

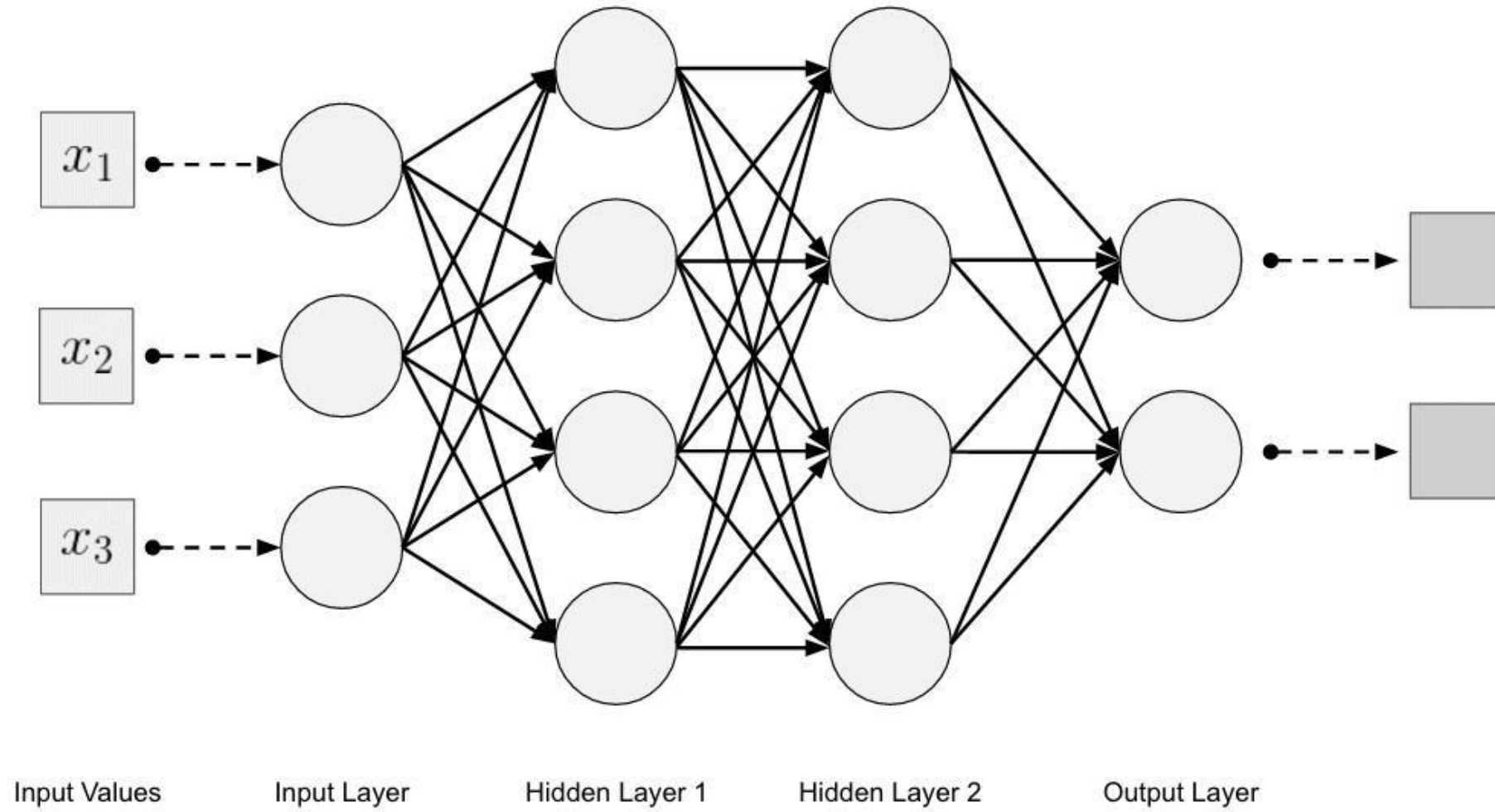
Module 2: Fundamentals of Deep Networks

Neural Networks

The behavior of neural networks is shaped by its network architecture

- Number of neurons
- Number of layers
- Types of connections between layers

Multilayer Neural Networks Topology



Multilayer neural network topology

Multilayer Neural Networks Topology

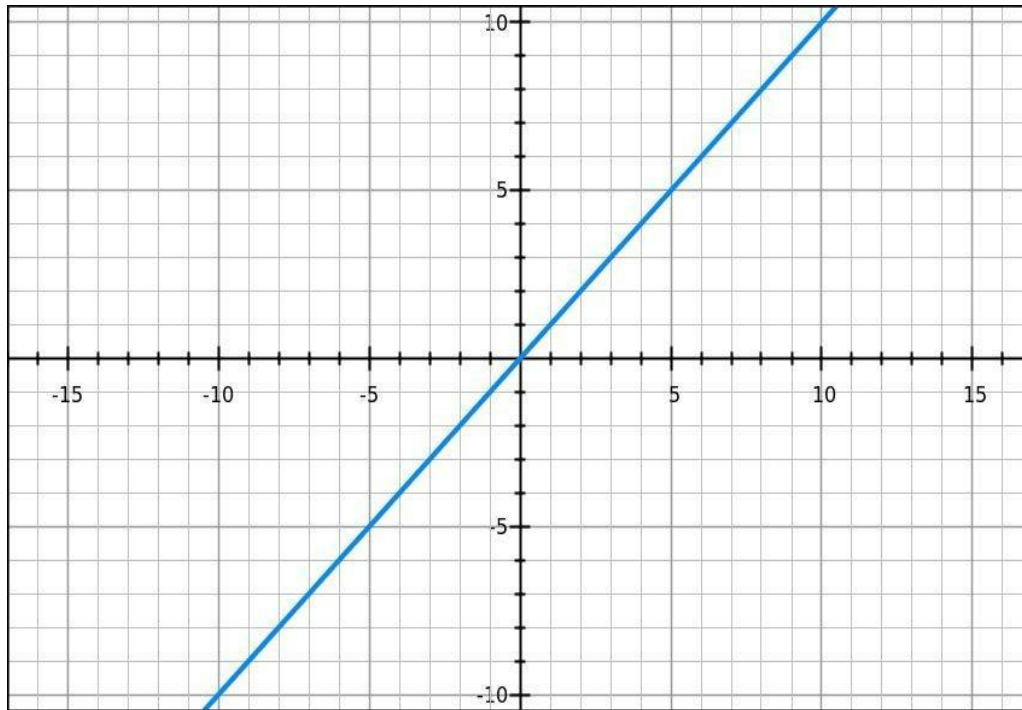
- *Connection weights*-coefficients that scale (amplify or minimize) the input signal to a given neuron in the network
- *Biases*-
 - Scalar values added to the input to ensure that at least few nodes per layer are activated regardless of signal strength.
 - Biases allow learning to happen by giving the network action in the event of low signal.
 - Allow the network to try new interpretations or behaviors.
 - Notated as b , and, like weights, biases are modified throughout the learning process
- *Activation functions*-The functions that govern the artificial neuron's behavior

Activation functions-

Activation functions-

- Propagate the output of one layer's nodes forward to the next layer (up to and including the output layer).
- a scalar-to-scalar function, yielding the neuron's activation.
- Used for hidden neurons in a neural network to introduce nonlinearity into the network's modeling capabilities.
- Many activation functions belong to a logistic class of transforms that (when graphed) resemble an S.
- This class of function is called *sigmoidal*.
 - The sigmoid family of functions contains several variations, one of which is known as the Sigmoid function.

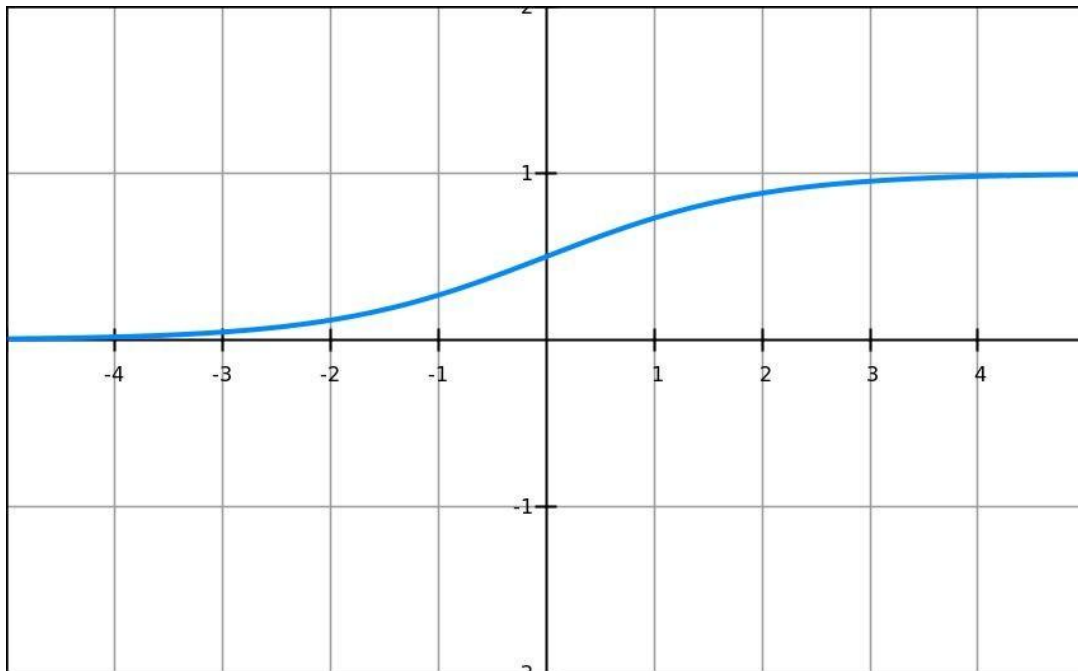
Activation functions Linear



- A linear transform is basically the identity function,
- $f(x) = W * x$
- Dependent variable has a direct, proportional relationship with the independent variable.
- It means the function passes the signal through unchanged.
- Used in the input layer of neural networks

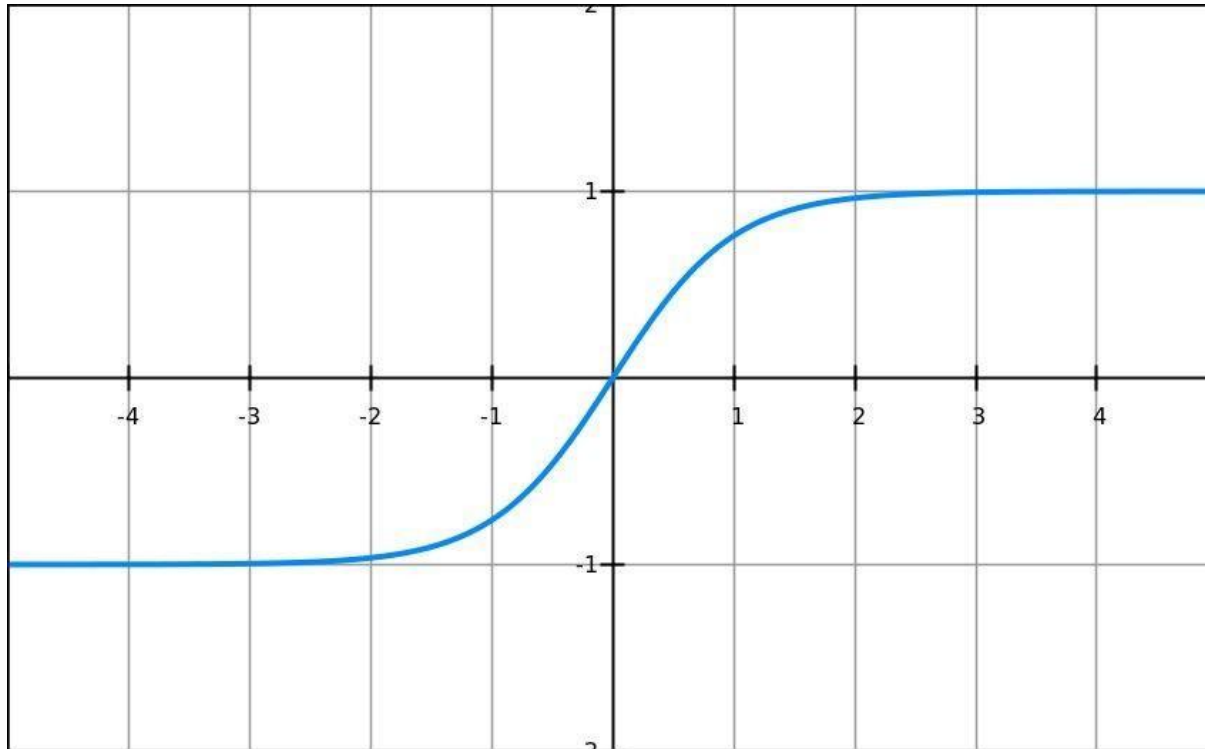
Linear activation function

Activation functions Sigmoid



- Sigmoids can reduce extreme values or outliers in data without removing them.
- A sigmoid activation function outputs an independent probability for each class.
- A sigmoid function is a machine that converts independent variables of near infinite range into simple probabilities between 0 and 1, and most of its output will be very close to 0 or 1.

Activation functions Tanh



- Hyperbolic trigonometric function
- Just as the tangent represents a ratio between the opposite and adjacent sides of a right triangle, tanh represents the ratio of the hyperbolic sine to the hyperbolic cosine:
 - $\tanh(x) = \sinh(x) / \cosh(x)$.
- normalized range of tanh is -1 to 1 .
- Advantage -it can deal more easily with negative numbers

Activation functions Hard tanh

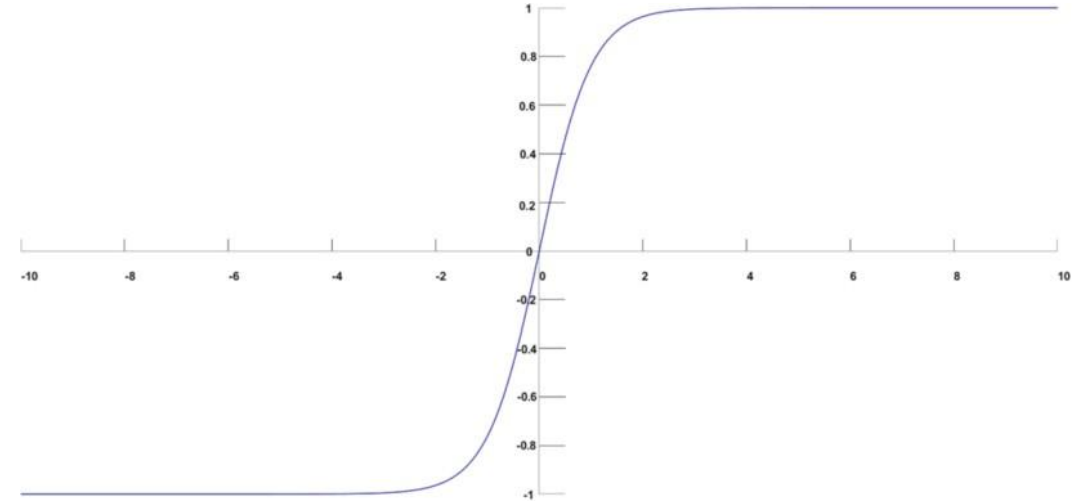
Hard Tanh-

- Similar to tanh, hard tanh simply applies hard caps to the normalized range.
- Anything more than 1 is made into 1, and anything less than -1 is made into -1 .
- This allows for a more robust activation function that allows for a limited decision boundary.

