

As other applications, there are **instance segmentation** and **panoptic segmentation**. Instance segmentation not only classifies each pixel into a specific class, but also detects if each cell belongs to one specific object (e.g. human, car) in a given image. Panoptic segmentation is the combination of instance segmentation and the regular segmentation. It applies instance segmentation to, for example, people or cars and uses regular segmentation for road and buildings.

8.3 CNN architectures

There are 10 types of famous architectures:

1. **LeNet5**: The first CNN model.
2. **AlexNet**: The CNN used in 2012.
3. **ZFNet**: A variant of AlexNet.
4. **VGGNet**: Small filter and deeper layers.

²⁵Image is a state and a sample is an action. Based on the result, the sample will get a reward.

²⁶Combinations of locations, scales, aspect ratio and so on.



Figure 5: The visualizations of semantic segmentation and object detection.



Figure 6: The visualizations of instance segmentation and panoptic detection.

5. **GoogLeNet**: Inception module and less parameters. Inception module uses multiple kernel size convolution and concatenate all outputs depth-wisely.
6. **ResNet**: Residual block (a solution for the gradient vanishing) and deeper layer. Each residual block doubles the number of filters and halves pixel size. It also uses bottleneck layer (1×1 convolution layers)
7. **WideResNet**: It reduces the layers and widens channel size, so parallelizable.
8. **ResNeXt**: It parallelizes residual blocks.
9. **DenseNet**: It handles gradient vanishing by taking each input at each layer in successive layers and encourages feature reuse.
10. **MobileNet**: It uses depthwise separable convolution, so faster computation.
11. **SENet (Squeeze-and-Excitation Networks)**: A feature recalibration module that learns to adaptively reweight feature maps

In the previous layer of fully connected layer, CNN typically requires flatten operation and this size depends on the input image size. However, models such as GoogLeNet and ResNet uses global pooling before the layer, so we do not need to consider the pixel size anymore.

8.4 Batch normalization

Batch normalization is one of the **regularizations** and makes training easier due to the **handling of the internal covariate shift** and zero-mean of the activated values which **prevent gradients vanishing/exploding** and allow higher learning rates and **faster convergence** ²⁷. Additionally, since batch normalization can activate dead neurons, it improves gradient flow. Another benefit is **less**

²⁷Since each weight is normalized, the higher weight values are not likely to be a problem; therefore, we can take a higher learning rate.

Table 2: The case study of fine-tuning

Similarity / size	Little	Large
Similar	Trains only final layer	Fine-tunes a few layers
Not Similar	Fine-tunes a larger number of layers	Tries a linear layer from a shallower part

sensitivity to the weight initialization, because it normalizes values at each layer. The operation of batch normalization is performed by the following:

$$\text{BN}(x) = \gamma \frac{x - \mu}{\sigma} + \beta$$

$$\text{where } \mu = \frac{1}{B} \sum_{b=1}^B x_b, \sigma^2 = \frac{1}{B} \sum_{b=1}^B (x_b - \mu)^2 \quad (45)$$

where β, γ are learnable parameters and when $\beta = \mu, \gamma = \sigma$, batch normalization becomes an identity function. In the case of MLP, each dimension is processed independently while each pixel is processed jointly, i.e. taking the mean and variance over the batch of whole images, in CNNs. Note that batch normalization is used after fully-connected or convolutional layer to undo the scale and before nonlinearity to reduce dead neurons. Additionally, during testing, batch normalization uses μ and σ obtained by averaging values during training. For this reason, batch normalization exhibits different behaviors depending on either during training or testing.

The major problems of batch normalization are the handling of **mini-batch learning**, because of the unstable statistics and **recurrent neural structures**, because each time step has different statistics and requires another batch normalization layer. To address this issue, the following normalizations are invented:

1. **Batch normalization** (batch, pixel): Basic normalization. In the case of MLP, handles each dimension independently.
2. **Layer normalization** (pixel, all channels): Used for RNN.
3. **Instance normalization** (pixel): Used for style transfer and GAN.
4. **Group normalization** (pixel, some channels): Composing a group by some channels

Since all of them do not take statistics over batch, they behave in the identical way during both training and test. Note that all the normalization techniques take the normalization with respect to each dimension independently in MLPs.

8.5 Transfer learning

Transfer learning or fine-tuning is to use pre-trained weights as parts of neural networks for new tasks. Since the shallower parts of weights are more generic and the deeper parts of weights are more specific, we only initialize deeper layers and train only the deeper parts of networks. Practically, we use lower learning rate and train more layers for bigger dataset. Table 2 lists case-by-case solutions. Fine-tuning is **effective especially when we do not have enough data**.

Another transfer learning is self-supervised learning and performed as follows:

1. creates pretext tasks using massive unlabeled data
2. solves the pretext tasks
3. evaluates the performance on a labeled dataset

Note that the pretext tasks are tasks generated from unlabeled data easily. For example, rotation detection can be generated easily, because we just applies rotation operation to each data and such operations are the answers.

8.6 Debugging strategies

To quickly debug our model, we first train on a tiny dataset and keep modifying the model until reaching good performance. Another strategy is to monitor quantitative metrics such as loss function and visualize the inference on feeded data. When the training does not work well, one should monitor gradients and activations to know how much neurons are dead and compare the gradients with finite differential equation.

9 Hyperparameter optimization (HPO)

9.1 Practical tips

When we apply HPO, we have to know there are **numerical**, **discrete**, **categorical** and **conditional** parameters and consider the scale (e.g. **log and uniform**) of discrete and numerical parameters. This knowledge is important when designing the searching space. When the model shows overfitting, we have to check **capacity of the model**, **HPO methods** and **degree of regularization**. On the other hand, when the model shows underfitting, we have to check **different regularizations**, **double descent phenomenon**²⁸ or **train on more data** if available.

In general, HPO is useless unless the model can overfit the training dataset. For this reason, one good strategy is to make sure that the model (i.e. training loss) can overfit the subset and then move to the HPO of validation loss on the full dataset. As another tip, the learning rate is typically important, so one heuristic is to increase learning rate until the loss diverges and reduce it a little bit.

