

AuE 8930 Deep Learning: Applications in Engineering

1 Optimization

Loss function

A loss function \mathcal{L} is a function that evaluate the difference between the prediction \hat{y} and the ground truth y . The common loss functions are shown below.

- Least squared error

$$\mathcal{L}(\hat{y}, y) = \frac{1}{2}(\hat{y} - y)^2$$

which is usually used in the regression problem.

- Cross-entropy

$$\mathcal{L}(\hat{y}, y) = -[y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})]$$

which is used for the classification problem.

Cost function

The cost function J is commonly used to assess the performance of a model, and is defined with the loss function \mathcal{L} as follows:

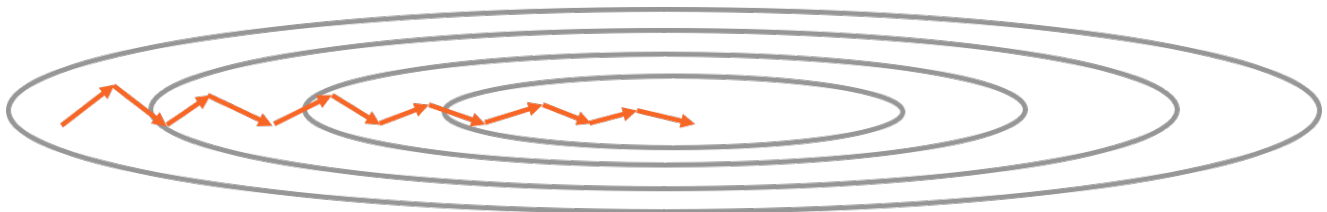
$$J(w, b) = \sum_{i=1}^m \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$

Gradient descent

By noting the learning rate $\alpha \in \mathbb{R}$, the update rule for gradient descent can be expressed as follows:

$$w := w - \alpha \cdot \frac{\partial J(w, b)}{\partial w}$$

$$b := b - \alpha \cdot \frac{\partial J(w, b)}{\partial b}$$



Mini-batch gradient descent (MBGD)

During the training phase, updating weights is usually not based on the whole training set at once due to computation complexities or one data point due to noise issues. Instead, the update step is done on mini-batches, where the batch size is a hyperparameter that we can tune.

- If batch size is set to be 1, we will get **stochastic gradient descent (SGD)**.
- If batch size equals the number of training examples, we are using **batch gradient descent (BGD)**.

Adaptive learning rate

Letting the learning rate vary when training a model can reduce the training time and improve the numerical optimal solution.

- Momentum
- RMSprop
- Adam (most commonly used technique)

Learning rate decay

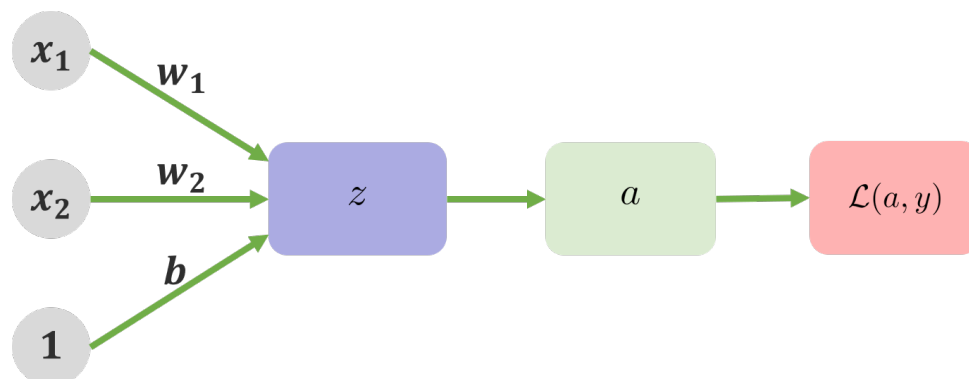
Learning rate decay is another method to speed up the training process.

- Exponential decay
- Discrete staircase

2 Logistic Regression

Logistic regression models the probabilities for classification problems with two possible outcomes. It's an extension of the linear regression model for classification problems.

Forward propagation

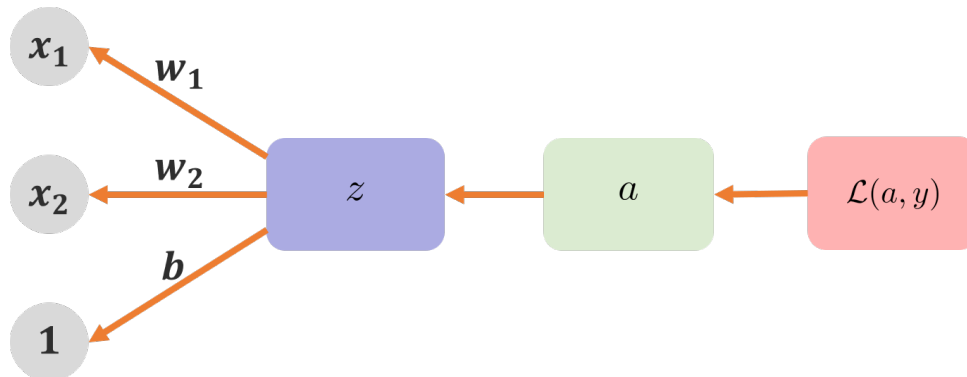


$$z = w_1 x_1 + w_2 x_2 + b = w^T x + b$$

$$\hat{y} = a = \sigma(z) = \frac{1}{1 + e^{-(w^T x + b)}}$$

$$\mathcal{L}(\hat{y}, y) = -(y \log(\hat{y}) + (1 - y) \log(1 - \hat{y})) = - \left(y \log \frac{1}{1 + e^{-(w^T x + b)}} + (1 - y) \log \left(1 - \frac{1}{1 + e^{-(w^T x + b)}} \right) \right)$$