Activation functions Softmax

Softmax-

- Also known as -softargmax or normalized exponential function
- Allows us to express our inputs as a discrete probability distribution
- a generalization of logistic regression in as much as it can be applied to continuous data (rather than classifying binary) and can contain multiple decision boundaries.
- It handles multinomial labeling systems.
- Generally used at the output layer of a classifier
- Defined as follows:
 - for each value in input vector, the Softmax value is the exponent of the individual input divided by a sum of the exponents of all the inputs.
- Returns the probability distribution over mutually exclusive output classes.



$$\sigma(\vec{z})_i = \frac{e^{z_i}}{\sum_{j=1}^K e^{z_j}}$$

σ = softmax

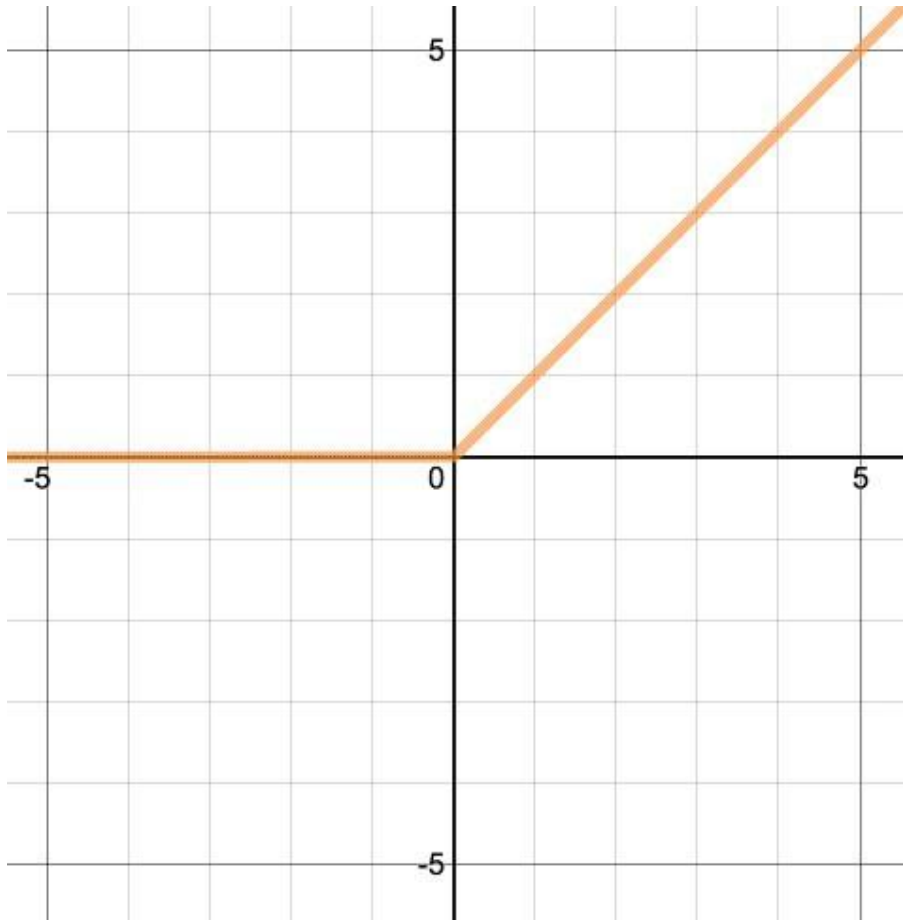
\vec{z} = input vector

e^{z_i} = standard exponential function for input vector

K = number of classes in the multi-class classifier

e^{z_j} = standard exponential function for output vector

Activation functions Rectified Linear

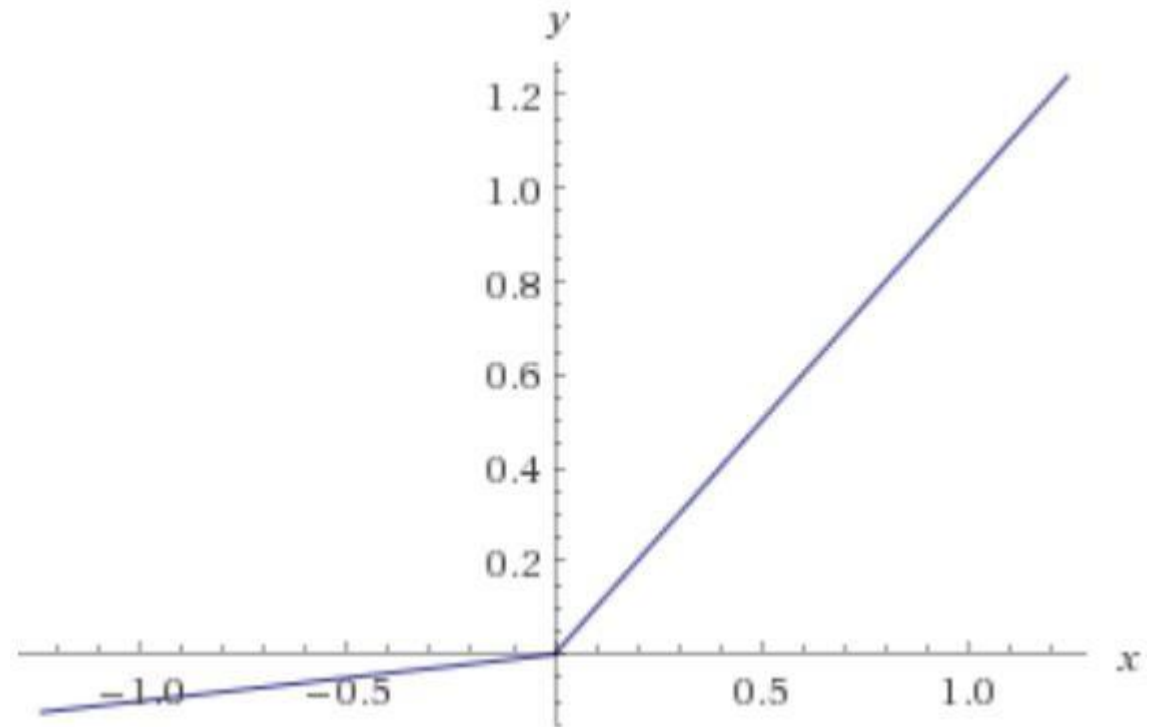


- More interesting transform that activates a node only if the input is above a certain quantity.
- While the input is below zero, the output is zero, but when the input rises above a certain threshold, it has a linear relationship with the dependent variable $f(x) = \max(0, x)$
- ReLU activation functions have shown to train better in practice than sigmoid activation functions
- Compared to the sigmoid and tanh activation functions, the ReLU activation function does not suffer from vanishing gradient issues

Activation functions Leaky Relu

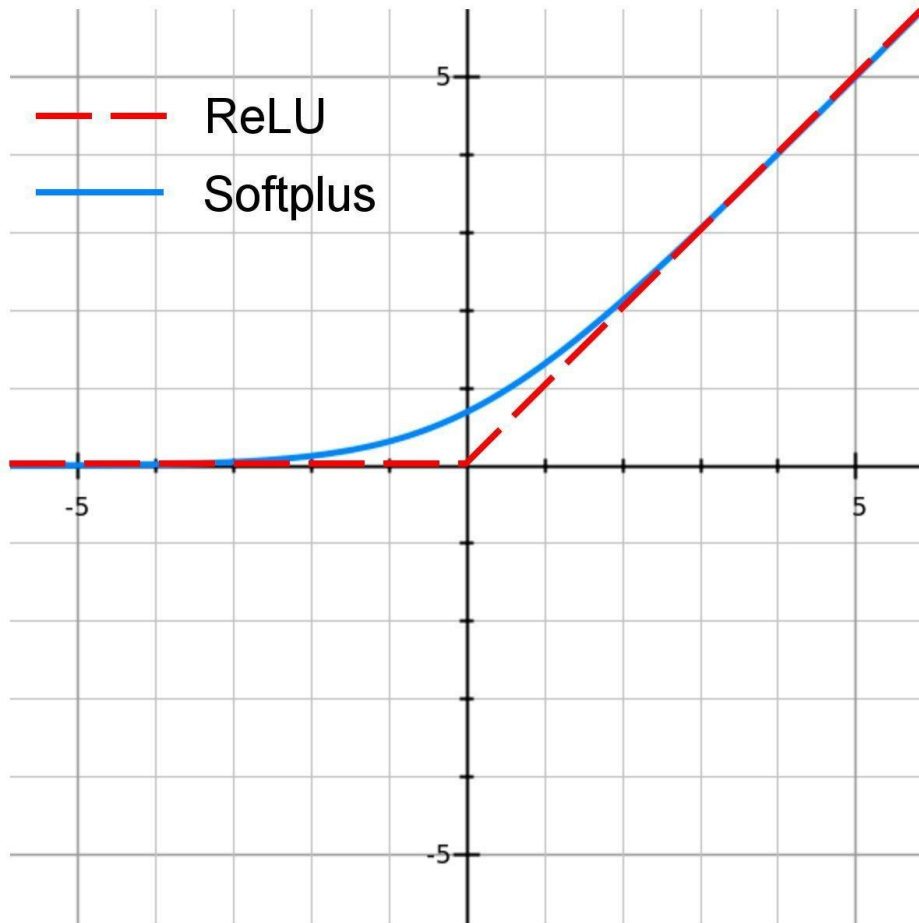
- A strategy to mitigate the “dying ReLU”
- As opposed to having the function being zero when $x < 0$, the leaky ReLU will instead have a small negative slope (e.g., “around 0.01”).
- Some success has been seen in practice with this ReLU variation, but results are not always consistent
- The equation is :

$$f(x) = \begin{cases} x & \text{if } x > 0 \\ 0.01x & \text{otherwise} \end{cases}$$



Activation functions Softplus

- This activation function is the “smooth version of the ReLU,”
- Compare this plot to the ReLU - shows that the softplus activation function ($f(x) = \ln[1 + \exp(x)]$) has a similar shape to the ReLU.



Loss Functions

- Loss Function Notation-
 - “N” - number of samples (set of inputs with corresponding outcomes) that have been gathered.
 - “P” - number of input features gathered
 - “M” - the number of output features that have been observed.
 - (X,Y) to denote the input and output data collected
 - N such pairs where the input is a collection of P values and the output Y is a collection of M values.
 - Denote the i^{th} pair in the dataset as X_i and Y_i .
- \hat{Y} we will use to denote the output of the neural net.
- \hat{Y} - network's guess at Y and therefore it will also have M features.

Loss Functions

- Use the notation $h(X_i) = \hat{Y}_i$ to denote the neural network transforming the input to give the output.
- When referring to j th output feature, use it as a subscript firmly linking the notation to a matrix where the rows are different data points and the columns are the different unique features. Thus $y_{i,j}$ refers to the j th feature observed in the i th sample collected.
- Represent the loss function by $L(W,b)$.

$$h_{w,b}(X) = \hat{Y}$$

Loss Functions

- Loss Functions for Regression-
 - Mean squared error loss

$$L(W, b) = \frac{1}{N} \sum_{i=1}^N (\hat{Y}_i - Y_i)^2$$

$$L(W, b) = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^M (\hat{y}_{ij} - y_{ij})^2$$

$$L(W, b) = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^M (\hat{y}_{ij} - y_{ij})^2$$

Loss Functions

- Loss Functions for Regression-

- Mean absolute error loss

$$L(W, b) = \frac{1}{2N} \sum_{i=1}^N \sum_{j=1}^M | \hat{y}_{ij} - y_{ij} |$$

- Mean squared log error loss

$$L(W, b) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M (\log \hat{y}_{ij} - \log y_{ij})^2$$

- Mean absolute percentage error loss

$$L(W, b) = \frac{1}{N} \sum_{i=1}^N \sum_{j=1}^M \frac{100 \times | \hat{y}_{ij} - y_{ij} |}{y_{ij}}$$

Loss Functions

- Loss Functions for Classification-

- **Hinge loss**- most commonly used loss function when the network must be optimized for a hard classification
- Also seen in a class of models called maximum-margin classification models (e.g., support vector machines)

- equation for hinge loss when data points must be categorized as -1 or 1

$$L(W, b) = \frac{1}{N} \sum_{i=1}^N \max(0, 1 - y_{ij} \times \hat{y}_{ij})$$

- Mostly used for binary classifications

Loss Functions

- Loss Functions for Classification-
 - **Logistic loss**- used when probabilities are of greater interest than hard classifications.

$$P(y_i = 1 \mid X_i; \mathbf{W}, \mathbf{b}) = h_{\mathbf{W}, \mathbf{b}}(X_i)$$

$$P(y_i = 0 \mid X_i; \mathbf{W}, \mathbf{b}) = 1 - h_{\mathbf{W}, \mathbf{b}}(X_i)$$

$$P(y_i \mid X_i; \mathbf{W}, \mathbf{b}) = (h_{\mathbf{W}, \mathbf{b}}(X_i))^{y_i} \times (1 - h_{\mathbf{W}, \mathbf{b}}(X_i))^{1-y_i}$$

$$L(W, b) = \prod_{i=1}^N \hat{y}_i^{y_i} \times (1 - \hat{y}_i)^{1-y_i}$$

Loss Functions

- Loss Functions for Classification-
 - Negative log likelihood

$$L(W, b) = - \sum_{i=1}^N y_i \times \log \hat{y}_i + (1 - y_i) \times \log (1 - \hat{y}_i)$$

- Extending the loss function from two classes to M classes gives us the following equation for logistic loss

$$L(W, b) = - \sum_{i=1}^N \sum_{j=1}^M y_{i,j} \times \log \hat{y}_{i,j}$$

Loss Functions

- Loss Functions for Reconstruction-
- equation for KL divergence

$$D_{KL}(Y \parallel \hat{Y}) = - \sum_{i=1}^N Y_i \times \log \left(\frac{Y_i}{\hat{Y}_i} \right)$$

What is Deep Learning?