9.2 Black-box optimization

Black-box optimization is to optimize functions without derivative information. The most basic algorithm is random search and it maintains global search and can be computed in parallel. Additionally, it can deal with low intrinsic dimensionality compared to grid search. Another famous example is **local search** and it fixes the best configuration and changes the configuration along only a single axis.

Those methods are naïve methods and not sample-efficient. Therefore, there are several methods using previous observations. One example is bayesian optimization and it fits a probabilistic model $\mathcal{N}(\mu(\theta), \sigma^2(\theta))$ where θ is a configuration. Typical bayesian optimization balances the trade-off between exploration and exploitation and finds the next configurations to be evaluated from the surrogate searching space. Since bayesian optimization itself is not cheap compared to random search, it is used when the objective function is expensive. One major example of bayesian optimization is TPE and it can deal with conditional and categorical parameters. Another example is a population-based method and it maintains diversity and improves the fitness of the whole population.

9.3 Speedup techniques of HPO

Since the objectives in HPO is usually expensive, many speedup techniques have been studied. In the lecture, meta-learning across datasets, extrapolation of learning curves and multi-fidelity optimization are introduced. Note that multi-fidelity optimization is to use cheap-to-evaluate proxy instead of the high-fidelity objective. The concrete examples are the early-stopping of poor performance and removing poor configurations based on the training on a tiny dataset.

²⁸The phenomenon where the loss goes down when increasing the model capacity and goes up at some points, but when we keep increasing the capacity, the loss goes down again.

9.4 Hyperparameter gradient descent

Recently, the gradient-based HPO method has been invented and it is formulated as follows:

$$\begin{aligned} & \min_{\theta} L_{\text{val}}(\mathbf{w}^*(\theta), \theta) \\ & \text{where } \mathbf{w}^*(\theta) \in \operatorname{argmin}_{\mathbf{w}} L_{\text{train}}(\mathbf{w}, \theta) \end{aligned} \quad (46)$$

Note that the optimization of weights is called **inner objective** and that of hyperparameters is called **outer objective**. Although it requires twice much computation and there is no guarantee of the convergence, it can potentially deal with more hyperparameters. Additionally, the bi-level optimization allows quicker optimization and it updates weights at training step and hyperparameters at validation step.

10 Neural architecture search (NAS)

In this section, we will discuss how to find the optimal neural architecture. Note that since NAS requires substantial amount of computation and neural architecture does not make big difference compared to hyperparameter configurations, we recommend to apply HPO first unless we have adequate resources.

10.1 Searching space

10.1.1 Basic search space

NAS optimizes the objective $f(a)$ where $a \in \mathcal{A}$ is architecture composed of **micro**, **macro** and **hierarchical** architectures. The most basic searching space is to align the components in the chain-structure. Another searching space is to build a graph by choosing such components. However, since the selection of such components leads to hard combinatorial optimization problems, it is often infeasible. Therefore, there have been many papers introducing searching spaces.

10.1.2 Cell search space

The cells are made up of B blocks ²⁹ (e.g. the connection of hidden layers and kernel size) and we obtain architectures by stacking such cells as shown in Figure 7. There are two types of cells. One is a normal cell that maintains the shape of inputs and the other is a reduction cell that reduces the shape of inputs. Since the searching of the optimal cell yields a smaller searching space compared to architecture space, it leads to **quicker searching**. Additionally, the cells can be **transferred to other tasks**. Although stacking repeating patterns are proven to be a useful design (e.g. ResNet), it can miss the possibilities that **the different cells work well in the different parts** of the architecture. Another drawback is that we need to determine the macro architecture.

10.1.3 Hierarchical representation

Hierarchical representation searches different levels of components and assembles these components as directed acyclic graph (DAG) in the end. The level-1 primitives are the component such as 3×3 convolution and max pooling. The level $n + 1$ motifs are the combinations of level n motifs. Since the hierarchical representation assembles the motifs, it **can eliminate conditional parameters and make them categorical parameters**. This searching space also allows repeated patterns in the architecture and to reuse various levels of motifs. The major advantage over the cell search space is more flexibility and it **allows different blocks in different parts** of the network. On the other hand, since it is more flexible, the **searching space is larger** than the cell searching space; therefore, the optimization may potentially much harder.

²⁹The number of hyperparameters is the choices of each block times B .

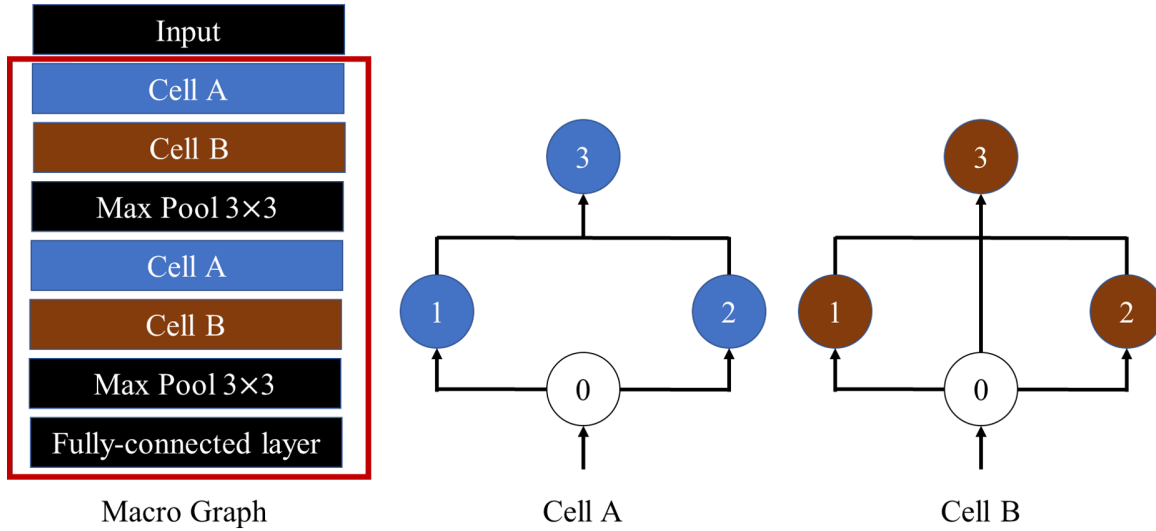


Figure 7: The example of macro graph and cells. Macro architecture is built by assembling cells and other components. Cell is composed of operations such as 3×3 convolution and connections.