Backward propagation



$$da = \frac{dL}{da} = \frac{-y}{a} + \frac{1-y}{1-a}$$

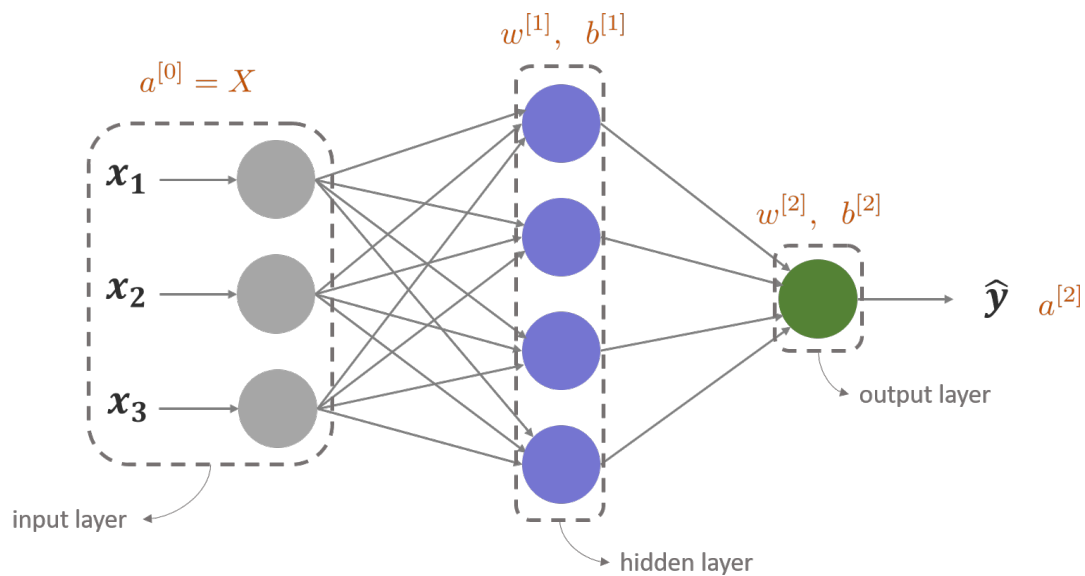
$$dz = \frac{dL}{dz} = \frac{dL}{da} \cdot \frac{da}{dz} = a - y$$

$$dw_1 = \frac{dL}{dw_1} = \frac{dL}{dz} \cdot \frac{dz}{dw_1} = (a - y)x_1$$

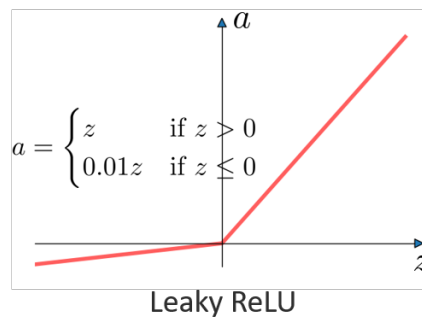
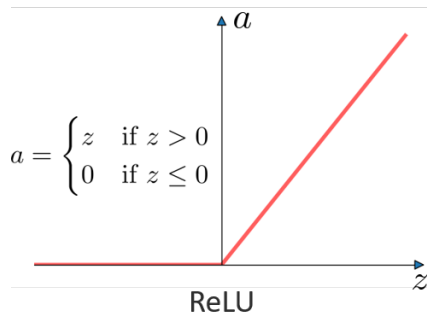
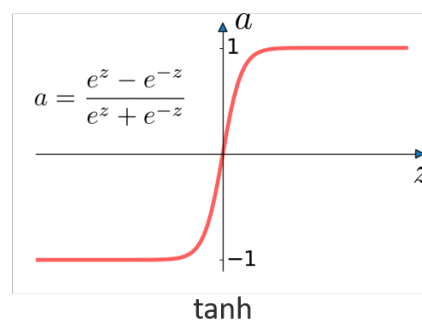
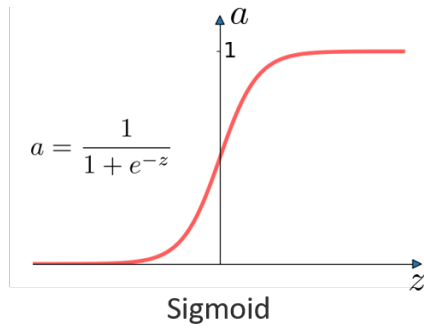
$$dw_2 = \frac{dL}{dw_2} = \frac{dL}{dz} \cdot \frac{dz}{dw_2} = (a - y)x_2$$

$$db = \frac{dL}{db} = \frac{dL}{dz} \cdot \frac{dz}{db} = a - y$$

3 Shallow neural network



Common activation functions



Weights initialization

Instead of initializing the weights in a purely random manner, we should have initial weights that take into account characteristics that are unique to the architecture.

- He method: $w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{2}{n^{(l-1)}})$
- Xavier method: $w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{1}{n^{(l-1)}})$
- Other method: $w^{[l]} = \text{np.random.randn(shape)} * \text{np.sqrt}(\frac{2}{n^{(l-1)} + n^{(l)}})$

Forward propagation

$$Z^{[1]} = W^{[1]}X + b^{[1]}$$

$$A^{[1]} = g^{[1]}(Z^{[1]})$$

$$Z^{[2]} = W^{[2]}A^{[1]} + b^{[2]}$$

$$A^{[2]} = g^{[2]}(Z^{[2]})$$

Backward propagation

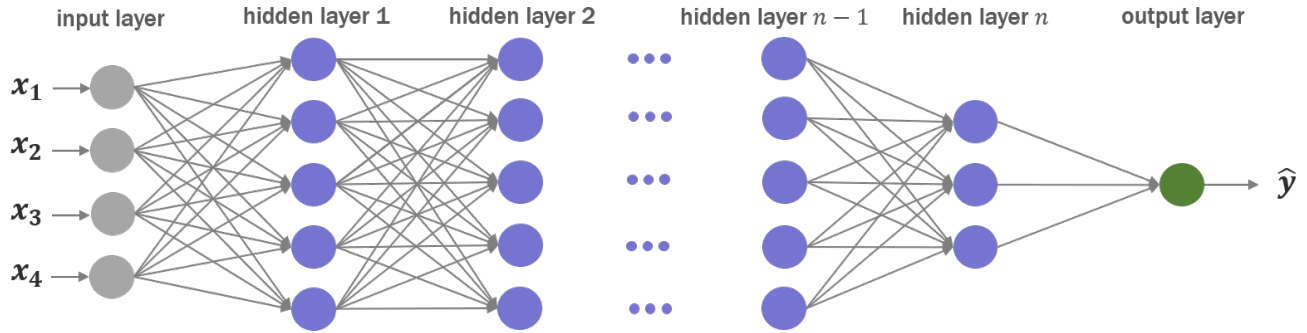
$$dZ^{[2]} = A^{[2]} - Y$$

$$dW^{[2]} = \frac{1}{m} dZ^{[2]} A^{[1]T}$$

$$db^{[2]} = \frac{1}{m} \text{np.sum}(dZ^{[2]}, \text{axis} = 1, \text{keepdims}=\text{True})$$