- More neurons than previous networks
- More complex ways of connecting layers
- Automatic feature extraction

Architectures of deep networks

- Four major architectures of deep networks-
 - Unsupervised Pretrained Networks
 - Convolutional Neural Networks
 - Recurrent Neural Networks
 - Recursive Neural Networks

Unsupervised Pretrained Networks (UPN):

- **Purpose:** UPNs are designed for unsupervised learning, where the model learns patterns and representations from unlabeled data.
- **Training:** Initially trained on a large dataset without labeled outputs, the model captures inherent structures and features.
- **Transfer Learning:** Pretrained UPNs can be fine-tuned for specific tasks with smaller labeled datasets, leveraging the learned representations.

Convolutional Neural Networks (CNNs):

- **Purpose:** Primarily used for visual data, such as images and videos.
- **Architecture:** Consists of convolutional layers for spatial hierarchies, pooling layers for down sampling, and fully connected layers for classification.
- **Feature Extraction:** CNNs automatically learn hierarchical features, recognizing patterns like edges, textures, and shapes.

Recurrent Neural Networks (RNNs):

- **Purpose:** Suited for sequential data, like time-series, speech, and natural language.
- **Architecture:** Incorporates recurrent connections to maintain a memory of previous inputs, allowing the model to capture temporal dependencies.
- **Applications:** Used in tasks like language modeling, speech recognition, and sentiment analysis.

Recursive Neural Networks (RecNNs):

- Purpose:** Applies to hierarchical structures, such as tree-structured data.
- Architecture:** Operates recursively on a hierarchical structure, capturing dependencies at different levels of abstraction.
- Applications:** Commonly used in natural language processing tasks that involve parse trees, syntactic structures, or hierarchical relationships.

Advances in Network Architecture

- Advances in layer types
- Advances in neuron types
- Hybrid architectures

From Feature Engineering to Automated feature learning

- Feature engineering-handcrafted features □ highly accurate models but take a lot of time and experience to produce.
- Feature learning-
 - Feature maps
 - CNN

Feature Engineering:

- **Traditional Approach:** In classical machine learning, engineers manually craft features from raw data. These features are selected or engineered based on domain knowledge and insights.
- **Purpose:** Effective feature engineering enhances model performance, providing relevant information to the algorithm.

Automated Feature Learning:

- **Modern Paradigm:** With the advent of deep learning and neural networks, there is a move towards automated feature learning.
- **Process:** Instead of relying on manually designed features, deep learning models can automatically learn relevant features from raw data during training.
- **Advantages:**
 - **Data Representation:** Neural networks create hierarchical representations of data, discovering intricate patterns and relationships.
 - **Reduced Human Intervention:** The need for manual feature engineering is reduced, allowing models to adapt to diverse and complex datasets.

1.Key Technologies:

- 1. Deep Learning:** Neural networks with multiple layers can automatically extract features from raw data, capturing both simple and complex patterns.
- 2. Convolutional Neural Networks (CNNs), Recurrent Neural Networks (RNNs), and Transformers:** Architectures designed for specific data types automate the learning of features, particularly in image, sequence, and language data.

2.Benefits:

- 1. Flexibility:** Automated feature learning allows models to adapt to varying datasets without extensive manual adjustments.
- 2. Performance:** Deep learning models, by automatically learning features, can often outperform traditional models on complex tasks.

Challenges:

- **Interpretability:** Automated feature learning can make models more complex, potentially making it challenging to interpret why certain decisions are made.
- **Data Requirements:** Deep learning models typically require larger amounts of data for effective feature learning.

Common Architectural Principles of Deep Networks

Core components-

- Parameters
- Layers
- Activation functions
- Loss functions
- Optimization methods
- Hyperparameters

Common Architectural Principles of Deep Networks

- Core components-
- Layers-
 - Fundamental architectural unit
 - Combinations of layers /activation function functions
 - Layer specific hyperparameters

Common Architectural Principles of Deep Networks

Core components-

■ Activation Functions-

- two main areas across all architectures:
 - Hidden layers
 - Output layers
- Hidden layer activation functions-
 - Sigmoid
 - Tanh
 - Hard tanh
 - Rectified linear unit (ReLU) (and its variants)
- A more continuous distribution of input data is generally best modeled with a ReLU activation function.
- Optionally use tanh activation function (if the network isn't very deep) in the event that the ReLU did not achieve good results

Common Architectural Principles of Deep Networks

Core components-

Activation Functions-

- **Output layer for regression-**
 - what type of answer we expect our model to output.
 - For e.g. a single real-valued number as o/p from our model, then use a linear activation function.
- **Output layer for binary classification-**
 - Use a sigmoid output layer with a single neuron to give us a real value in the range of 0.0 to 1.0 (excluding those values) for the single class.
 - This real-valued output is typically interpreted as a probability distribution.
- **Output layer for multiclass classification-**
 - If we only care about the best score across these classes, use a softmax output layer with an `arg-max()` function to get the highest score of all the classes.
 - The softmax output layer gives us a probability distribution over all the classes.

Common Architectural Principles of Deep Networks

Core components-

Loss functions-

- Regression
- Classification
- Reconstruction-
 - Involved in unsupervised feature extraction
 - In certain architectures of deep networks reconstruction loss functions help the network extract features more effectively when paired with the appropriate activation function.
 - Example - use the multiclass cross-entropy as a loss function in a layer with a softmax activation function for classification output

Common Architectural Principles of Deep Networks

Core components-

■ Optimization Algorithms-

■ First-order

- calculate the Jacobian matrix.
- The Jacobian has one partial derivative per parameter
- The algorithm then takes one step in the direction specified by the Jacobian.
- Calculate a gradient (Jacobian) at each step to determine which direction to go in next.
- At each iteration, or step, algorithm try to find the next best possible direction to go, as defined by objective function.
- For e.g- GD

$$\mathbf{H}_f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \cdots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix},$$

Common Architectural Principles of Deep Networks

Core components-

- **Optimization Algorithms-**

- Second-order-

- Calculate or approximate the Hessian
 - Calculate the derivative of the Jacobian (i.e., the derivative of a matrix of derivatives) by approximating the Hessian.
 - Take into account interdependencies between parameters when to choosing how much modify each parameter

- Methods-

- Limited-memory BFGS (L-BFGS) Broyden-Fletcher-Goldfarb-Shanno
 - Conjugate gradient
 - Hessian-free

Common Architectural Principles of Deep Networks

Core components-

- Hyperparameters-any configuration setting free to be chosen by the user that might affect performance.
- Categories-
 - Layer size
 - Magnitude (momentum, learning rate)
 - Regularization (dropout, drop connect, L1, L2)
 - Activations (and activation function families)
 - Weight initialization strategy
 - Loss functions
 - Settings for epochs during training (mini-batch size)
 - Normalization scheme for input data (vectorization)

Common Architectural Principles of Deep Networks

Hypermeters-

1. Layer size

- Number of neurons in a given layer.
- For the input layer- match up to the number of features in the input vector.
- For the output layer- either be a single output neuron or a number of neurons matching the number of classes we are trying to predict.
- Challenge in hyperparameter tuning- Deciding on neuron counts for each hidden layer
- Complexity of a problem -directly correlated to how many neurons to be used in the hidden layers of networks
- More parameters -increase in the amount of effort needed to train the network.
- Large parameter counts can lead to long training times and models that struggle to find convergence.

Common Architectural Principles of Deep Networks

Hypermeters-

2. Magnitude hyperparameters-

- Learning rate-
 - One of the key hyperparameters in neural networks.
 - How fast we change the parameter vector as we move through search space.
 - Too high--move toward goal faster (but we might also take a step so large that we shoot right past the best answer to the problem)
 - Too small-longer time for training process to complete, make our learning algorithm inefficient.
 - Specific to the dataset and to other hyperparameters.
 - Schedule learning rates to decrease over time according to some rule

Common Architectural Principles of Deep Networks

Hypermeters-

2. Magnitude hyperparameters-

- Sparsity –
 - Recognizes that for some inputs only a few features are relevant
- Momentum-
 - Momentum helps the learning algorithm get out of spots in the search space where it would otherwise become stuck.
 - Momentum is to the learning rate what the learning rate is to weights, and it helps us produce better quality models

Common Architectural Principles of Deep Networks

Hypermeters-

2. Magnitude hyperparameters-

▪ **M o m e n t u m -**

▪ **Nesterov's m o m e n t u m**

- Momentum is a factor between 0.0 and 1.0 that is applied to the rate of change of the weights over time.
- Typical value for momentum between 0.9 and 0.99.

▪ **AdaGrad**

- One technique developed to help augment finding the “right” learning rate.
- Named in reference to how it “adaptively” uses sub-gradient methods to dynamically control the learning rate of an optimization algorithm.
- Monotonically decreasing and never increases the learning rate above whatever the base learning rate was set at initially.
- Square root of the sum of squares of the history of gradient computations.
- Speeds training in the beginning and slows it appropriately toward convergence, allowing for a smoother training process.

Common Architectural Principles of Deep Networks

Hypermeters-

2. Magnitude hyperparameters-

- **Momentum -**
 - **RMSProp-**
 - Very effective, but currently unpublished adaptive learning rate method
 - **Adam -**
 - ADAM (a more recently developed updating technique from the University of Toronto) derives learning rates from estimates of first and second moments of the gradients
- **AdaDelta-**
 - Variant of AdaGrad that keeps only the most recent history rather than accumulating it like AdaGrad does.

Common Architectural Principles of Deep Networks

Regularization-

- a measure taken against overfitting
- Regularization for hyperparameters helps modify the gradient so that it doesn't step in directions that lead it to overfit
 - Dropout
 - DropConnect
 - L1 penalty
 - L2 penalty
- Regularization works by adding an extra term to the normal gradient computed.

Common Architectural Principles of Deep Networks

Regularization-

- *Dropout-*

- Mechanism used to improve the training of neural networks by omitting a hidden unit.
- Speeds training.
- Driven by randomly dropping a neuron so that it will not contribute to the forward pass and backpropagation.

- *DropConnect*

- Does the same thing as Dropout, but instead of choosing a hidden unit, it mutes the connection between two neurons.
- **Dropout and DropConnect mute parts of the input to each layer, such that the neural network learns other portions.**
- **Zeroing-out parts of the data causes a neural network to learn more general representations**

Common Architectural Principles of Deep Networks

Regularization-

- The penalty methods L1 and L2, in contrast, are a way of preventing the neural network parameter space from getting too big in one direction.
- They make large weights smaller.
- **L1-**
 - Computationally inefficient in the nonsparse case, has sparse outputs, and includes built-in feature selection.
 - Multiplies the absolute value of weights rather than their squares.
 - This function drives many weights to zero while allowing a few to grow large, making it easier to interpret the weights

Common Architectural Principles of Deep Networks

Regularization-

- **L2** –
 - Computationally efficient due to analytical solutions and nonsparse outputs, but it does not do feature selection automatically.
 - Common and simple hyperparameter, adds a term to the objective function that decreases the squared weights.
 - Multiply half the sum of the squared weights by a coefficient called the weight-cost.
 - Improves generalization, smooths the output of the model as input changes, and helps the network ignore weights it does not use.

Common Architectural Principles of Deep Networks

Mini-batching-

- We send more than one input vector (a group or batch of vectors) to be trained in the learning system.
- Also allows us to compute certain linear algebra operations (specifically matrix-to-matrix multiplications) in a vectorized fashion.
- In this scenario we also have the option of sending the vectorized computations to GPUs if they are present.

Building Blocks of deep networks

Specific building blocks

- Feed-forward multilayer neural networks
 - Restricted Boltzmann Machines (RBMs)
 - Autoencoders
-
- Both RBMs and autoencoders are characterized by an extra layer-wise step for training.
 - They are often used for the pretraining phase in other larger deep networks.

Building Blocks of deep networks

■ RBMs-

- Model probability and are great at feature extraction.
- Feedforward networks in which data is fed through them in one direction with two biases rather than one bias as in traditional backpropagation feed-forward networks.
- Used in deep learning for the following:
 - Feature extraction
 - Dimensionality reduction
- The “restricted” part of the name “Restricted Boltzmann Machines” - connections between nodes of the same layer are prohibited (e.g., there are no visible-visible or hidden-hidden connections along which signal passes).
- Geoff Hinton—” A network of symmetrically connected, neuron-like units that make stochastic decisions about whether to be on or off.”

Building Blocks of deep networks

■ RBMs-

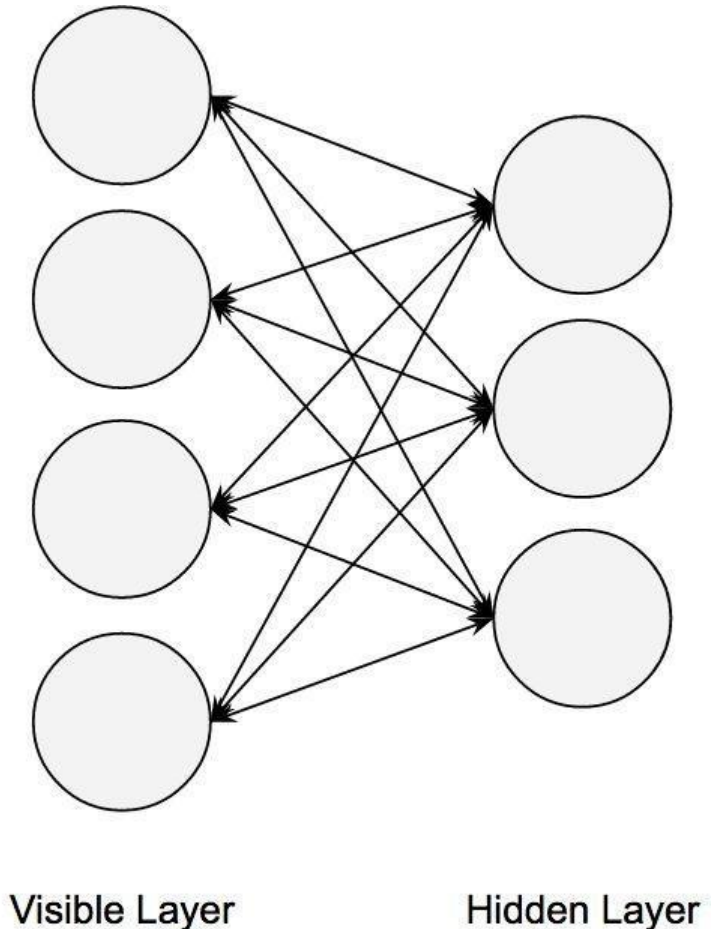
- Also, a type of autoencoder
- Used for pretraining layers in larger networks such as Deep Belief Networks.