10.2 Black-box optimization

Neural architecture can also be optimized by black-box optimization. The major difference from hyperparameter optimization is the existence of many conditional parameters and categorical parameters.

One of the earliest works is **neuroevolution**. Neuroevolution **evolves architecture jointly with its weights**. It first initializes a population with N random architectures. Then it repeats sampling with replacement from the population based on the fitness and while applying mutation to the N individuals. Note that mutations include adding, changing or removing a layer and we can choose elitism to preserve the elites. Another method is **Bayesian optimization** such as TPE, SMAC and NASBOT³⁰. They **jointly optimize architecture and hyperparameters**. **Regularized evolution** also shows good performance. It promotes the exploration of architectures by dropping older architectures from the population and updates the weights by SGD. The regularized evolution encodes the architecture space by categorical hyperparameters and it is known that it also works for Bayesian optimization. Additionally, the combination of **reinforcement learning (RL) and recurrent neural networks (RNN)** has been studied actively these days. The RNN called **controller** generates a neural architecture (called **child net**) as a sequence³¹ and evaluates the child net on a validation dataset. Then the weights of RNN is updated using RL based on the reward that is the performance on the validation dataset.

10.3 Speed-up techniques

Since NAS requires substantial amount of time, it is essential to reduce computational time. The major methods are **multi-fidelity optimization**, **learning curve prediction** and **meta-learning**³². Other examples are **network morphisms** and **weight sharing and one-shot models**. Network morphisms adds new operations $o(\mathbf{x}; \mathbf{w}_{\text{new}})$ on the existing architecture $\mathcal{A}(\mathbf{x}; \mathbf{w})$ to satisfy the following property:

$$\mathcal{A}(\mathbf{x}; \mathbf{w}) = o(\mathcal{A}(\mathbf{x}; \mathbf{w}); \mathbf{w}_{\text{new}}) \quad (47)$$

³⁰It uses arc-kernel that is invented for NAS and can handle conditional parameters.

³¹Each state is the hidden states and each action is to take a component. RNN approximates the Q-value function taking a hidden state and a component as arguments.

³²Since there are not many datasets for DL, there are only little work on meta-learning for NAS.

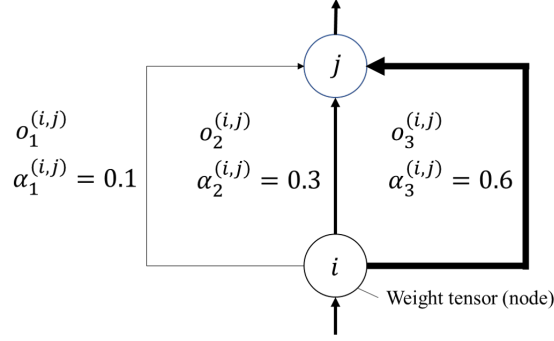


Figure 8: One-shot models set all the possible operations ($o_1^{(i,j)}, o_2^{(i,j)}, o_3^{(i,j)}$) as shown in the figure. Other NAS algorithms chooses only one of the operations (hard choice). However, DARTS chooses the operation by soft choice. In the implementation, the DARTS takes the weighted sum of the outputs from each operation during searching and chooses the best operation in the end.

In other words, the $o(\mathbf{x}; \mathbf{w}_{\text{new}})$ is an identity function. By transferring the weights from the previous architecture and training jointly with the new operation $o(\cdot)$, we can save the large amount of training time. The drawback of this method is the growing size of the architecture³³. Therefore, the multi-objective (the performance and the computational time) optimization is required for the network morphisms. The last method is weight sharing and one-shot models. One-shot models first construct the architecture graph, called **supergraph**, with all the possible operations, that are **edges**, and the weights in between those operations that are called **nodes**. Each subgraph shares the same weights with the supergraph. Then we search the optimal subgraph among the supergraphs while training the supergraph. Since we need to train the supergraph only once, The total computational cost can be significantly reduced, Note that since the HPO problems do not handle architecture parameters, network morphisms and one-shot models are not used in HPO.

10.4 Differentiable architecture search (DARTS)

Since the choices of operation and conditional parameters are not differentiable, it is not possible to apply gradient methods directly. Therefore, the DARTS introduces the following **stochastic relaxation**:

$$\mathbf{z}^{(j)} = \sum_{i < j} \tilde{o}^{(i,j)}(\mathbf{z}^{(i)}) = \sum_{i < j} \sum_{o \in \mathcal{O}^{(i,j)}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}^{(i,j)}} \exp(\alpha_{o'}^{(i,j)})} o(\mathbf{z}^{(i)}) \quad (48)$$

where $\mathbf{z}^{(j)}$ is the feature vector after passing the j -th node weight tensor, $\mathcal{O}^{(i,j)}$ is the set of operations between the i - and j -th nodes and $\alpha_o^{(i,j)}$ is the weight of the operation o . Roughly speaking, DARTS uses soft-choice shown in Figure 8 of operations instead of hard-choice and that is why the architecture is differentiable as in the case of softmax loss. In the end, this method chooses the operation that has the largest $\alpha_o^{(i,j)}$ among the set $\mathcal{O}^{(i,j)}$ and discretizes the architecture. The major issues of DARTS are the **hyperparameter sensitivity** and **memory consumption**. The first issue is addressed by early-stopping based on the curvature³⁴ of the validation loss or adversarial training. The memory issue is addressed by keeping only a single architecture or performing the search on a subset of the channels in the one-shot model.