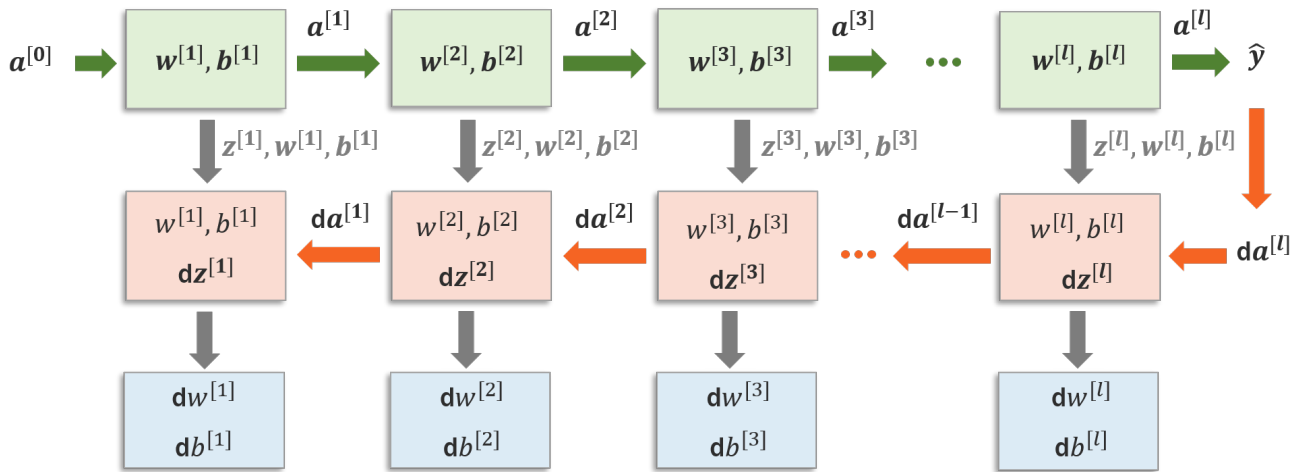
$$dZ^{[1]} = W^{[2]T} dZ^{[2]} * g^{[1]'}(Z^{[1]})$$

4 Deep neural network



$n + 1$ layer NN

Building blocks



Forward propagation

$$Z^{[l]} = W^{[l]} A^{[l-1]} + b^{[l]}$$

$$A^{[l]} = g^{[l]}(Z^{[l]})$$

Backward propagation

$$dZ^{[l]} = dA^{[l]} * g^{[l]'}(Z^{[l]})$$

$$dW^{[l]} = \frac{1}{m} dZ^{[l]} \cdot A^{[l-1]T}$$

$$db^{[l]} = \frac{1}{m} \text{np.sum}(dZ^{[l]}, \text{axis} = 1, \text{keepdims}=\text{True})$$

$$dA^{[l-1]} = W^{[l]T} dZ^{[l]}$$

Updating parameters

$$w^{[l]} := w^{[l]} - \alpha \cdot dw^{[l]}$$

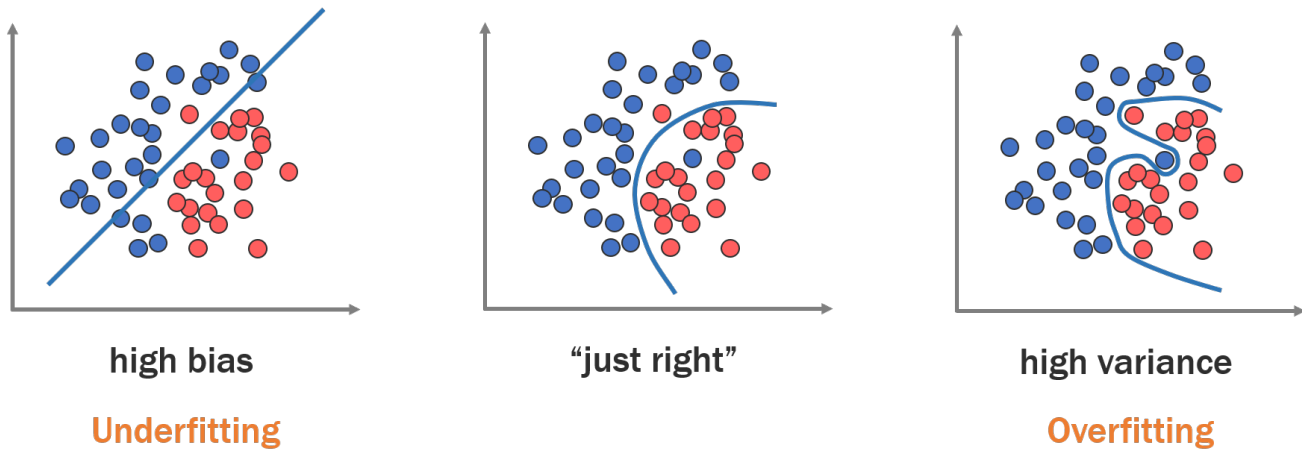
$$b^{[l]} := b^{[l]} - \alpha \cdot db^{[l]}$$

Hyperparameters

1. learning rate α
2. # of iterations of gradient descent process
3. # of hidden layers L
4. # of hidden units in each layer
5. choice of activation function

5 Tips and tricks in NN

Bias and variance



Some cases of bias and variance:

Error	High bias (underfitting)	High variance (overfitting)	High bias and high variance
training set	14%	2%	15%
dev set	15%	12%	25%

Methods to reduce high bias:

- increase model complexity
 - more # of hidden layers
 - more neurons in each layer
- increase the number of features
- training longer (more epochs)
- use better optimization techniques
- remove noise from data

Approaches to reduce high variance:

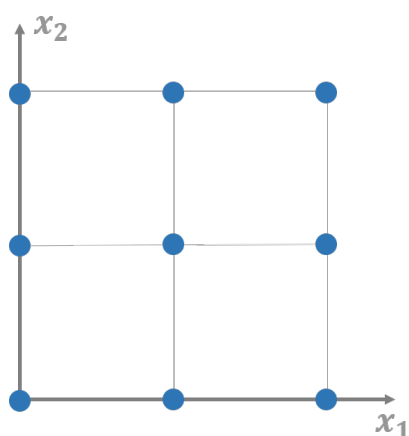
- increase the training data set
- reduce model complexity
- regularization techniques
- early stopping during training phase

Regularization

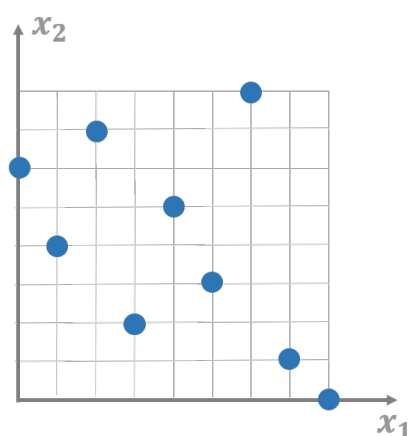
- Dropout
- Weight regularization
 - L1
 - L2
- Early stopping
- Input normalization

Hyperparameter tuning

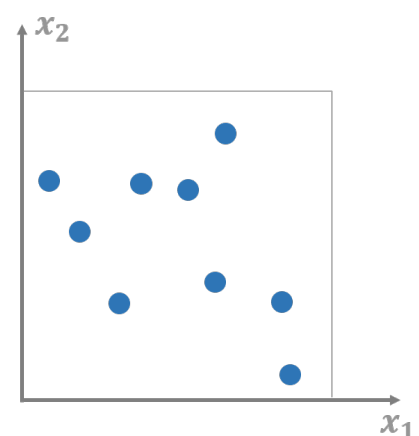
Strategies to sample the hyperparameter design space:



Grid search



Latin hypercube sampling



Random sampling

Batch normalization is a method used to make NNs faster and more stable through normalization of the input layer by re-centering and re-scaling. The algorithm is shown below.

compute $z^{(1)}, z^{(2)}, \dots, z^{(m)}$

$$\mu = \frac{1}{m} \sum_{i=1}^m z^{(i)}$$

$$\sigma^2 = \frac{1}{m} \sum_{i=1}^m \left(z^{(i)} - \mu \right)^2$$

$$z_{\text{normalization}}^{(i)} = \frac{z^{(i)} - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

$$\tilde{z}_{\text{normalization}}^{(i)} = \gamma z_{\text{normalization}}^{(i)} + \beta \quad \gamma \text{ and } \beta \text{ are learnable parameters in the algorithm}$$

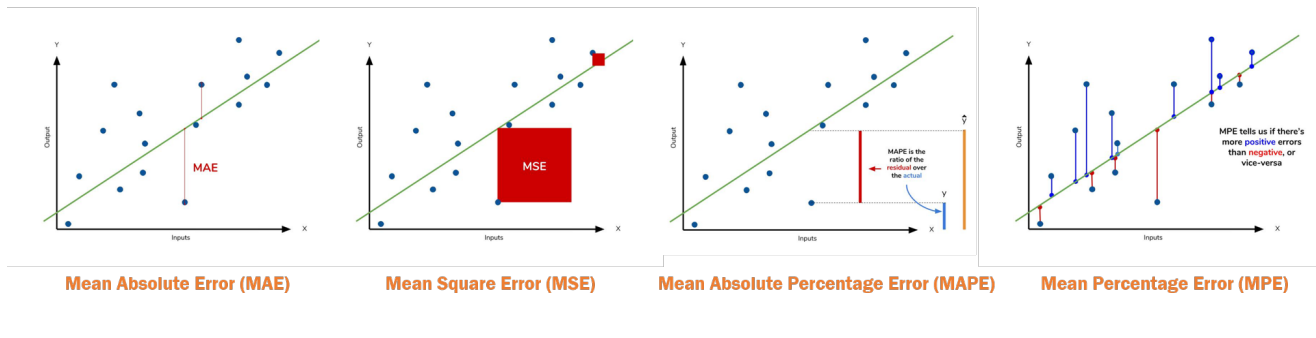
→ Avoid the mean value of $z_{\text{normalization}}^{(i)}$ close to 0

Confusion matrix is the most commonly used error measurement in classification problem.

- A row-normalized row summary displays the number of correctly and incorrectly classified observations for each true class as percentages of the number of observations of the corresponding true class.
- A column-normalized column summary displays the number of correctly and incorrectly classified observations for each predicted class as percentages of the number of observations of the corresponding predicted class.

airplane	923	4	21	8	4	1	5	5	23	6	92.3%	7.7%
automobile	5	972	2					1	5	15	97.2%	2.8%
bird	26	2	892	30	13	8	17	5	4	3	89.2%	10.8%
cat	12	4	32	826	24	48	30	12	5	7	82.6%	17.4%
deer	5	1	28	24	898	13	14	14	2	1	89.8%	10.2%
dog	7	2	28	111	18	801	13	17		3	80.1%	19.9%
frog	5		16	27	3	4	943	1	1		94.3%	5.7%
horse	9	1	14	13	22	17	3	915	2	4	91.5%	8.5%
ship	37	10	4	4		1	2	1	931	10	93.1%	6.9%
truck	20	39	3	3			2	1	9	923	92.3%	7.7%
	88.0%	93.9%	95.8%	79.0%	91.4%	89.7%	91.6%	94.1%	94.8%	95.0%		
	12.0%	6.1%	14.2%	21.0%	8.6%	10.3%	8.4%	5.9%	5.2%	5.0%		
	airplane	automobile	bird	cat	deer	dog	frog	horse	ship	truck		

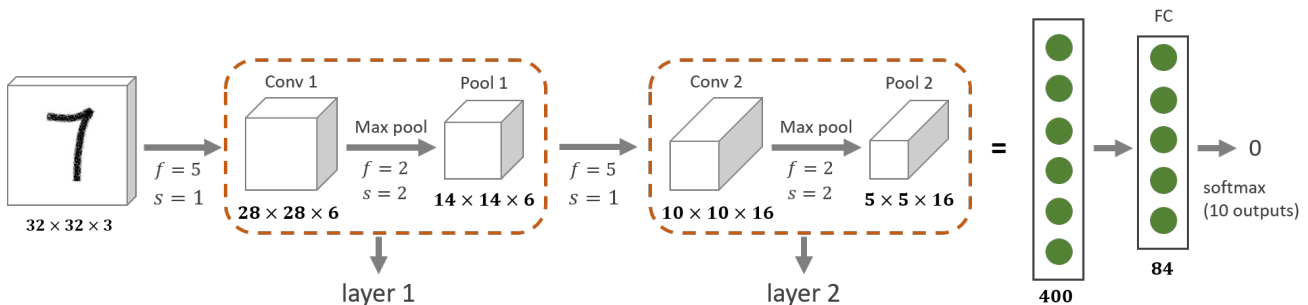
Regressor performance



6 Convolutional neural network

Architecture of CNN

Convolutional neural networks, also known as CNNs, are a specific type of neural networks that are generally composed of the following layers:

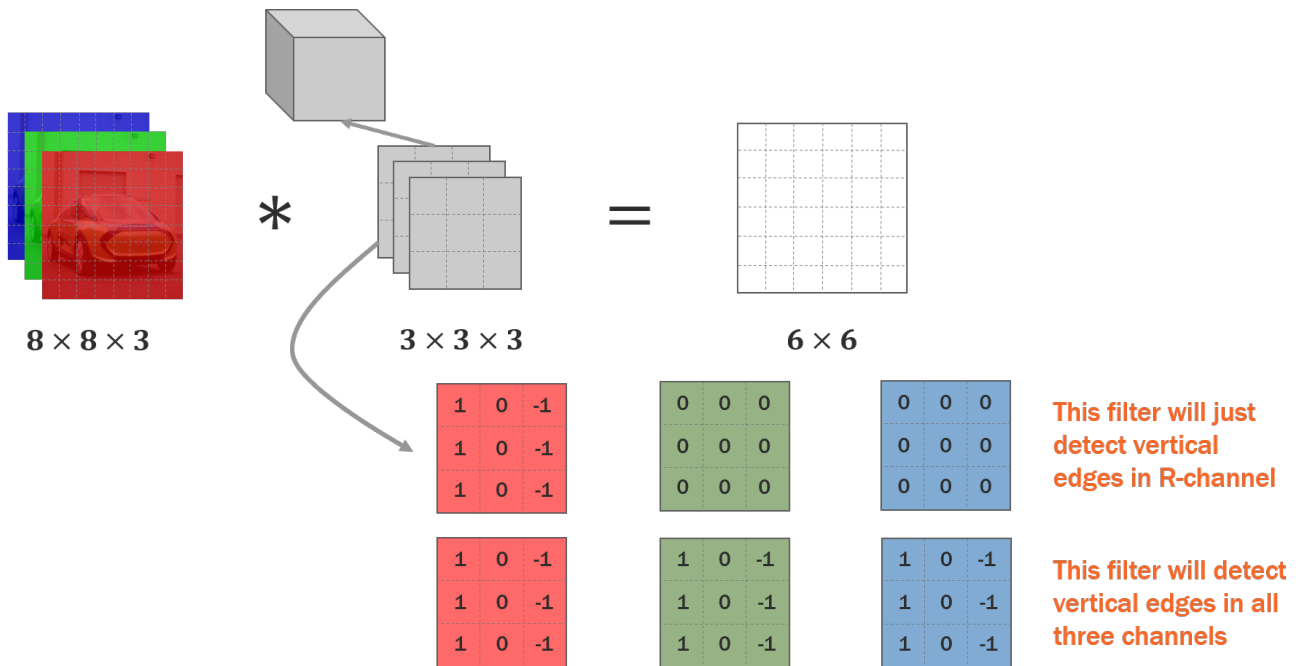


Convolution layer

The convolution layer (CONV) uses filters that perform convolution operations as it is scanning the input I with respect to its dimensions. The resulting output is called *feature map*.

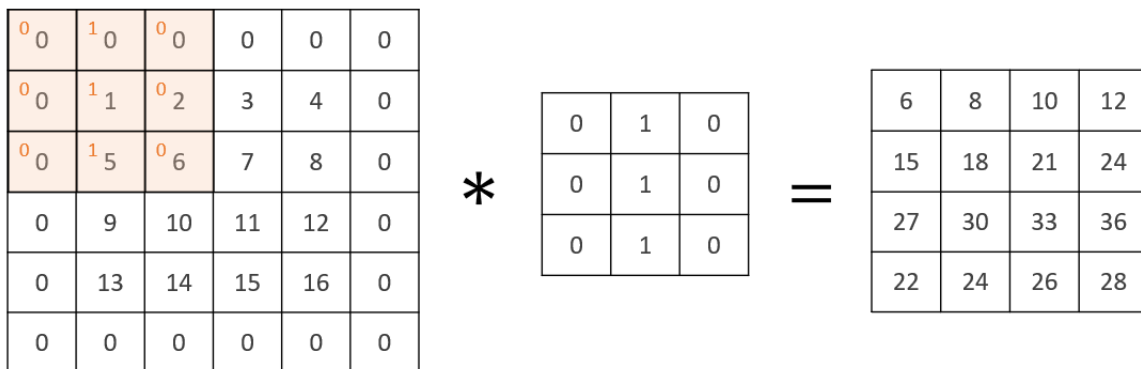
$$\begin{bmatrix} 0 & 1 & 0 & 3 \\ 0 & 5 & 1 & 6 & 0 & 7 \\ 0 & 9 & 1 & 10 & 0 & 11 \\ 13 & 14 & 15 & 16 \end{bmatrix} * \begin{bmatrix} 0 & 1 & 0 \\ 0 & 1 & 0 \\ 0 & 1 & 0 \end{bmatrix} = \begin{bmatrix} 18 & 21 \\ 30 & 33 \end{bmatrix}$$

Convolution on volume



Padding

The issue when applying convolutional layers is that we tend to lose pixels on the edges of the image. A straightforward solution to this problem is to add extra pixels of filler around the boundary of our input image.

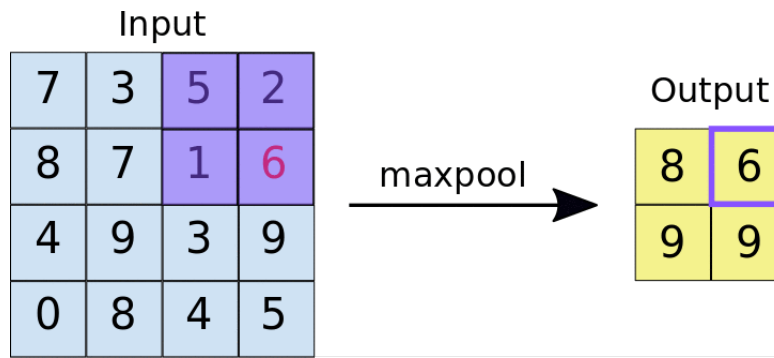


Computing output shape

$$\left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor \times \left\lfloor \frac{n + 2p - f}{s} + 1 \right\rfloor$$

Max pooling

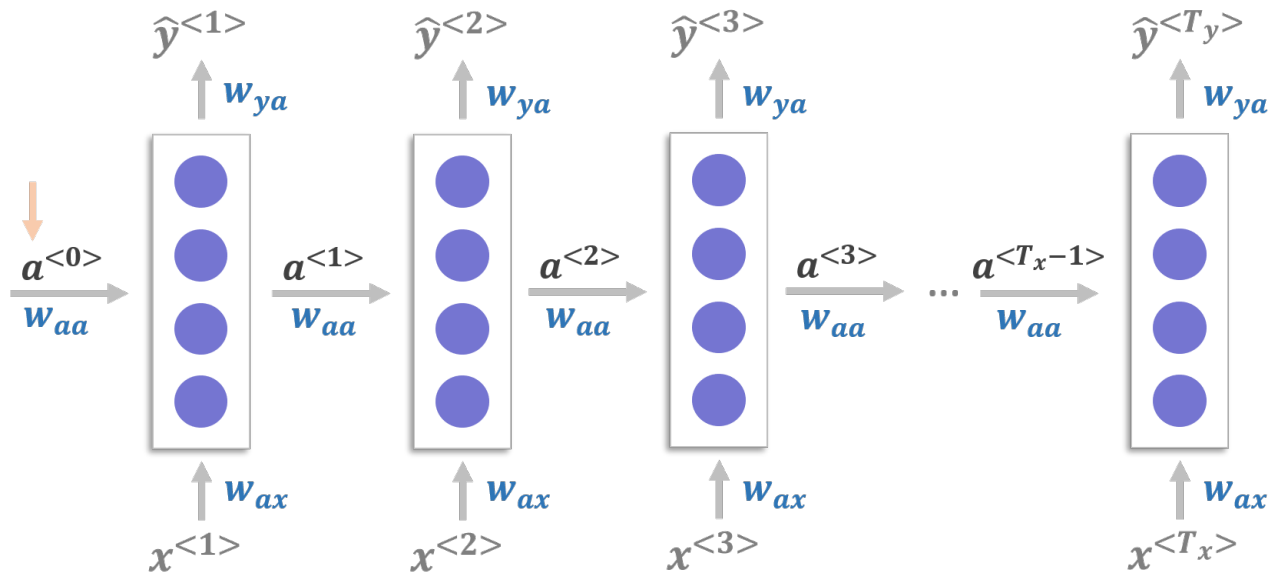
The max pooling layer is a downsampling operation, typically applied after a convolution layer, which does some spatial invariance.