■ Network layout

■ 5 main parts of a basic RBM:

- Visible units
- Hidden units
- Weights
- Visible bias units
- Hidden bias units

- A standard RBM has a visible layer and a hidden layer
- Every visible unit is connected to every hidden unit,
- Each layer of an RBM can be imagined as a row of nodes.
- The nodes of the visible and hidden layers are connected by connections with associated weights

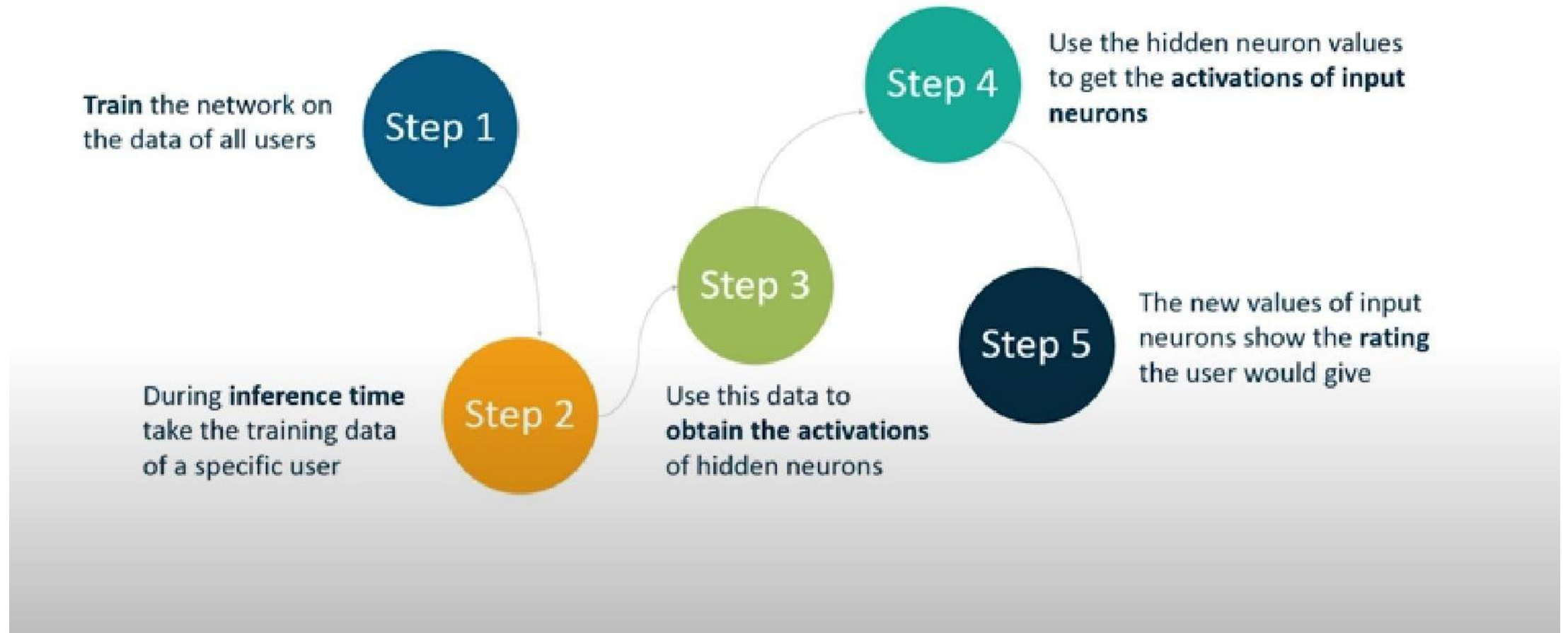


Building Blocks of deep networks

RBM-

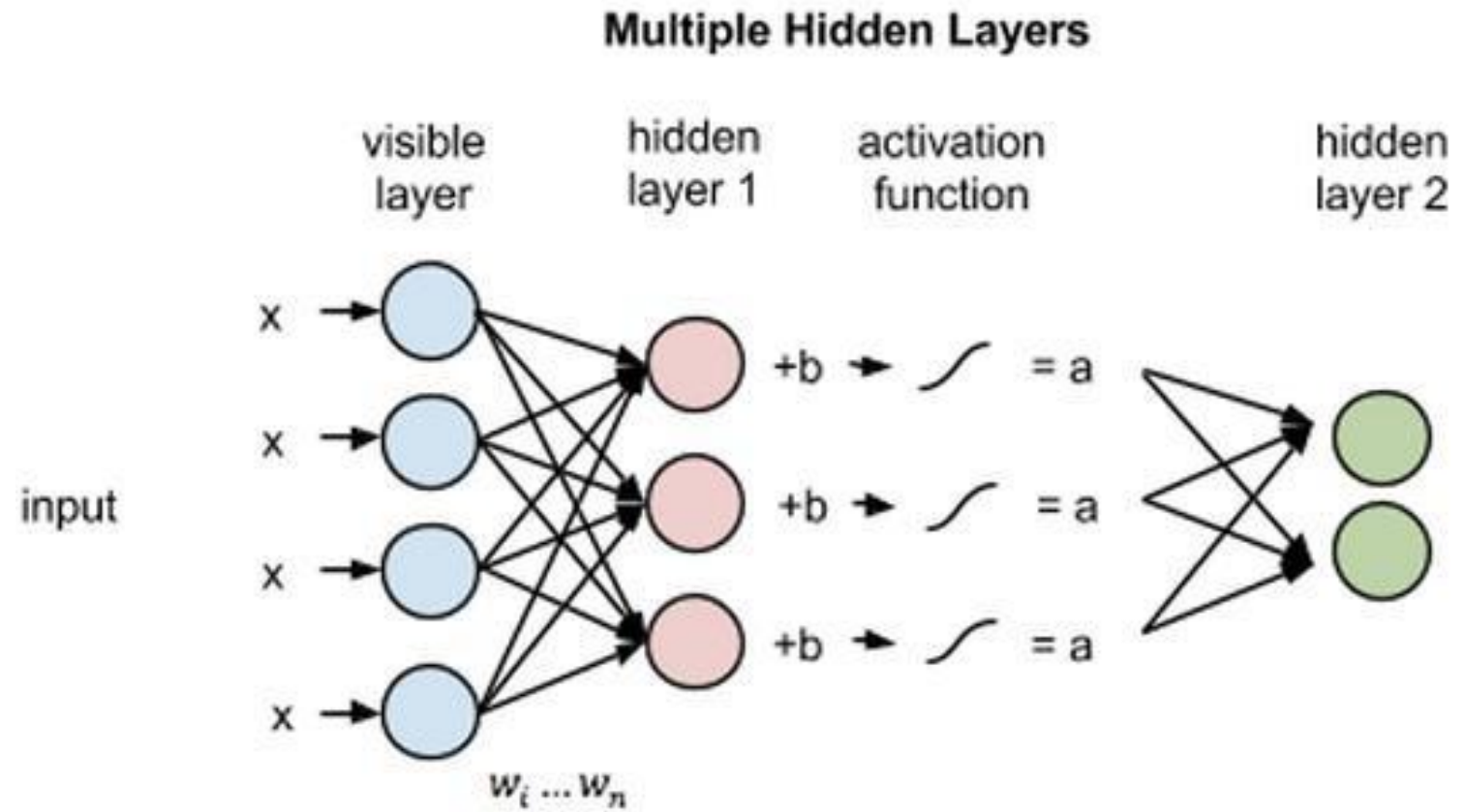
- **Network layout**
 - *Visible and hidden layers-*
 - *Connections and weights-*
 - *Biases-*
- **Training-**
- **Reconstruction**