³³Since it is hard to yield an identity mapping by reducing the size, the network morphisms realize the mapping by adding components.

³⁴The curvature is computed from eigenvalues of the Hessian matrix. Since DARTS discretizes in the end and the architecture may stay close to the sharp local minima, it can lead to degenerate architecture.

Algorithm 1 DARTS (Bi-level optimization)

```

1: function DARTS                                     ▷ Each update is performed by SGD.
2:   while The budget is left do
3:      $\mathbf{w} = \mathbf{w} - \frac{\partial L_{\text{train}}(\mathbf{w}, \boldsymbol{\alpha})}{\partial \mathbf{w}}$                                      ▷ Training (Inner loop)
4:      $\boldsymbol{\alpha} = \boldsymbol{\alpha} - \frac{\partial L_{\text{valid}}(\mathbf{w}, \boldsymbol{\alpha})}{\partial \boldsymbol{\alpha}}$                                      ▷ Architecture search (Outer loop)
5:   return  $\text{argmax}_{o \in \mathcal{O}(i,j)} \alpha_o^{(i,j)}$  for each edge  $(i, j)$ 

```

11 Generative models

11.1 Auto-encoder

Auto-encoder is a generative model that reconstruct an input \mathbf{x} by two neural networks. Auto-encoder first encodes an input \mathbf{x} into latent representation $\mathbf{z} = f^e(\mathbf{x}; \mathbf{w}_e)$. This representation is compressed lower dimensional representation of the input. Then it decodes the latent representation and reconstructs the original input $\hat{\mathbf{x}} = f^d(\mathbf{z}; \mathbf{w}_d)$. The neural networks learn the discrepancy between the original input and the reconstructed input. Therefore, it does not require labels and it is called **unsupervised learning**. The conceptual visualization is shown in Figure 9.

11.1.1 Regularized auto-encoders

Regularized auto-encoder has higher dimension of the latent space than the original input (called **over-complete representation**) and it puts a penalty on the magnitude of the variables by regularization. By the regularization, it throws out not important information and prevents pure copy. This model can choose the capacity based on the complexity of the data distribution.

11.1.2 Sparse auto-encoder

Sparse auto-encoder has the sparsity penalty as a regularization. The goal is to maximize the likelihood of the joint distribution of the reconstructed input and the latent vector. The joint distribution is formulated as follows:

$$p(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p(\mathbf{x}|\mathbf{z}) \quad (49)$$

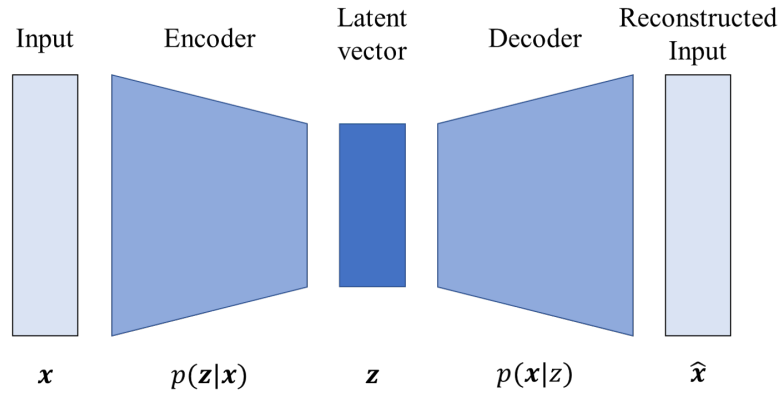


Figure 9: The conceptual visualization of auto-encoder. Both encoder and decoder are neural networks.

where $p(\mathbf{z})$ is a prior distribution and $p(\mathbf{x}|\mathbf{z})$ is a reconstructed distribution or likelihood given \mathbf{z} ³⁵. We obtain the following log-likelihood by marginalization:

$$\log p(\mathbf{x}) = \log \int p(\mathbf{x}, \mathbf{z}) d\mathbf{z} \quad (50)$$

In essence, the latent vector is determined by a parametric model, so we maximize the following:

$$\log p(\mathbf{x}, \mathbf{z}) = \log p(\mathbf{x}|\mathbf{z}) + \log p(\mathbf{z}) \quad (51)$$

We can achieve the sparsity by introducing Laplace prior $p(\mathbf{z}) = \frac{\lambda}{2} \exp(-\lambda|\mathbf{z}|)$. In general, the sparsity often increases generalization and avoids overfitting. This auto-encoder is often used as an unsupervised pre-processing step for downstream tasks. Note that downstream tasks are to transfer pre-trained weights to supervised learning tasks such as classification tasks.

11.1.3 Denoising auto-encoder

Denoising auto-encoder learns to cancel out the effect of noise corruption rather than adding an explicit regularization. It minimizes the loss function $L(\mathbf{x}, f^d(f^e(\tilde{\mathbf{x}})))$ where $\tilde{\mathbf{x}}$ is a corrupted input. To realize the corruption, we introduce a **corruption process** $C(\tilde{\mathbf{x}}|\mathbf{x})$ after the input. The optimization is performed by the gradient descent with respect to the following negative log-likelihood:

$$-\mathbb{E}_{\mathbf{x} \sim \mathcal{D}, \tilde{\mathbf{x}} \sim C(\tilde{\mathbf{x}}|\mathbf{x})} \log p(\mathbf{x}|\mathbf{z} = f^e(\tilde{\mathbf{x}})) \quad (52)$$

As a result, we obtain the reconstructed distribution $p(\mathbf{x}|\tilde{\mathbf{x}})$. Note that the latent vector is deterministic in this method, but **if we make it stochastic, it is equivalent to variational auto-encoder**.

11.1.4 Variational auto-encoder (VAE)

In VAE, we assume that each data has important factors and VAE explicitly models such factors. In other words, \mathbf{z} reflects the **causal factors** $p(\mathbf{x}|\mathbf{z})$ that makes it easier to reconstruct $p(\mathbf{x}) = \int p(\mathbf{x}|\mathbf{z})p(\mathbf{z})d\mathbf{z}$ ³⁶. To achieve it, we minimize the KL-divergence between the posterior and the parametric distribution rather than the optimization of the log-likelihood. Once we have a good approximation of $p(\mathbf{x})$, we can generate new data that are similar, but different from the original data.