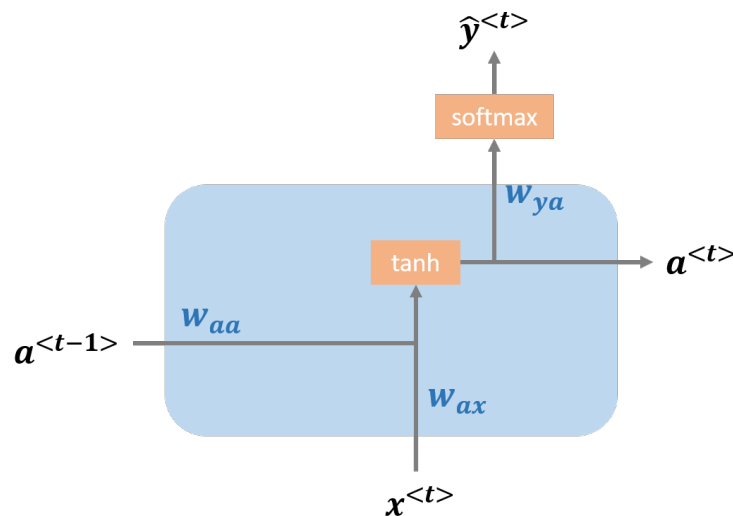
7 Recurrent neural network

Architecture of RNN

Recurrent neural networks, also known as RNNs, are a class of neural networks that allow previous outputs to be used as inputs while having hidden states.



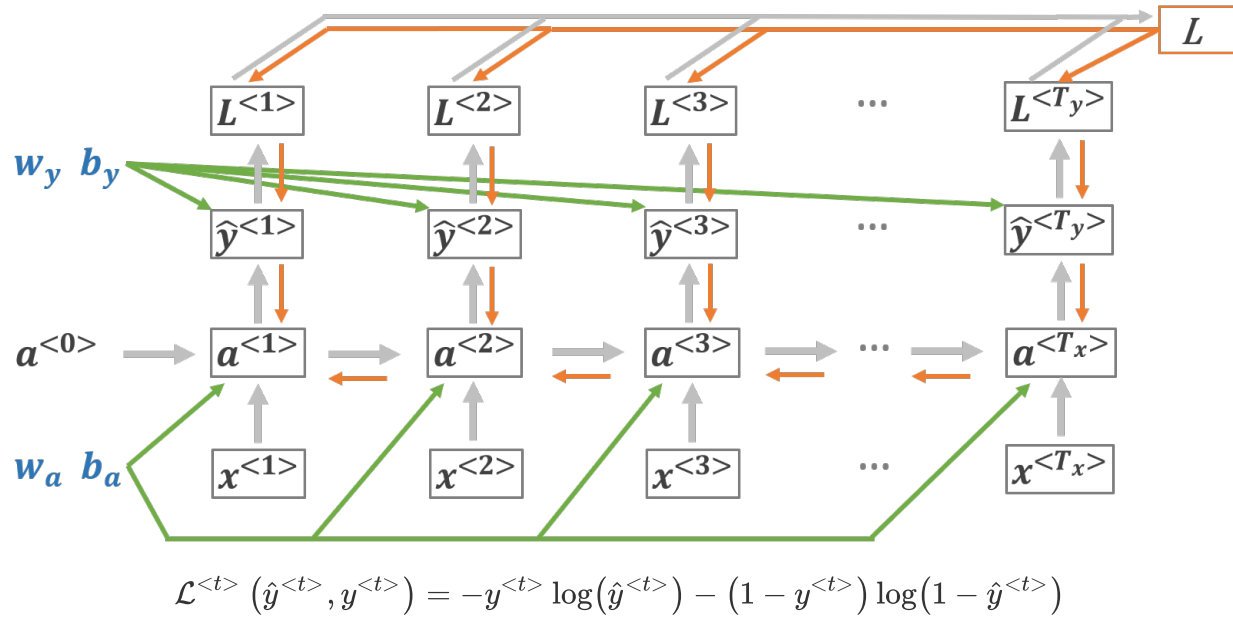
Forward propagation



$$a^{<t>} = g(w_{aa} \cdot a^{<t-1>} + w_{ax} \cdot x^{<t>} + b_a)$$

$$\hat{y}^{<t>} = g(w_{ya} \cdot a^{<t>} + b_y)$$

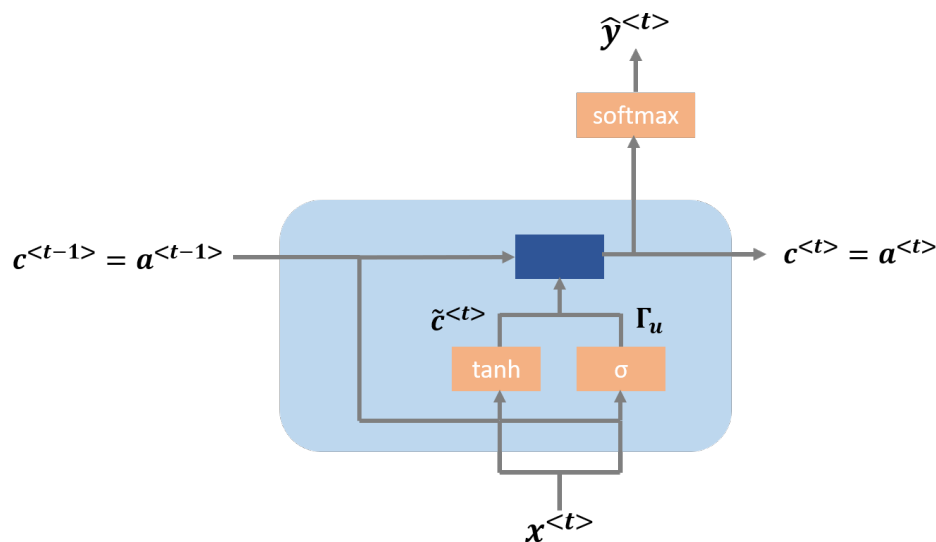
Backward propagation



$$\mathcal{L}(\hat{y}, y) = \sum_{t=1}^{T_x} \mathcal{L}^{<t>}(\hat{y}^{<t>}, y^{<t>})$$

Gated recurrent unit (GRU)

The key distinction between vanilla RNNs and GRUs is that the latter support gating of the hidden state. This means that we have dedicated mechanisms for when a hidden state should be *updated* and also when it should be *reset*.



$$\tilde{c}^{<t>} = \tanh(w_c [c^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u [c^{<t-1>}, x^{<t>}] + b_u)$$