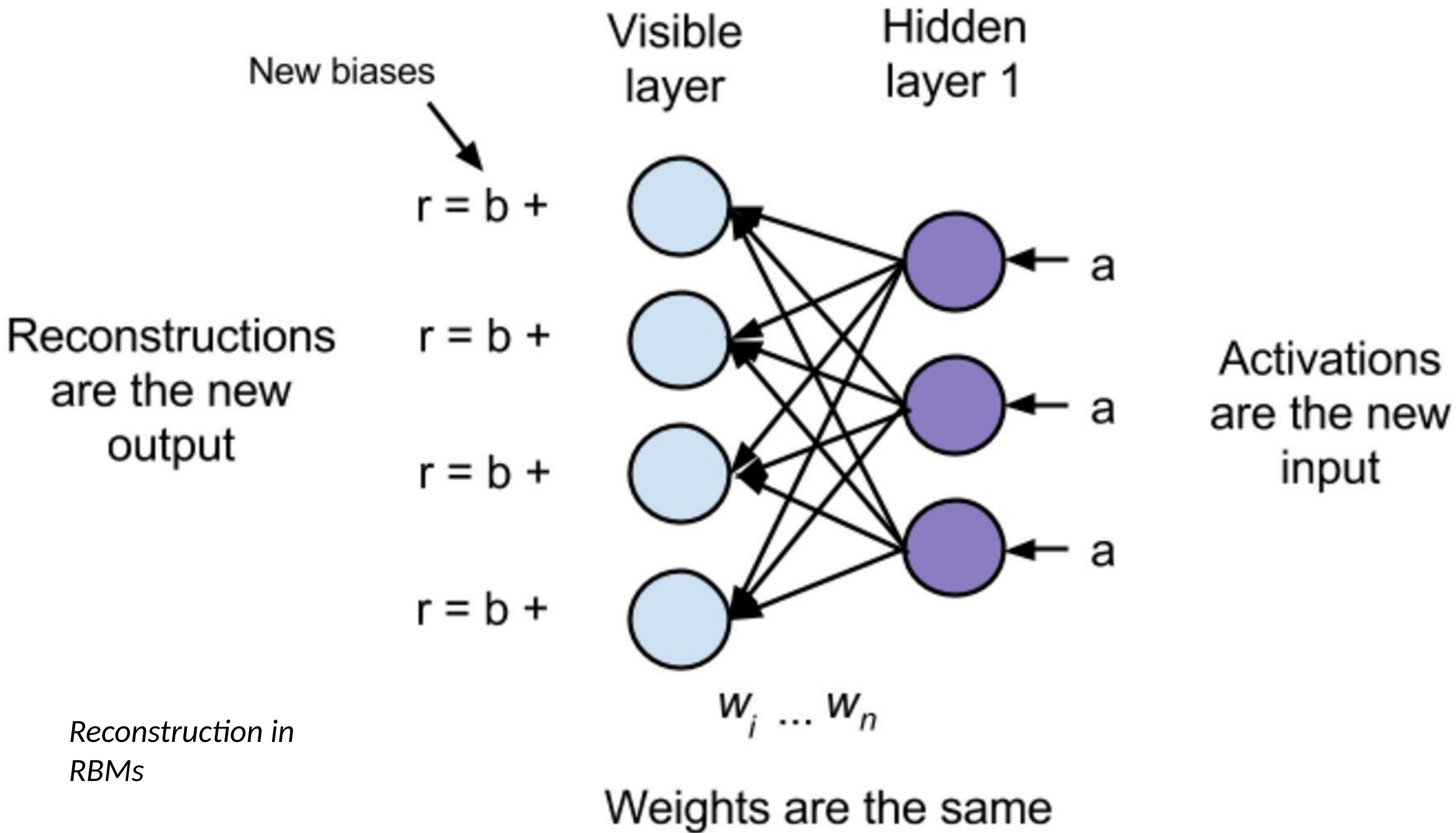
RBMS-



Building Blocks of deep networks

- RBMs-

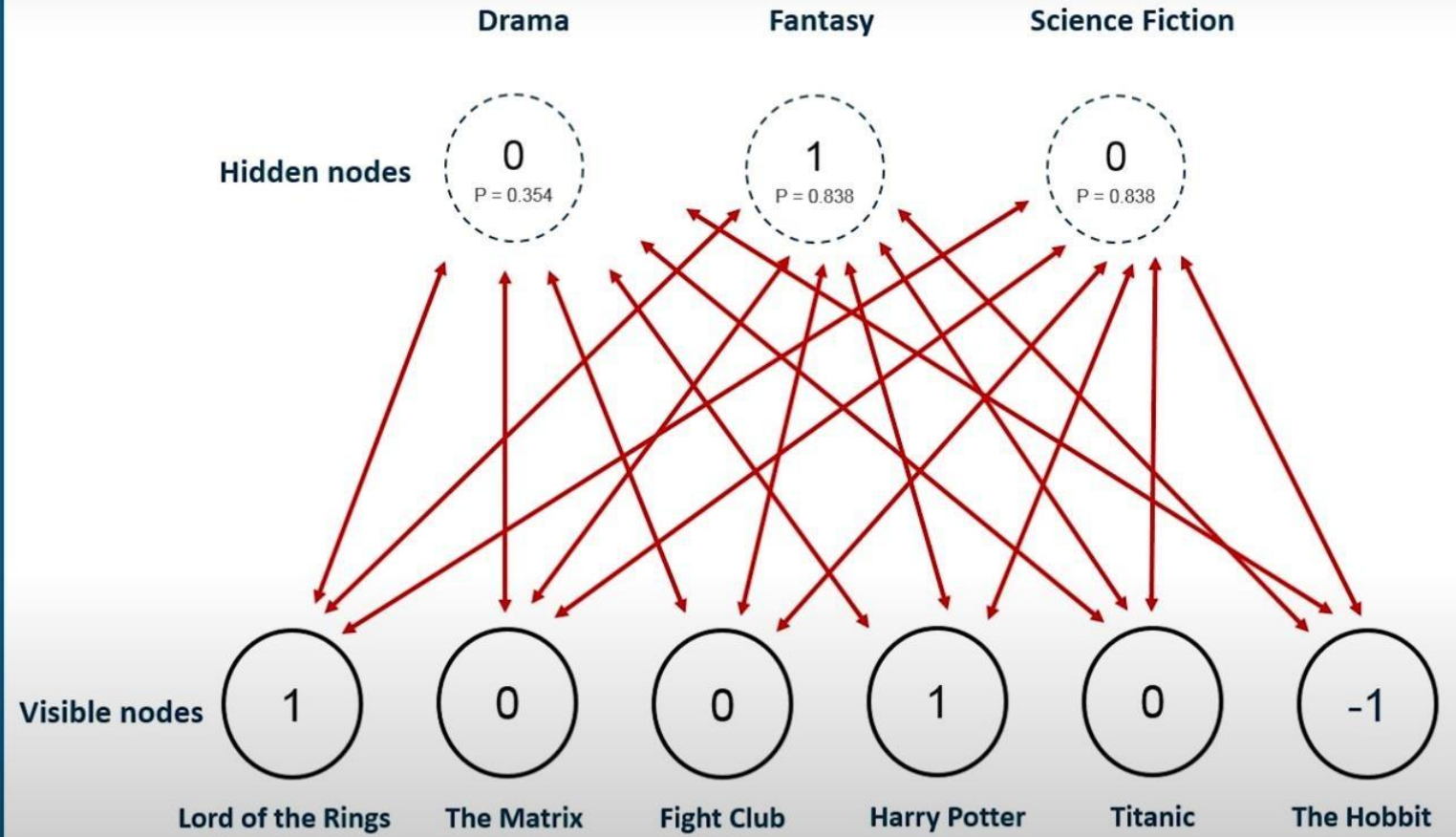




RBMS-

RBM: Collaborative Filtering

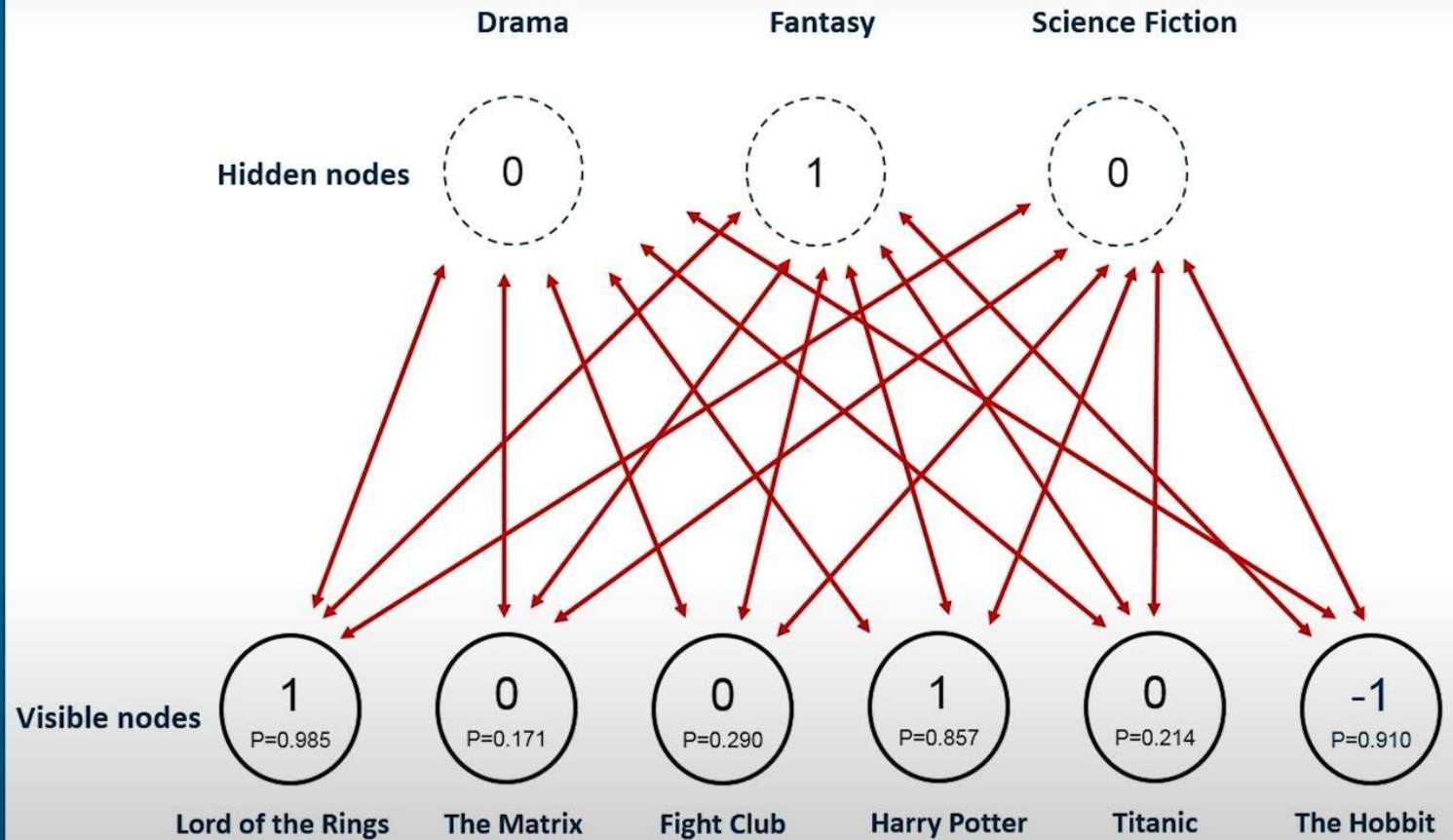
- 1 Recognizing Latent Factors in The Data
- 2 Using Latent Factors for Prediction



RBMS-

RBM: Collaborative Filtering

- 1 Recognizing Latent Factors in The Data
- 2 Using Latent Factors for Prediction



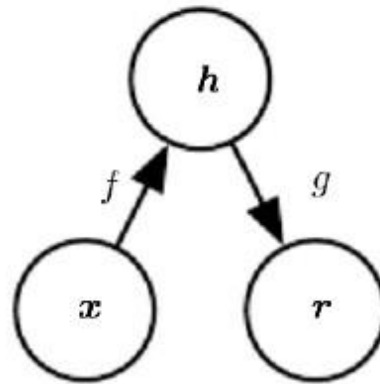
Building Blocks of deep networks

■ **Other uses of RBMs-**

- Dimensionality reduction
- Classification
- Regression
- Collaborative filtering
- Topic modelling
- Imbalance data problem
- Unstructured data
- Noisy label problems

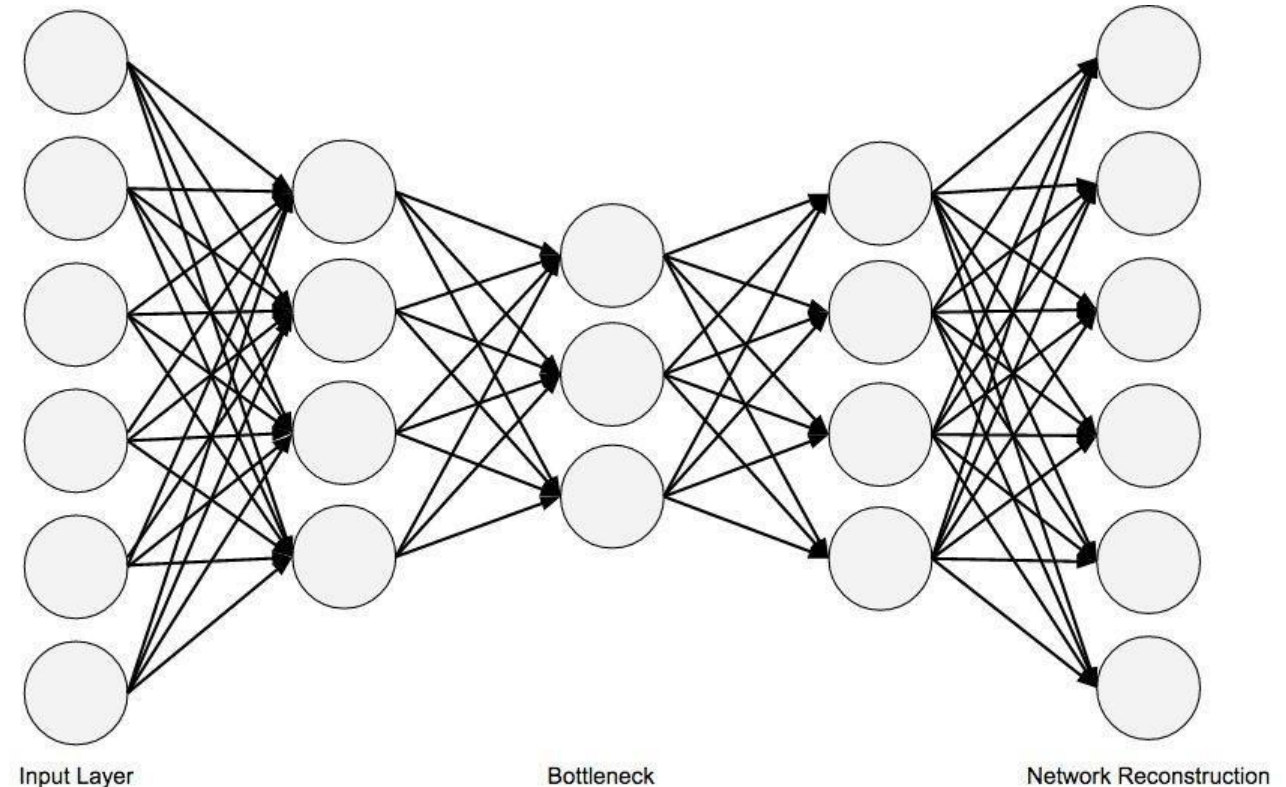
Building Blocks of deep networks

- Autoencoder is a neural network, which is trained to reproduce its input.
 - Encoder: $h = f(x)$
 - Decoder: $r = g(h)$

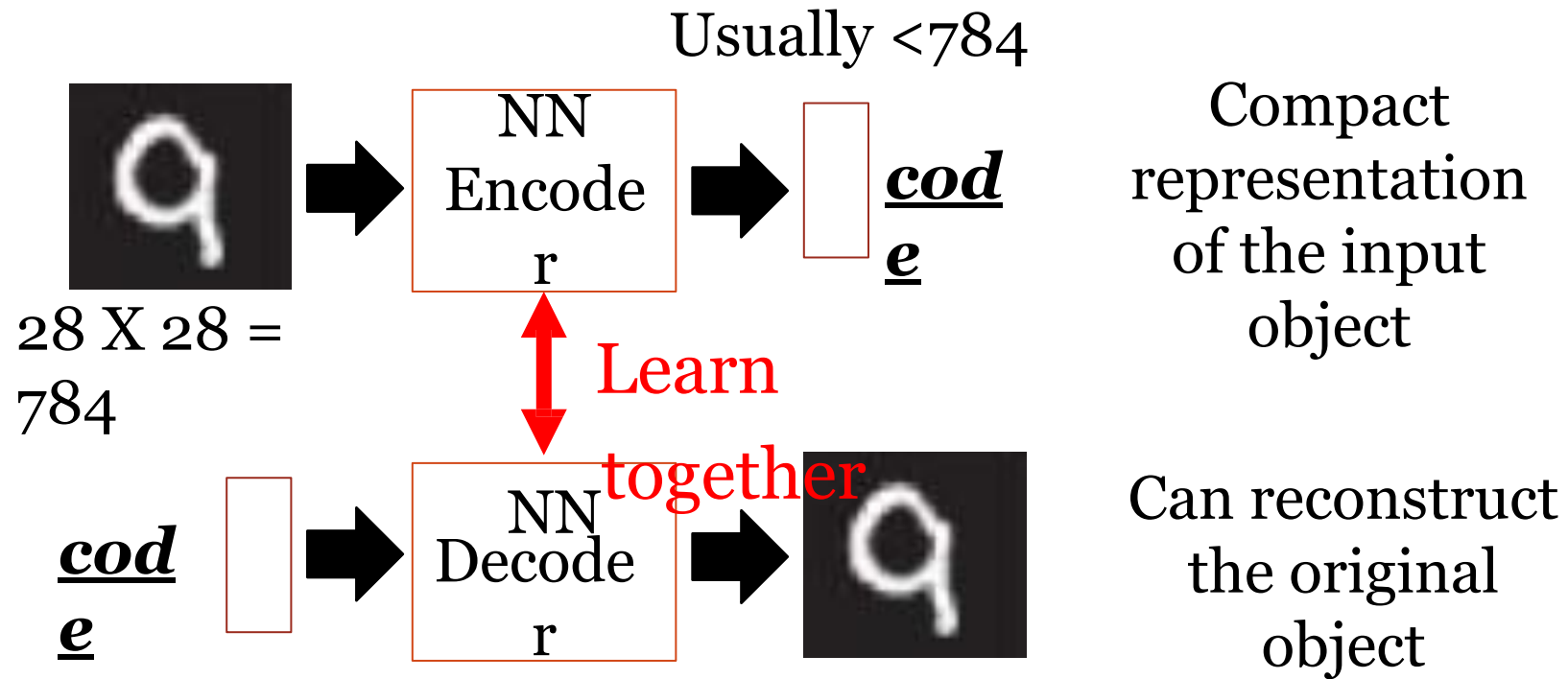


Autoencoders

- Used to learn compressed representations of datasets.
- To reduce a dataset's dimensionality.
- The output of the autoencoder network is a reconstruction of the input data in the most efficient form.
- **Similarities to multilayer perceptrons**
- Input layer, hidden layers of neurons, and then an output layer.
- Key diff- output layer in an autoencoder has the same number of units as the input layer does



Auto-Encoder



Deep Auto-Encoder

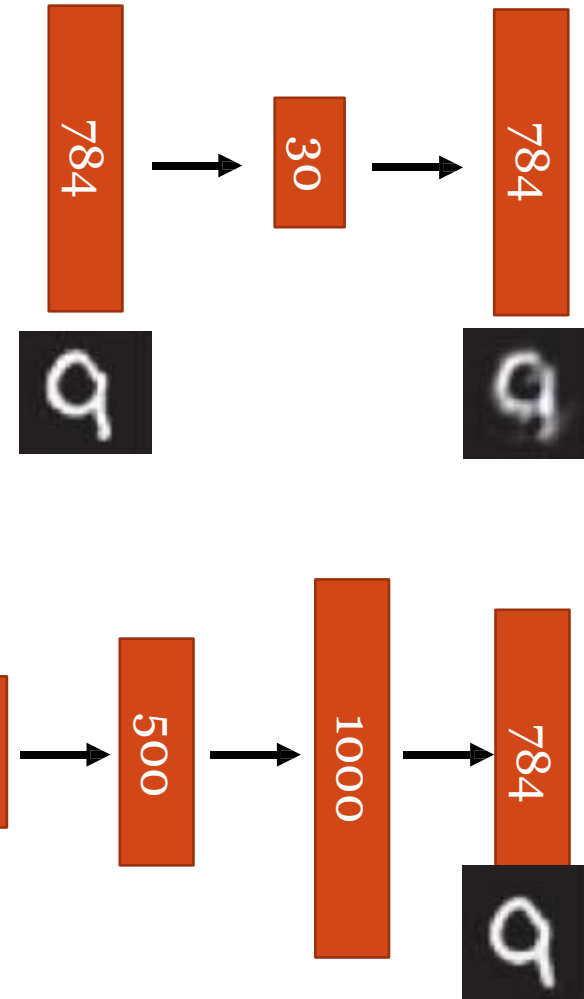
Original Image



PCA



Deep Auto-encoder



Autoencoders-

Defining features of autoencoders

- Autoencoders differ from multilayer perceptrons in a couple of ways:
 - They use unlabeled data in unsupervised learning.
 - They build a compressed representation of the input data.
- ***Unsupervised learning of unlabeled data***
 - The autoencoder learns directly from unlabeled data.
 - Connected to the second major difference between multilayer perceptrons and autoencoders.
- ***Learning to reproduce the input data***
 - The goal of a multilayer perceptron network is to generate predictions over a class (e.g., fraud versus not fraud).
 - An autoencoder is trained to reproduce its own input data.