To obtain such a model, we first need to infer the posterior distribution of the latent vector $p(\mathbf{z}|\mathbf{x}) = \frac{p(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{\int p(\mathbf{x}, \mathbf{z})d\mathbf{z}}$. Since the integral in the denominator is intractable, we use **approximate inference**. The approximate inference uses a simpler parametric distribution $q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e)$ and minimizes the KL-divergence between $q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e)$ and $p(\mathbf{z}|\mathbf{x})$. The equation is formulated as follows:

$$\begin{aligned} D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e)||p(\mathbf{z}|\mathbf{x})) &= \int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \log \frac{q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e)}{p(\mathbf{z}|\mathbf{x})} d\mathbf{z} \\ &= \int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \left(\log q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) - \log \frac{p(\mathbf{x}, \mathbf{z})}{p(\mathbf{x})} \right) d\mathbf{z} \\ &= \underbrace{\int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \left(\log q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) - \log p(\mathbf{x}, \mathbf{z}) \right) d\mathbf{z}}_{= -\mathcal{L}(\mathbf{x}, q), (\log p(\mathbf{x}) \geq \mathcal{L}(\mathbf{x}, q))} + \log p(\mathbf{x}) \end{aligned} \quad (53)$$

Since $\mathcal{L}(\mathbf{x}, \mathbf{w}_e)$ is the lower bound of $\log p(\mathbf{x})$, it is called **evidence lower bound (ELBO)**. The goal is to minimize the KL-divergence and it is achieved when we maximize the ELBO. This is because

³⁵Strictly speaking, these \mathbf{x} is $\tilde{\mathbf{x}}$. However, we unite the notation to \mathbf{x} and consider the distribution as the likelihood.

³⁶As mentioned previously, this is equivalent to the likelihood of the data distribution given the latent vector.

$\log p(\mathbf{x})$ is a constant number and when we maximize the ELBO, KL divergence approaches zero. The ELBO is reformulated as follows:

$$\begin{aligned}
\mathcal{L}(\mathbf{w}_e, \mathbf{w}_d, \mathbf{x}) &= \int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \left(\log p(\mathbf{x}, \mathbf{z}; \mathbf{w}_d) - \log q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \right) d\mathbf{z} \\
&= \int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \left(\log p(\mathbf{x}|\mathbf{z}; \mathbf{w}_d) + \log p(\mathbf{z}) - \log q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \right) d\mathbf{z} \\
&= \underbrace{-D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) || p(\mathbf{z}))}_{\text{regularization term}} + \underbrace{\int q(\mathbf{z}|\mathbf{x}; \mathbf{w}_e) \log p(\mathbf{x}|\mathbf{z}; \mathbf{w}_d) d\mathbf{z}}_{\text{reconstruction term}}
\end{aligned} \tag{54}$$

where $p(\mathbf{z})$ is a prior distribution. Note that the deterministic VAE with respect to the latent vector does not have the self-regularization term. The naïve implementation requires sampling to take the expectation. To handle the issue, the **reparametrization trick** has been invented. Reparametrization trick parametrizes the parameters of q . For example, if q is a gaussian distribution, the parameters are mean vector $\boldsymbol{\mu}$ and covariance matrix $\boldsymbol{\sigma}$, so we parametrize them as follows:

$$\mathbf{z} = \boldsymbol{\mu}(\mathbf{x}) + \boldsymbol{\sigma}(\mathbf{x}) \odot \boldsymbol{\epsilon} \tag{55}$$

where $\boldsymbol{\epsilon}$ is a random number sampled from $\mathcal{N}(\mathbf{0}, \mathbf{I})$. In this case, the neural network for the encoder returns the mean vector and the covariance matrix. Since the forward pass is continuous and differentiable, we can apply the back propagation. The loss function is the following:

$$\begin{aligned}
-\mathcal{L}(\mathbf{w}_e, \mathbf{w}_d, \mathbf{X}) &= \sum_{i=1}^B \left(D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}_i) || p(\mathbf{z})) - \int q(\mathbf{z}|\mathbf{x}_i; \mathbf{w}_d) \log p(\mathbf{x}_i|\mathbf{z}; \mathbf{w}_d) d\mathbf{z} \right) \\
&= \sum_{i=1}^B \left(D_{\text{KL}}(q(\mathbf{z}|\mathbf{x}_i) || p(\mathbf{z})) - \sum_{d=1}^D \left(x_{i,d} \log \hat{x}_{i,d} + (1 - x_{i,d}) \log (1 - \hat{x}_{i,d}) \right) \right)
\end{aligned} \tag{56}$$

where B is the batch size and D is the dimension of inputs. Intuitively speaking, the likelihood $p(\mathbf{x}_i|\mathbf{z}; \mathbf{w}_d)$ becomes smaller when the reconstruction error, i.e. the second term, becomes larger. It corresponds to the property of the negative log-likelihood. Note that the KL-divergence can be analytically computed in the most cases of exponential family.