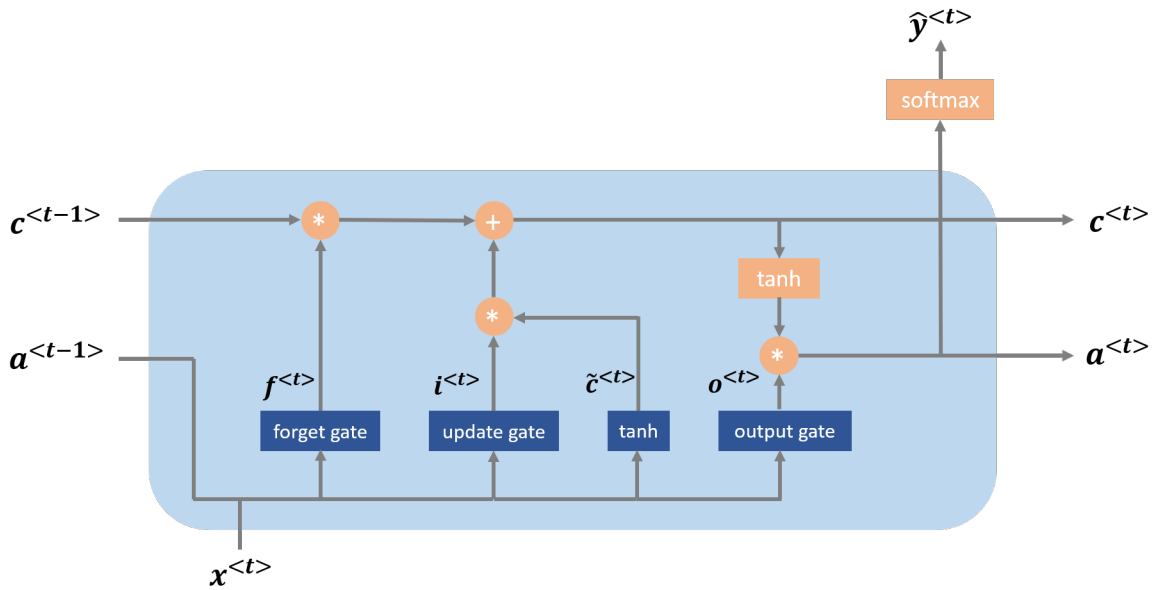
$$\Gamma_r = \sigma(W_r [c^{<t-1>}, x^{<t>}] + b_r)$$

$$a^{<t>} = c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + (1 - \Gamma_u) * c^{<t-1>}$$

Long short-term memory (LSTM)

LSTM introduces a *memory cell* (or *cell* for short) that has the same shape as the hidden state (some literatures consider the memory cell as a special type of the hidden state), engineered to record additional information.

- **output gate**: read out the entries from the cell.
- **update gate** (input gate): decide when to read data into the cell.
- **forget gate**: reset the content of the cell.



$$\tilde{c}^{<t>} = \tanh(W_c [a^{<t-1>}, x^{<t>}] + b_c)$$

$$\Gamma_u = \sigma(W_u [a^{<t-1>}, x^{<t>}] + b_u)$$

$$\Gamma_f = \sigma(W_f [a^{<t-1>}, x^{<t>}] + b_f)$$

$$\Gamma_o = \sigma(W_o [a^{<t-1>}, x^{<t>}] + b_o)$$

$$c^{<t>} = \Gamma_u * \tilde{c}^{<t>} + \Gamma_f * c^{<t-1>}$$

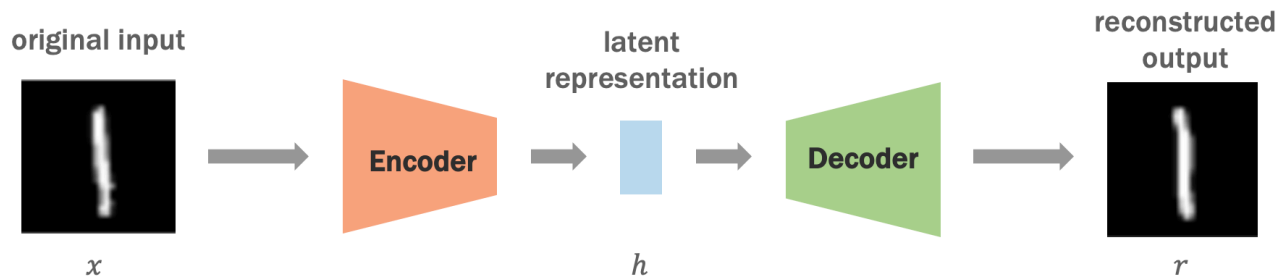
$$a^{<t>} = \Gamma_o * c^{<t>}$$

8 Autoencoders and Variational Autoencoders

Unsupervised learning

- clustering
- data compression
- representation learning

Autoencoders

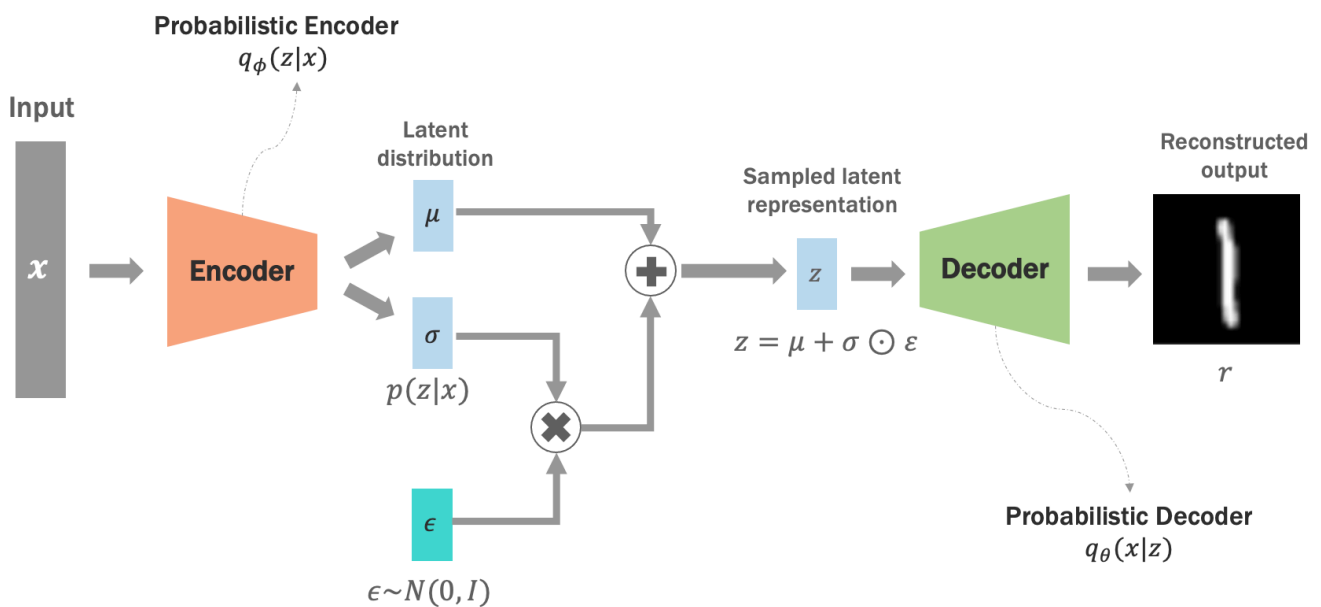


- Encoder: Compact/compress input features into a latent-space features of usually smaller dimension.
- Decoder: Reconstruct original input from the latent space.