Building Blocks of deep networks

Autoencoders-

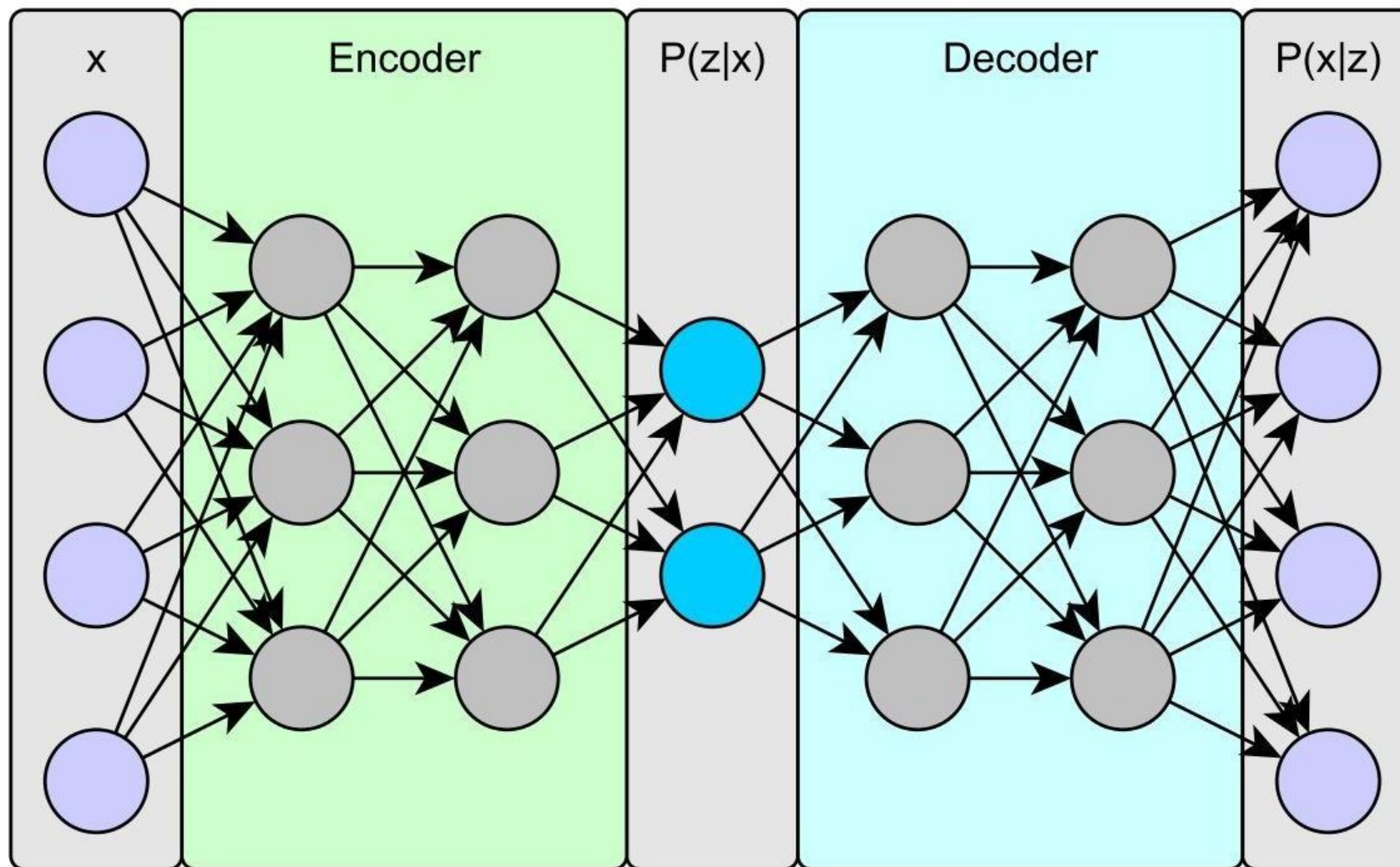
- **Training autoencoders**
- Autoencoders rely on backpropagation to update their weights.
- The main difference between RBMs and the more general class of autoencoders is in how they calculate the gradients

Building Blocks of deep networks

Common variants of autoencoders

- *Compression autoencoders*
- *Denoising autoencoders.*
- *Compression autoencoders*
 - The network input must pass through a bottleneck region of the network before being expanded back into the output representation.
- *Denoising autoencoders*
 - Autoencoder is given a corrupted version (e.g., some features are removed randomly) of the input and the network is forced to learn the uncorrupted output.

Building Blocks of deep networks



VAE network architecture

Building Blocks of deep networks

- Applications:
 - dimensionality reduction
 - visualization
 - feature extraction
 - \uparrow prediction accuracy
 - \uparrow speed of prediction
 - \downarrow memory requirements
 - semantic hashing
 - unsupervised pretraining

Major Architectures of deep networks

- Unsupervised Pretrained Networks (UPNs)
- Convolutional Neural Networks (CNNs)
- Recurrent Neural Networks
- Recursive Neural Networks

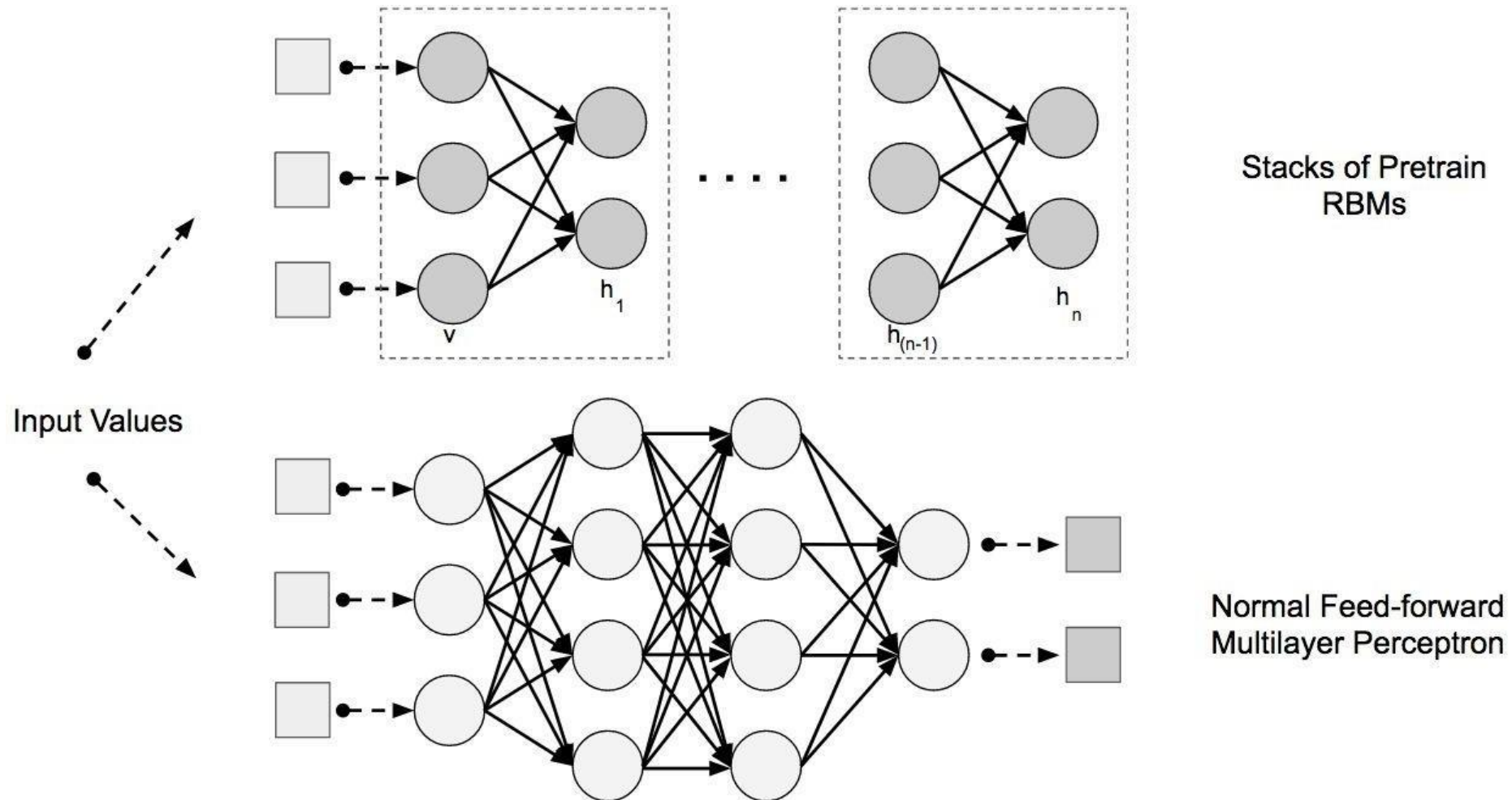
UNSUPERVISED PRETRAINED NETWORKS

Specific architectures:

- Autoencoders
- Deep Belief Networks (DBNs)
- Generative Adversarial Networks (GANs)

deep belief networks

- Composed of layers of Restricted Boltzmann Machines (RBMs) for the pretrain phase and then a feed-forward network for the fine-tune phase

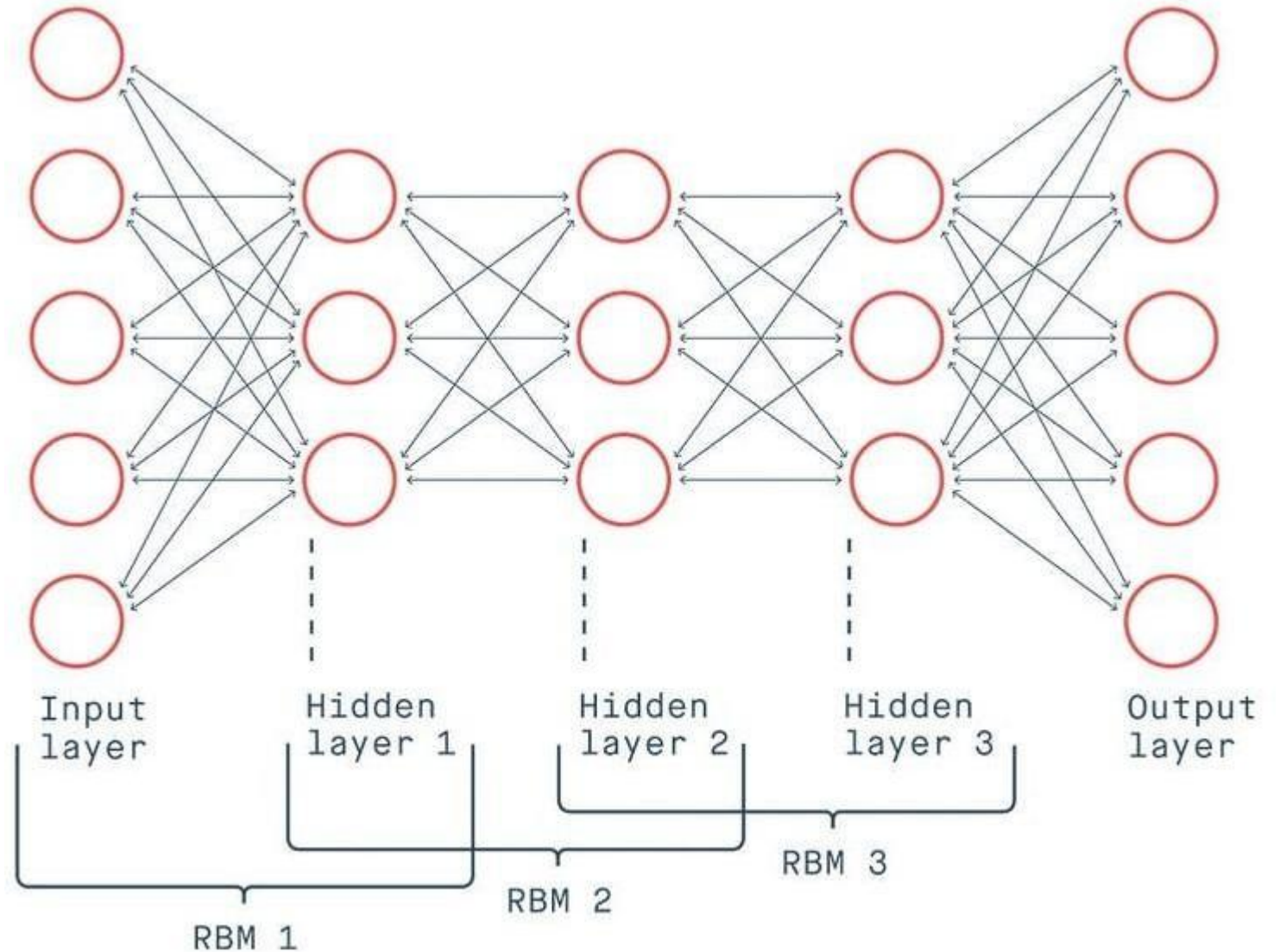


deep belief networks

- Backpropagation-local minima
- Solution-Pre-training
- Pre-training-Back propagation-reduced error rate

deep belief networks

- Network structure- DBN identical to MLP
- Training- DBN differs with MLP
- Role of Hidden nodes



deep belief networks

Training-

- Greedy learning algorithm- from bottom layer moving upside, fine tuning generative weights
- Greedy learning algorithms- Quick and efficient, Optimize weights at each layer
- Each layer learns entire input
- DBN work globally and regulate each layer in order, as the model slowly improves

deep belief networks

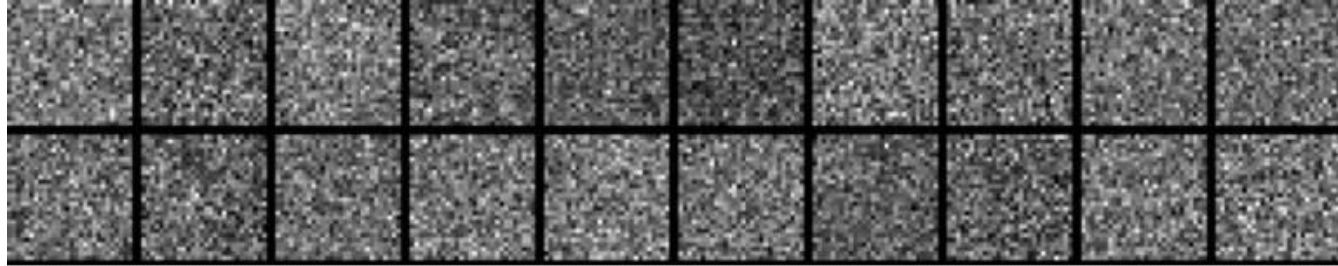
■ **Feature Extraction with RBM Layers-**

- RBMs -to extract higher-level features from the raw input vectors.
- The fundamental purpose of RBMs in the context of deep learning and DBNs -to learn these higher-level features of a dataset in an unsupervised training fashion.