11.2 Generative Adversarial networks (GAN)

GAN is the model that jointly trains both generator and discriminator. The generator tries to fool the discriminator by the generated data and the discriminator tries to detect the real data as much as possible. When we define the output data by the generator as $f^g(\cdot; \mathbf{w}_g)$ and the likelihood of the real data as $f^d(\cdot; \mathbf{w}_d)$, the formulation is the following:

$$\begin{aligned}
\text{Overall} : & \min_{\mathbf{w}_g} \max_{\mathbf{w}_d} \left[\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \log f^d(\mathbf{x}; \mathbf{w}_d) + \mathbb{E}_{\mathbf{r} \sim p(\mathbf{r})} \log (1 - f^d(f^g(\mathbf{r}; \mathbf{w}_g); \mathbf{w}_d)) \right] \\
\text{Discriminator} : & \max_{\mathbf{w}_d} \left[\mathbb{E}_{\mathbf{x} \sim \mathcal{D}} \log f^d(\mathbf{x}; \mathbf{w}_d) + \mathbb{E}_{\mathbf{r} \sim p(\mathbf{r})} \log (1 - f^d(f^g(\mathbf{r}; \mathbf{w}_g); \mathbf{w}_d)) \right] \\
\text{Generator} : & \min_{\mathbf{w}_g} \left[\mathbb{E}_{\mathbf{r} \sim p(\mathbf{r})} \log (1 - f^d(f^g(\mathbf{r}; \mathbf{w}_g); \mathbf{w}_d)) \right]
\end{aligned} \tag{57}$$

where \mathbf{r} is a random noise vector and it is sampled each iteration. Note that the generative tasks are much more difficult than the classification tasks; therefore, we **train from the discriminator**. Since the discriminator requires the maximization, we apply gradient ascent for the update. The major issue of the naïve GAN is that the generator does not often work well in this setting. This is due to the

small gradients when the generator does not work well, i.e. $f^d(f^g(\mathbf{r}; \mathbf{w}_g); \mathbf{w}_d)$ is close to 0. To alleviate the problem, we replace the original objective for the generator with the following:

$$\max_{\mathbf{w}_g} \left[\mathbb{E}_{\mathbf{r} \sim p(\mathbf{r})} \log f^d(f^g(\mathbf{r}; \mathbf{w}_g); \mathbf{w}_d) \right] \quad (58)$$

This objective promotes the learning when the generator works poorly.

12 Uncertainty handling

In some applications such as medical areas, more reliable decision making is required. In the same way, deep learning also requires uncertainty handling in some cases. Even though typical classification tasks use softmax and it can be interpreted probabilistically, it shows underconfidence for typically seen data and can show overconfidence for unrecognizable data. In other words, we would like to set the confidence higher for seen data and lower for unseen data. Therefore, we need additional uncertainty measures. There are two types of uncertainties: **Aleatoric uncertainty**, that is the intrinsic observation noise, and **Epistemic uncertainty**, that is the model uncertainty and can be reduced by more data. We discuss the measure of the combination of epistemic and aleatoric.

12.1 Properties of the Gaussian distribution

The Gaussian distribution has major four properties and we introduce them for the discussion later. Throuout this section, the Gaussian distribution takes the following form:

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{1,1} & \boldsymbol{\Sigma}_{1,2} \\ \boldsymbol{\Sigma}_{2,1} & \boldsymbol{\Sigma}_{2,2} \end{bmatrix} \right) \quad (59)$$

where $\boldsymbol{\mu}_i \in \mathbb{R}^{d_i}$, $\boldsymbol{\Sigma}_{i,j} \in \mathbb{R}^{d_i \times d_j}$ for $i, j = 1, 2$. The first property is realized by the marginalization $p(\mathbf{x}_1) = \int p(\mathbf{x}_1, \mathbf{x}_2) d\mathbf{x}_2$ and it is called **closed under marginalization**³⁷.

Theorem 4

Let \mathbf{x}_1 and \mathbf{x}_2 be jointly Gaussian distributed:

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{1,1} & \boldsymbol{\Sigma}_{1,2} \\ \boldsymbol{\Sigma}_{2,1} & \boldsymbol{\Sigma}_{2,2} \end{bmatrix} \right) \quad (60)$$

Then $\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_1, \boldsymbol{\Sigma}_{1,1})$

The second property is **closed under conditioning** and it is realized by $p(\mathbf{x}_2|\mathbf{x}_1) = \frac{p(\mathbf{x}_1, \mathbf{x}_2)}{p(\mathbf{x}_1)}$.

Theorem 5

Let \mathbf{x}_1 and \mathbf{x}_2 be jointly Gaussian distributed:

$$\begin{bmatrix} \mathbf{x}_1 \\ \mathbf{x}_2 \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \boldsymbol{\mu}_1 \\ \boldsymbol{\mu}_2 \end{bmatrix}, \begin{bmatrix} \boldsymbol{\Sigma}_{1,1} & \boldsymbol{\Sigma}_{1,2} \\ \boldsymbol{\Sigma}_{2,1} & \boldsymbol{\Sigma}_{2,2} \end{bmatrix} \right) \quad (61)$$

Then $\mathbf{x}_2|\mathbf{x}_1 \sim \mathcal{N}(\boldsymbol{\mu}_2 + \boldsymbol{\Sigma}_{1,2}^\top \boldsymbol{\Sigma}_{1,1}^{-1}(\mathbf{x}_1 - \boldsymbol{\mu}_1), \boldsymbol{\Sigma}_{2,2} - \boldsymbol{\Sigma}_{1,2}^\top \boldsymbol{\Sigma}_{1,1}^{-1} \boldsymbol{\Sigma}_{1,2})$

Other properties are **closed under multiplication** and **linear maps**. Note that the Gaussian distribution is often used due to the convenient mathematical properties, but not because sampled data follow the distribution.

³⁷It means we yield Gaussian distribution even after applying the operation.