Common applications of autoencoder:

- denoising
- image colorization
- image compression
- image search and retrieval

Variational Autoencoders (Optional)



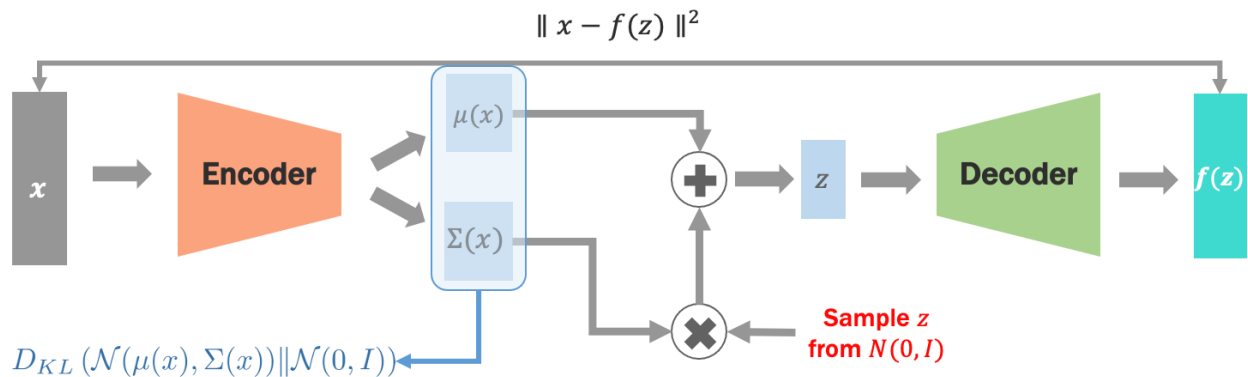
- KL divergence

$$D_{KL}(Q(z|x)||P(z|x)) = - \sum_z Q(z|x) \log \frac{P(x,z)}{Q(z|x)} + \log P(x)$$

- ELBO loss function

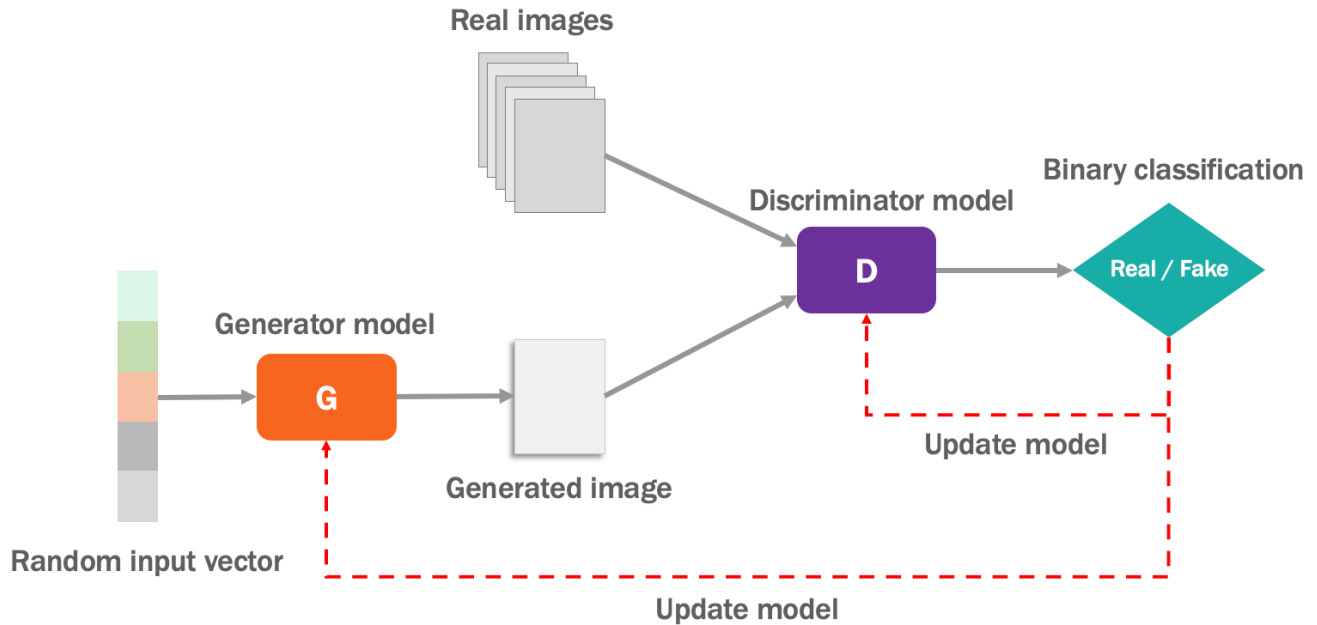
$$\underbrace{\mathbb{E}_{Q(z|x)} \log P(x|z)}_{\text{Maximize}} - \underbrace{D_{KL}(Q(z|x)||P(z))}_{\text{Minimize}}$$

- reparameterization

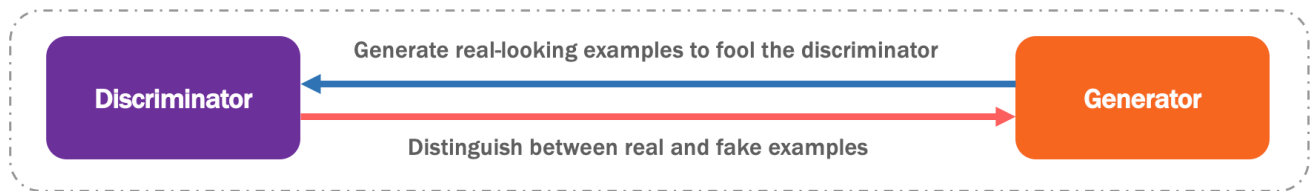


9 Generative Adversarial Networks (GANs)

Architecture of GAN



Main components



Mini-max Game

- Generator and discriminator have “opposite” targets
- We’ll start with the discriminator which is a classifier
- The discriminator aims to distinguish between real and fake images

Cost function

GANs can be considered as a zero sum game or minmax problem with the cost function:

$$\theta_G^* = \arg \min_{\theta_G} \max_{\theta_D} -J_D$$

- discriminator cost function

$$J_D = - \left(\sum_x \log D(x) + \sum_z \log(1 - D(G(z))) \right)$$

- generator cost function

$$J_G = \sum_z \log(1 - D(G(z)))$$