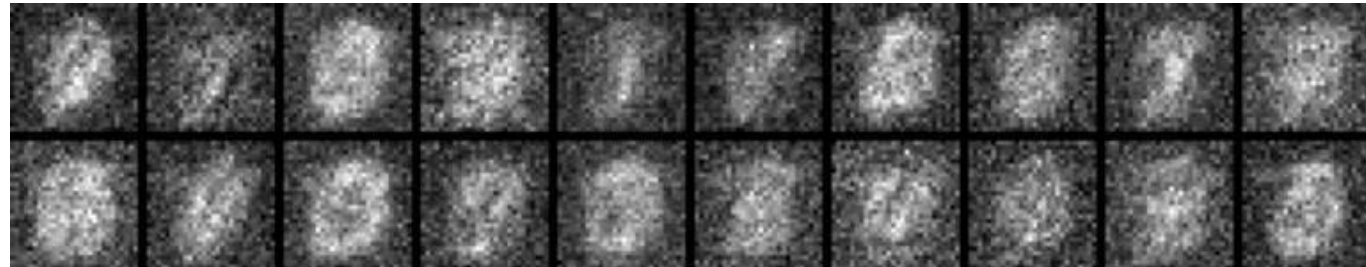
Learning higher-order features automatically-

Initializing the feed-forward network-fine-tune phase of DBNs

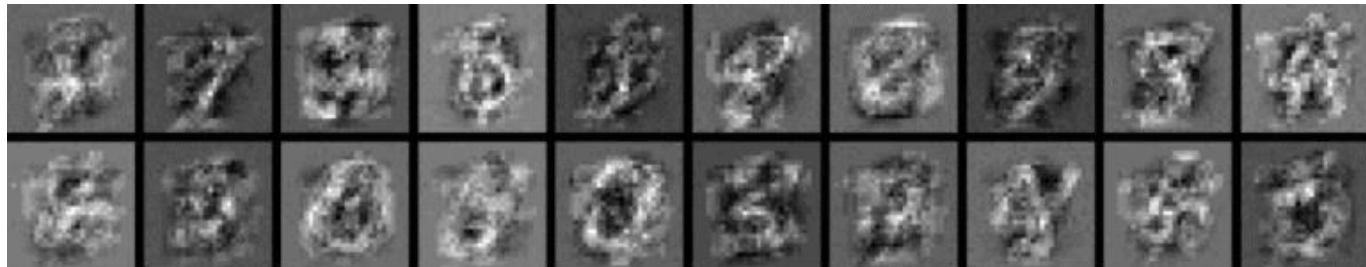
deep belief networks



Activation render at the beginning of training



Features emerge in later activation render Figure



Portions of MNIST digits emerge towards end of training

deep belief networks

- **Fine-tuning a DBN with a feed-forward multilayer neural network-**

- *Gentle backpropagation*
- *The output layer*

Applications-

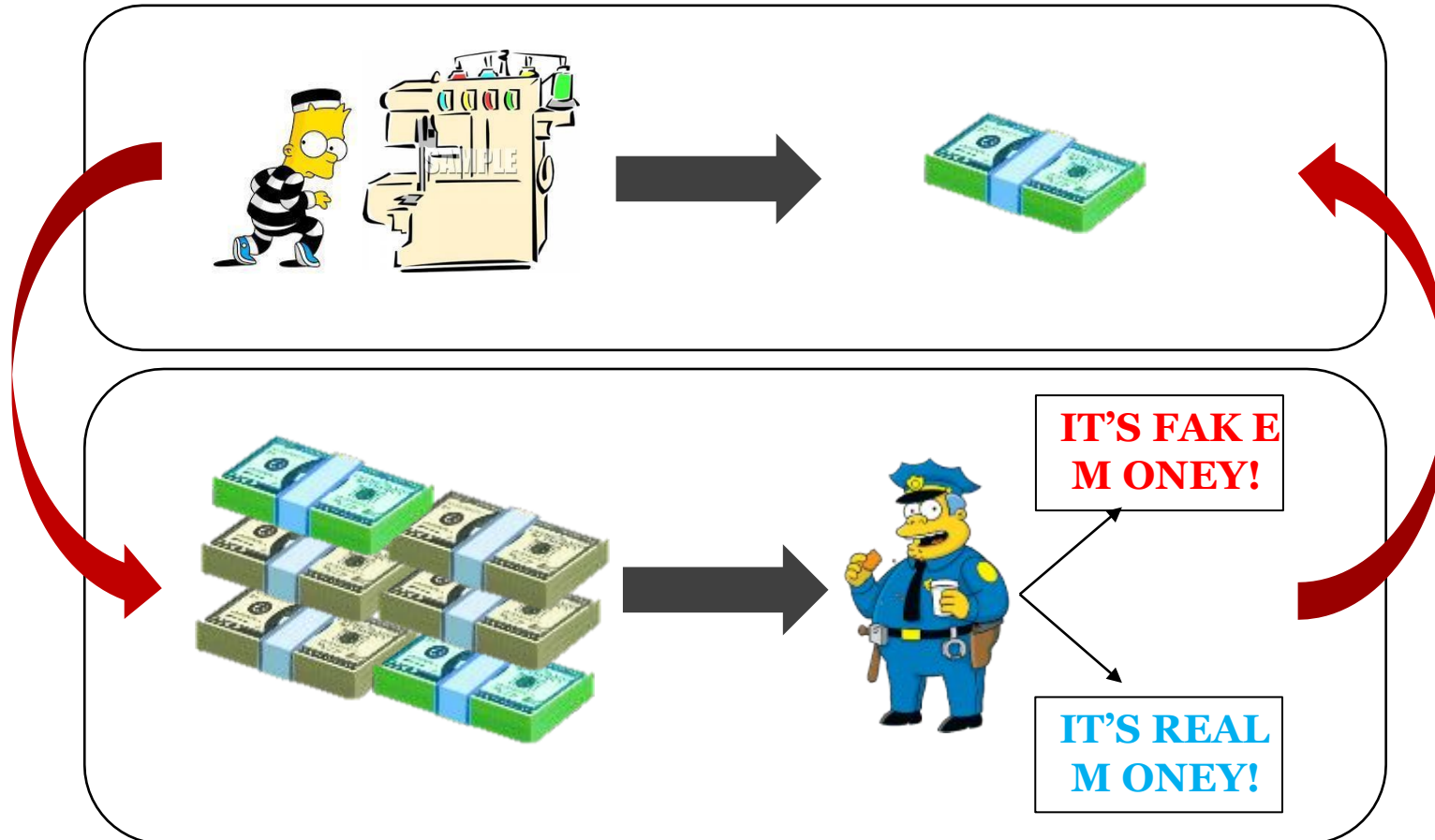
- *Image recognition*
- *Video recognition*
- *Motion capture data*

Generative Adversarial Networks (GAN)

- System of two neural networks competing against each other in a zero-sum game framework.
- They were first introduced by [Ian Goodfellow](#) *et al.* in 2014.
- Can learn to draw samples from a model that is similar to data that we give them.

The intuition behind GAN: Counterfeiters vs Police Game

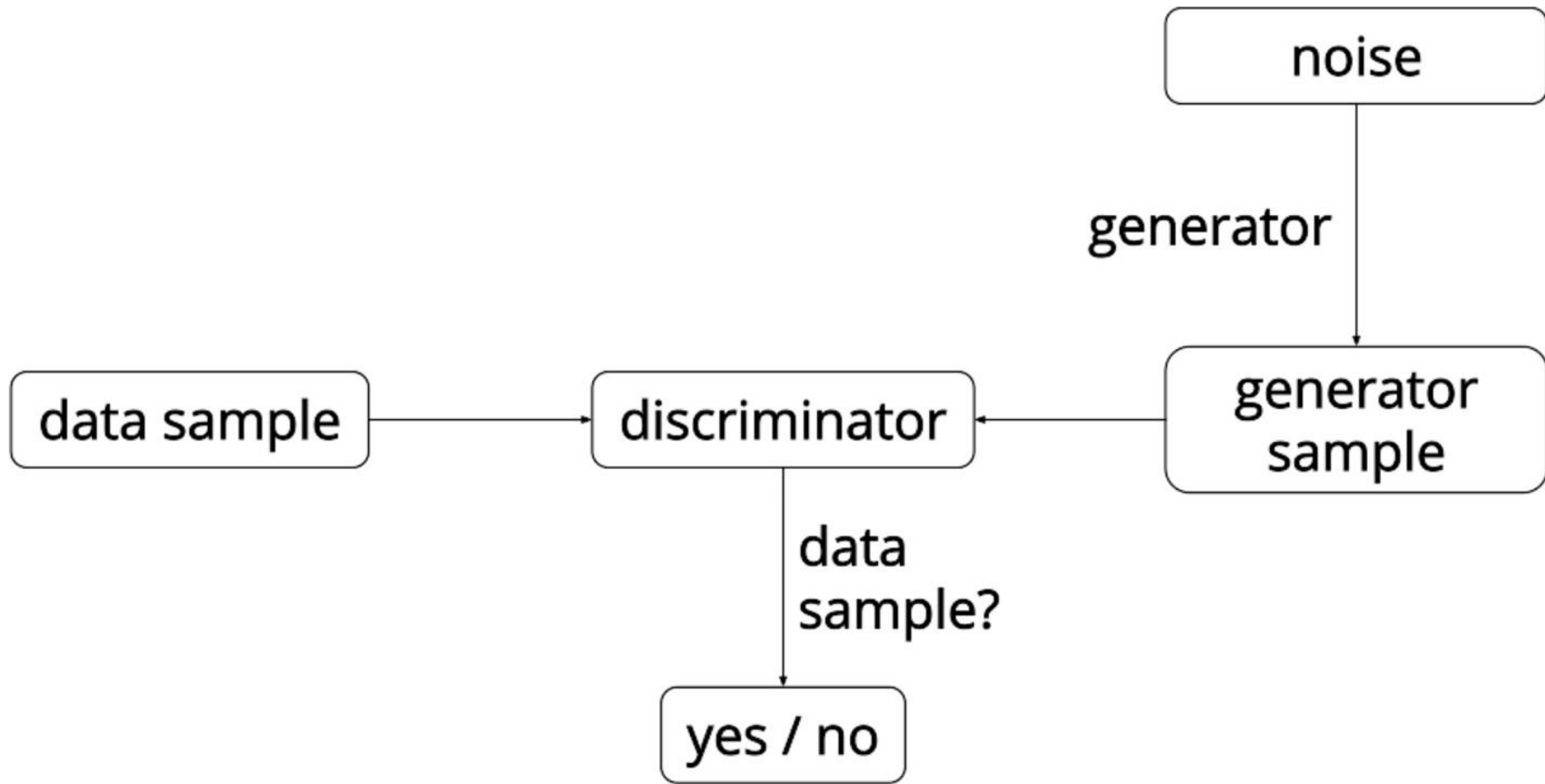
- The generative model can be thought of as analogous to a team of counterfeiters, trying to produce fake currency, while the discriminative model is analogous to the police, trying to detect the counterfeit currency. Competition in this game drives both teams to improve their methods until the counterfeits are indistinguishable from the genuine articles.



-Goodfellow et. al., "Generative Adversarial Networks" (2014)

What are GANs?

- **GAN is a new framework for estimating generative models via an adversarial process, in which we simultaneously train two models:**
 - **a generative model G that captures the data distribution.**
 - **a discriminative model D that estimates the probability that a sample came from the training data rather than G .**
- **The training procedure for G is to maximize the probability of D making a mistake. This framework corresponds to a minimax two-player game.**
- **In the case where G and D are defined by multilayer perceptrons, the entire system can be trained with backpropagation.**

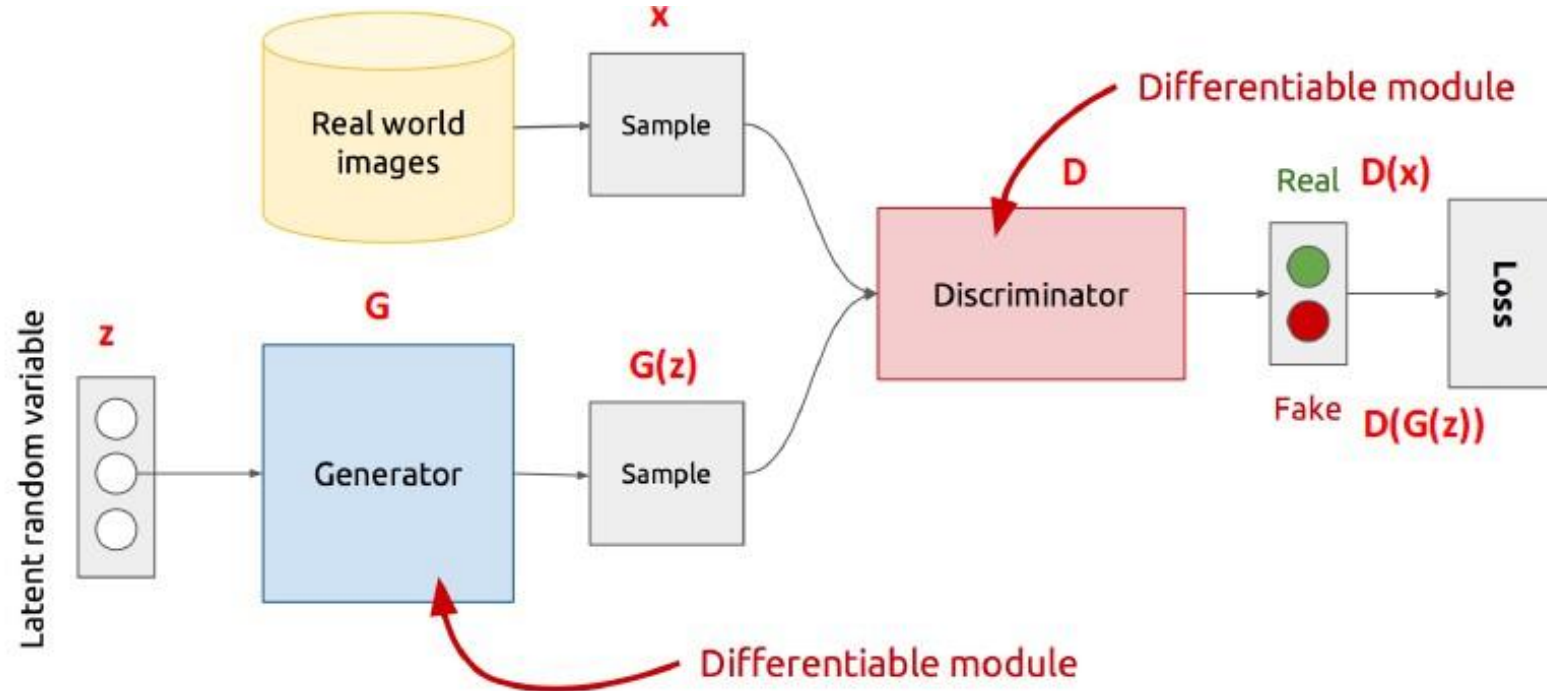


Overview of GANs

Source:

<https://ishmaelbelghazi.github.io/ALI>

GAN'S ARCHITECTURE

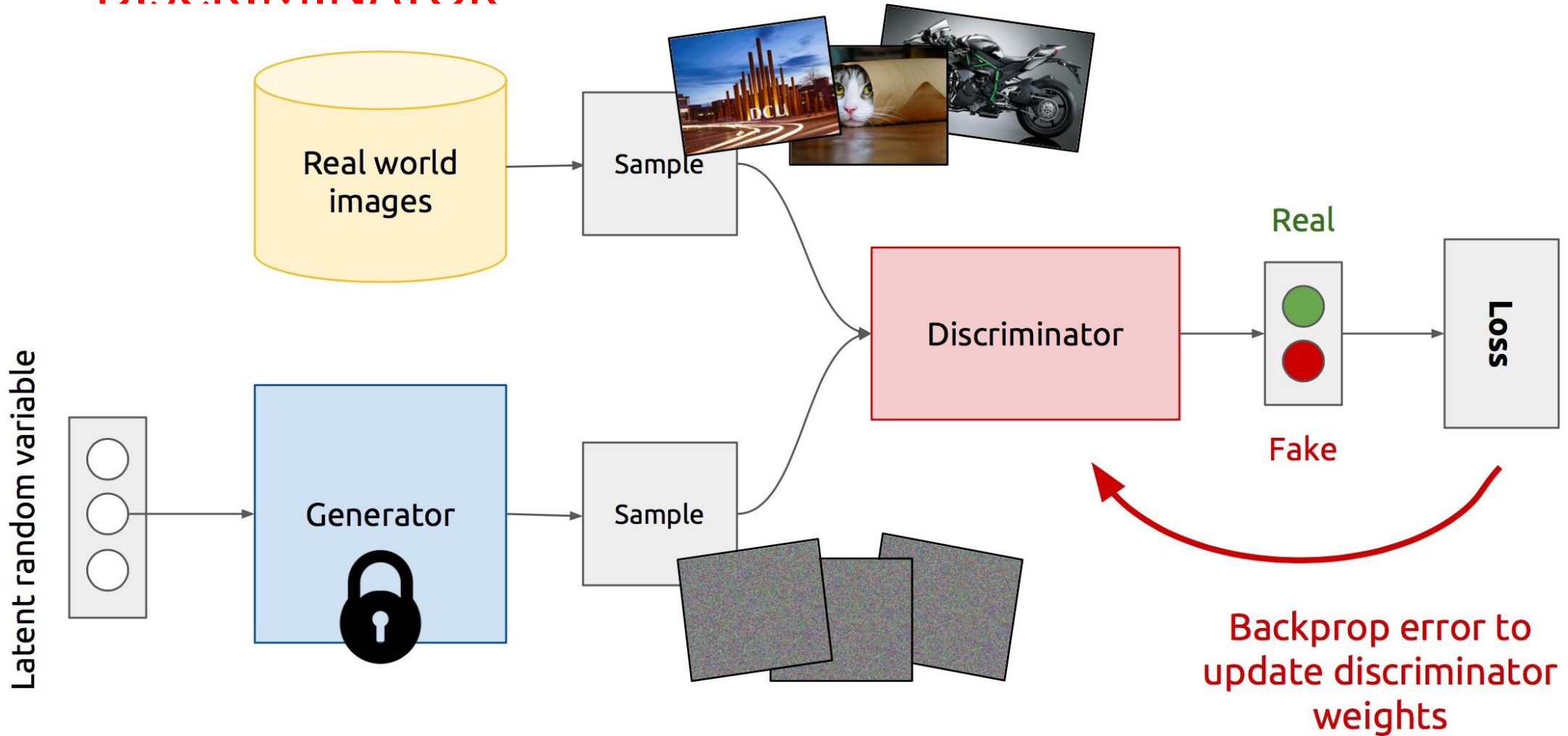


- **Z** is some random noise (Gaussian/Uniform).
- **Z** can be thought as the latent representation of the image.

Discriminative Models

- A **discriminative** model learns a function that maps the input data (x) to some desired output class label (y).
- In probabilistic terms, they directly learn the conditional distribution $P(y|x)$.

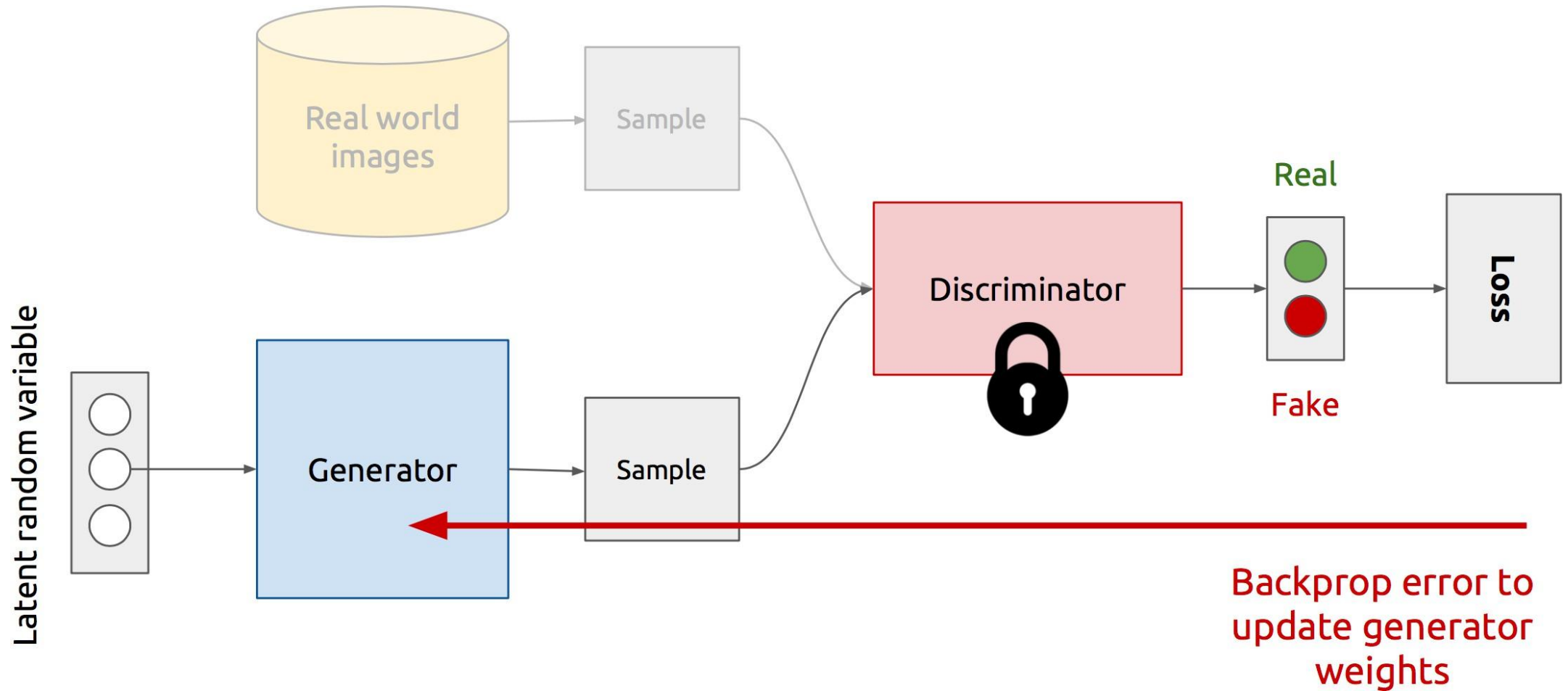
TRAINING DISCRIMINATOR



Generative Models

- A **generative** model tries to learn the joint probability of the input data and labels simultaneously i.e. $P(x,y)$.
- Potential to understand and explain the underlying structure of the input data even when there are no labels.

TRAINING GENERATOR

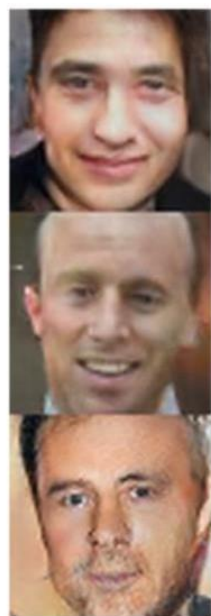


How GANs are being used?

- Applied for modelling natural images.
- Performance is fairly good in comparison to other generative models.
- Useful for unsupervised learning tasks.



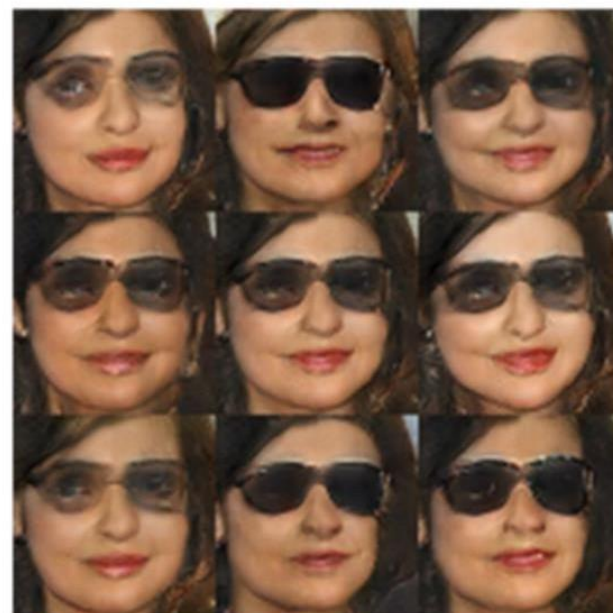
man
with glasses



man
without glasses



woman
without glasses



woman
with glasses

How to train GANs?

- Objective of generative network - increase the error rate of the discriminative network.
- Objective of discriminative network – decrease binary classification loss.
- Discriminator training - backprop from a binary classification loss.
- Generator training - backprop the negation of the binary classification loss of the discriminator.

Loss Functions

$$\mathcal{L}(\hat{x}) = \min_{x \in data} (x - \hat{x})^2$$

Generator

$$D_G^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_g(x)}$$

Discriminator

Improved Techniques for testing of GANs by Salimans et. al

- One-sided Label smoothing - replaces the 0 and 1 targets for a classifier with smoothed values, like .9 or .1 to reduce the vulnerability of neural networks to adversarial examples.
- Virtual batch Normalization - each example x is normalized based on the statistics collected on a reference batch of examples that are chosen once and fixed at the start of training, and on x itself.

Variations of GANs

- Several new concepts built on top of GANs have been introduced –
 - InfoGAN – Approximate the data distribution and learn interpretable, useful vector representations of data.
 - Conditional GANs - Able to generate samples taking into account external information (class label, text, another image). Force G to generate a particular type of output.

Major Difficulties

- Networks are difficult to converge.
- Ideal goal – Generator and discriminator to reach some desired equilibrium but this is rare.
- GANs are yet to converge on large problems (E.g. Imagenet).

Common Failure Cases

- The discriminator becomes too strong too quickly and the generator ends up not learning anything.
- The generator only learns very specific weaknesses of the discriminator.
- The generator learns only a very small subset of the true data distribution.

Solution

- Normalize the inputs
- A modified loss function
- BatchNorm
- Avoid Sparse Gradients: ReLU, MaxPool
- Use Soft and Noisy Labels
- DCGAN / Hybrid Models
- Track failures early (D loss goes to 0: failure mode)
- If you have labels, use them
- Add noise to inputs, decay over time

Applications of GANs

Labels to Street Scene



input



output

Aerial to Map

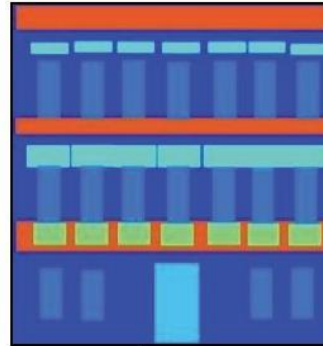


input



output

Labels to Facade



input



output

BW to Color



input



output

Day to Night



input



output

Edges to Photo



input



output

Isola, P., Zhu, J. Y., Zhou, T., & Efros, A. A. "Image-to-image translation with conditional adversarial networks". arXiv preprint arXiv:1611.07004. (2016).

Applications of GANs

- Motivation
- Given a text description, generate images closely associated.
- Uses a conditional GAN with the generator and discriminator being condition on “dense” text embedding.

this small bird has a pink breast and crown, and black primaries and secondaries.



this magnificent fellow is almost all black with a red crest, and white cheek patch.



the flower has petals that are bright pinkish purple with white stigma



this white and yellow flower have thin white petals and a round yellow stamen



WHY USE GANS FOR GENERATION?

- Can be trained using back-propagation for Neural Network based Generator/Discriminator functions.
- Sharper images can be generated.
- Faster to sample from the model distribution: *single* forward pass generates a *single* sample.

Conclusions

- GANs are generative models that are implemented using two stochastic neural network modules: **Generator** and **Discriminator**.
- **Generator** tries to generate samples from random noise as input
- **Discriminator** tries to distinguish the samples from Generator and samples from the real data distribution.
- Both networks are trained adversarially (in tandem) to fool the other component. In this process, both models become better at their respective tasks.

Conclusions

- Train GAN – Use discriminator as base model for transfer learning and the fine-tuning of a production model.
- A well-trained generator has learned the true data distribution well - Use generator as a source of data that is used to train a production model.

- <https://tryolabs.com/blog/2016/12/06/major-advancements-deep-learning-2016/>
- <https://blog.waya.ai/introduction-to-gans-a-boxing-match-b-w-neural-nets-b4e5319cc935#.6l7zh8u50>
- https://en.wikipedia.org/wiki/Generative_adversarial_networks
- <http://blog.aylien.com/introduction-generative-adversarial-networks-code-tensorflow/>
- <https://github.com/soumith/ganhacks>
- <https://www.youtube.com/watch?v=i64KpxyaLpo>