12.2 Bayesian linear regression

Bayesian linear regression gives an uncertainty measure to the linear regression. The linear regression assumes that the output \mathbf{y} follows the Gaussian distribution with a fixed variance σ and uses **maximum likelihood estimation** on $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$. On the other hand, bayesian linear regression computes the posterior distribution of the weights $p(\mathbf{w}|\mathbf{X}, \mathbf{y})$ using Bayes' theorem as follows:

$$\begin{aligned} \text{posterior} &= \frac{\text{likelihood} \times \text{prior}}{\text{marginal likelihood}} \\ p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &= \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} \end{aligned} \quad (62)$$

Since the denominator $p(\mathbf{y}|\mathbf{X})$ does not depend on the weight \mathbf{w} and the denominator generally requires a complicated integral, we use the maximum a posteriori (**MAP**) estimation³⁸. In the MAP estimation, we maximize the following:

$$\begin{aligned} \text{posterior} &\propto \text{likelihood} \times \text{prior} \\ p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &\propto p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w}) \end{aligned} \quad (63)$$

In this case, if we choose a prior that is conjugate to the likelihood, the posterior can often be analytically obtained. In the case of the likelihood following Gaussian distribution, the **conjugate prior** is Gaussian distribution. In other words, the following holds:

$$\begin{aligned} \text{prior} : \mathbf{w} &\sim \mathcal{N}(\mathbf{0}, \Sigma_{\text{prior}}) \\ \text{posterior} : p(\mathbf{w}|\mathbf{X}, \mathbf{y}) &\sim \mathcal{N}(\boldsymbol{\mu}_{\text{post}}, \Sigma_{\text{post}}) \end{aligned} \quad (64)$$

where Σ_{prior} is a covariance matrix for the prior distribution, σ is a pre-defined noise, $\boldsymbol{\mu}_{\text{post}} = \frac{1}{\sigma^2} \Sigma_{\text{post}} \mathbf{X}^\top \mathbf{y}$ and $\Sigma_{\text{post}} = \left(\frac{1}{\sigma^2} \mathbf{X}^\top \mathbf{X} + \Sigma_{\text{prior}}^{-1} \right)^{-1}$. These values can be analytically derived using log function on Eq (63). Note that since Σ_{prior} controls the regularization effect and affects the performance, the choice of Σ_{prior} is called **model selection** and it is usually performed by cross validation. Using the weights sampled from the posterior, the prediction of Bayesian linear regression is computed as follows:

$$\begin{aligned} p(\hat{y}|\mathbf{x}, \mathbf{X}, \mathbf{y}) &= \int p(\hat{y}|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{y})d\mathbf{w} \\ &= \mathcal{N}(\mathbf{x}^\top \boldsymbol{\mu}_{\text{post}}, \mathbf{x}^\top \Sigma_{\text{post}} \mathbf{x}) \end{aligned} \quad (65)$$

Note that we can realize the Bayesian non-linear regression by replacing the \mathbf{x} with non-linearly mapped elements such as $\Phi = [1, x, x^2, \dots, x^d]$ as long as **the computation is linear in \mathbf{w}** . Furthermore, although the Laplace prior is not conjugate to the Gaussian likelihood and it requires to solve non-linear system, we can introduce **the Laplace prior to promote the sparsity** of the model.

12.3 Neural networks with Bayesian

We discuss how we can use the linear regression in a Bayesian manner. We will expand the discussion and consider how we use the Bayesian concept for neural networks.

12.3.1 DNGO

DNGO first trains a neural network on a given dataset and use the feature vector obtained in the last hidden layer as basis functions of the Bayesian linear regression.

³⁸The maximum likelihood estimation is equivalent to the MAP estimation with the marginal likelihood that follows uniform distribution.

12.3.2 Bayesian neural networks

Bayesian neural networks learn the posterior of weights $p(\mathbf{w}|\mathcal{D})$ and predict by the following marginalization:

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathcal{D})d\mathbf{w} \quad (66)$$

The core part of Bayesian neural networks are to compute the posterior:

$$p(\mathbf{w}|\mathcal{D}) = p(\mathbf{w}|\mathbf{X}, \mathbf{y}) = \frac{p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})}{\int p(\mathbf{y}|\mathbf{X}, \mathbf{w})p(\mathbf{w})d\mathbf{w}} \quad (67)$$

Since the Gaussian distribution is not closed under non-linear mapping and neural networks include non-linear mapping, the likelihood $p(\mathbf{y}|\mathbf{X}, \mathbf{w})$ is complicated and the posterior approximation requires Markov-chain Monte-carlo method or variational inference. In the next section, we will discuss both methods.

12.4 Posterior inference for non-linear models

12.4.1 Variational inference

As in VAE, the variational inference minimizes the KL-divergence between the posterior $p(\mathbf{w}|\mathcal{D})$ and a simple parametric model $q(\mathbf{w}; \theta)$ as follows:

$$\begin{aligned} D_{\text{KL}}(q(\mathbf{w}; \theta) || p(\mathbf{w}|\mathcal{D})) &= \int q(\mathbf{w}; \theta) \log \frac{q(\mathbf{w}; \theta)}{p(\mathbf{w}|\mathcal{D})} d\mathbf{w} \\ &= \int q(\mathbf{w}; \theta) \left(\log q(\mathbf{w}; \theta) - \log \frac{p(\mathcal{D}|\mathbf{w})p(\mathbf{w})}{p(\mathcal{D})} \right) d\mathbf{w} \\ &= \int q(\mathbf{w}; \theta) \left(\log q(\mathbf{w}; \theta) - \log p(\mathcal{D}|\mathbf{w}) - \log p(\mathbf{w}) \right) d\mathbf{w} + \log p(\mathcal{D}) \\ &= D_{\text{KL}}(q(\mathbf{w}; \theta) || p(\mathbf{w})) - \underbrace{\mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \theta)} [\log p(\mathcal{D}|\mathbf{w})]}_{\text{log-likelihood}} + \underbrace{\log p(\mathcal{D})}_{\text{constant}} \\ &= - \underbrace{(\mathbb{E}_{\mathbf{w} \sim q(\mathbf{w}; \theta)} [\log p(\mathcal{D}|\mathbf{w})] - D_{\text{KL}}(q(\mathbf{w}; \theta) || p(\mathbf{w})))}_{\text{ELBO: } \mathcal{L}(\theta)} + \text{const} \end{aligned} \quad (68)$$

Since the lower bound of KL-divergence is zero and the sum of the KL-divergence and the ELBO is constant, the minimization of the KL-divergence is realized by the maximization of the ELBO $\mathcal{L}(\theta)$. Furthermore, the KL-divergence term works as regularization, because we have to make the distribution closer to the prior. In contrast to MCMC, the variational inference avoids sampling by taking a parametric model and approximates the model by the maximization of ELBO. Note that the implementation requires the computation over the whole dataset and the expectation requires Monte-carlo sampling.

12.4.2 Markov-chain Monte-carlo method (MCMC)

Another posterior estimation is performed by MCMC. MCMC samples each weights according to a **proposal distribution** or **transition distribution** $p(\mathbf{w}_{t+1}|\mathbf{w}_t)$ and moves around the space while accepting or rejecting the proposal from the distribution. Since the next state is determined only by the current state, it is called Markov-chain. For the final sampling, we use every t -th sample from the history to avoid the correlation between samples close to each other in terms of time steps. The goal of the sampling is to identify the **stationary distribution** π that is ideally the posterior distribution $p(\mathbf{w}|\mathcal{D})$. The sufficient condition for a stationary distribution to exist is that the **detailed balance** $\pi(\mathbf{w}_t)p(\mathbf{w}_{t+1}|\mathbf{w}_t) = \pi(\mathbf{w}_{t+1})p(\mathbf{w}_t|\mathbf{w}_{t+1})$. In practice, even when the detailed balance is satisfied, it

Algorithm 2 Metropolis Hastings

$p(\mathbf{w}'|\mathbf{w})$ \triangleright Proposal distribution. This is typically Gaussian distribution with mean = \mathbf{w} .
 $\pi(\mathbf{w}) = p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w})$ \triangleright Since we take the ratio, $p(\mathcal{D})$ is canceled out.
1: **function** METROPOLIS HASTINGS
2: **for** $t = 0, 1, \dots, T$ **do**
3: $\mathbf{w}' \sim p(\mathbf{w}'|\mathbf{w}_t)$
4: $\mathbf{w}_{t+1} = \mathbf{w}'$ with the probability of $\min\left(1, \frac{\pi(\mathbf{w}')p(\mathbf{w}_t|\mathbf{w}')}{\pi(\mathbf{w}_t)p(\mathbf{w}'|\mathbf{w}_t)}\right)$ otherwise \mathbf{w}_t

may still take time to reach the stationary distribution. This time (from the beginning) is called **mixing time**. Intuitively, the distribution reaches the stationary distribution when the chain forgets the beginning states. Therefore, we take samples after a burning-in phase and the length of the burning-in phase is a hyperparameter. MCMC can sample from multi-modal distributions, but it often fails if each peak is far away from each other.

The major algorithm for MCMC is Metropolis-Hasting shown in Algorithm 2. Since it takes the ratio of the distribution, we do not have to explicitly compute the constant of proportionality $p(\mathcal{D})$. Therefore, we can transform the posterior as follows:

$$p(\mathbf{w}|\mathcal{D}) \propto p(\mathcal{D}|\mathbf{w})p(\mathbf{w}) \text{ (likelihood} \times \text{prior)} \quad (69)$$

We use the right hand side of this equation for the computation of the acceptance rate. Another example is Hamiltonian Monte Carlo (HMC) that takes into account gradients of weights. The special case is the following Langevin dynamics:

$$\mathbf{w}_{t+1} = \mathbf{w}_t + \frac{\alpha}{2} \left(\underbrace{\frac{\partial \log p(\mathbf{w}_t)}{\partial \mathbf{w}}}_{\text{prior gradient}} + \underbrace{\frac{\partial \log p(\mathcal{D}|\mathbf{w}_t)}{\partial \mathbf{w}}}_{\text{likelihood gradient}} \right) + \mathbf{r}_t \quad (70)$$

where \mathbf{r}_t are random numbers that follow Gaussian distribution and α is a learning rate. The major issue of this gradient is to take the whole dataset to compute the likelihood gradient. To ease this issue, the stochastic gradient Langevin dynamics (SGLD) has been invented. Another solution is to use mini-batches and this is SGHMC.

12.5 Alternative uncertainty treatments for neural networks

12.5.1 Output ensembling

The naïve way to measure the uncertainty of regression tasks is the following:

$$\hat{y}(\mathbf{x}) \sim \mathcal{N}(\mu(\mathbf{x}), \sigma(\mathbf{x})^2)$$

$$\text{where } \mu(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M f_i(\mathbf{x}; \mathbf{w}_i), \sigma^2(\mathbf{x}) = \frac{1}{M} \sum_{i=1}^M (f_i(\mathbf{x}; \mathbf{w}_i) - \mu)^2 \quad (71)$$

where $f_i(\cdot; \mathbf{w}_i)$ is the i -th neural network and we take only the average, but not the variance, of softmax in the case of classification tasks. Depending on how we collect the M neural networks, the performance and the computational overhead will change. The simple methods are to collect models by **randomly dropping weights (MC dropout)** of a trained model or collect M **different weights from MCMC**. The **snapshot ensemble**, that uses a cyclic learning rate scheduling and collects weights at each end of the cycles, is another option. Those three methods **require only 1 training time**.

Another effective method is to train different M models by different initial weights using SGD. It requires M training, but it achieves the state-of-the-art performance. Other than that, models with different hyperparameter configurations can take advantage of hyperparameter optimization and M different architectures can diversify the neural architectures and outperform the naïve ensembles.

12.5.2 The estimation of parametric models

The most common model is the Gaussian distribution and we estimate the parameters μ, σ . The performance metric of the training is the loglikelihood as follows:

$$\log p(\mathcal{D}|\mathbf{w}) = \frac{1}{N} \sum_{i=1}^N \log p(y_i|\mu(\mathbf{x}_i; \mathbf{w}), \sigma^2(\mathbf{x}_i; \mathbf{w})) \quad (72)$$

where $p(y|\mu, \sigma^2)$ is the probability density value of y in the Gaussian distribution with the mean μ and the variance σ^2 . We can combine the parametric model and ensembling as follows:

$$\mu = \frac{1}{M} \sum_{i=1}^M \mu_i, \quad \sigma^2 = \frac{1}{M} \sum_{i=1}^M ((\mu_i - \mu)^2 + \sigma_i^2) \quad (73)$$

Note that this variance is the total variance.