

Introduction to Neural Networks

Kun Hu

April 13, 2020

Overview

Overview of (Feedforward) Neural Networks

Structure of Neural Networks

Components of Neural Networks

Train Neural Networks

Neural Networks In Action

Build a DNN from Scratch

Software Packages

Limitation of Neural Networks

Hyper-parameter Tuning

Perceptron (Single Neuron) Flow Chat

Inputs Weights

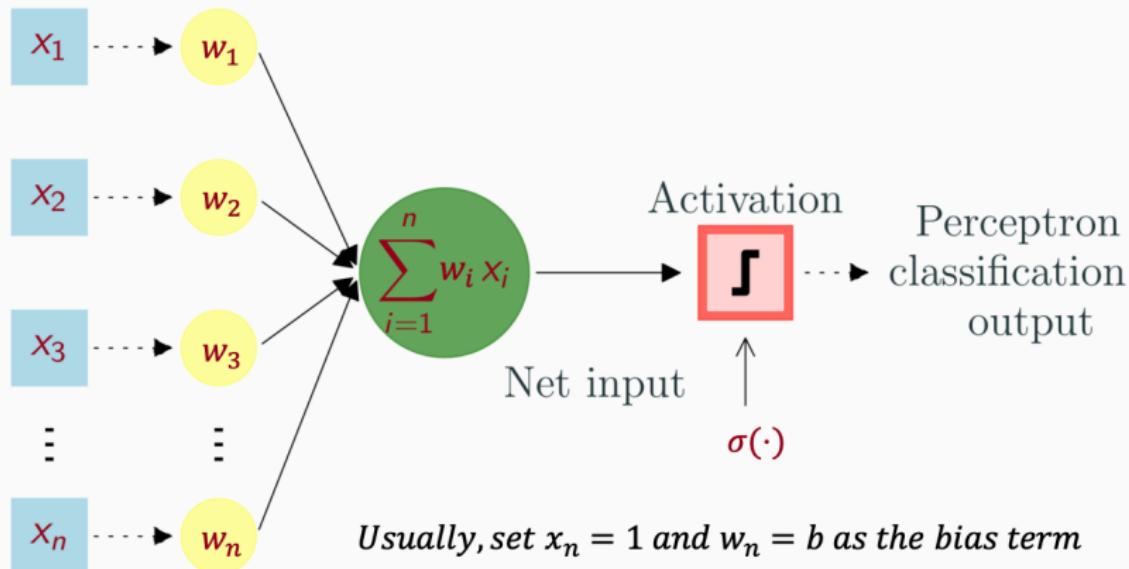
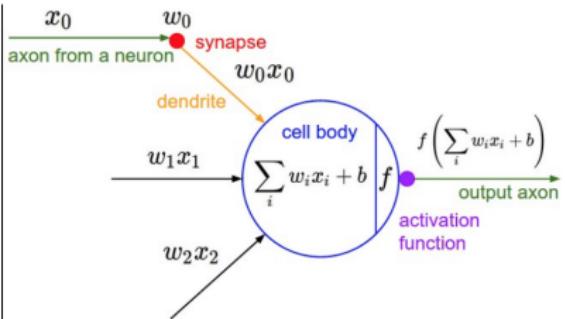
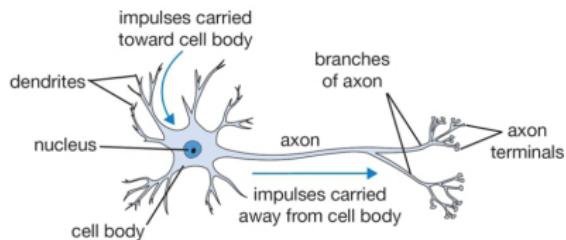


Figure: neuron: a mapping from \mathbb{R}^N to \mathbb{R}

Biological Analog



A cartoon drawing of a biological neuron (left) and its mathematical model (right).

MultiLayer Perceptrons

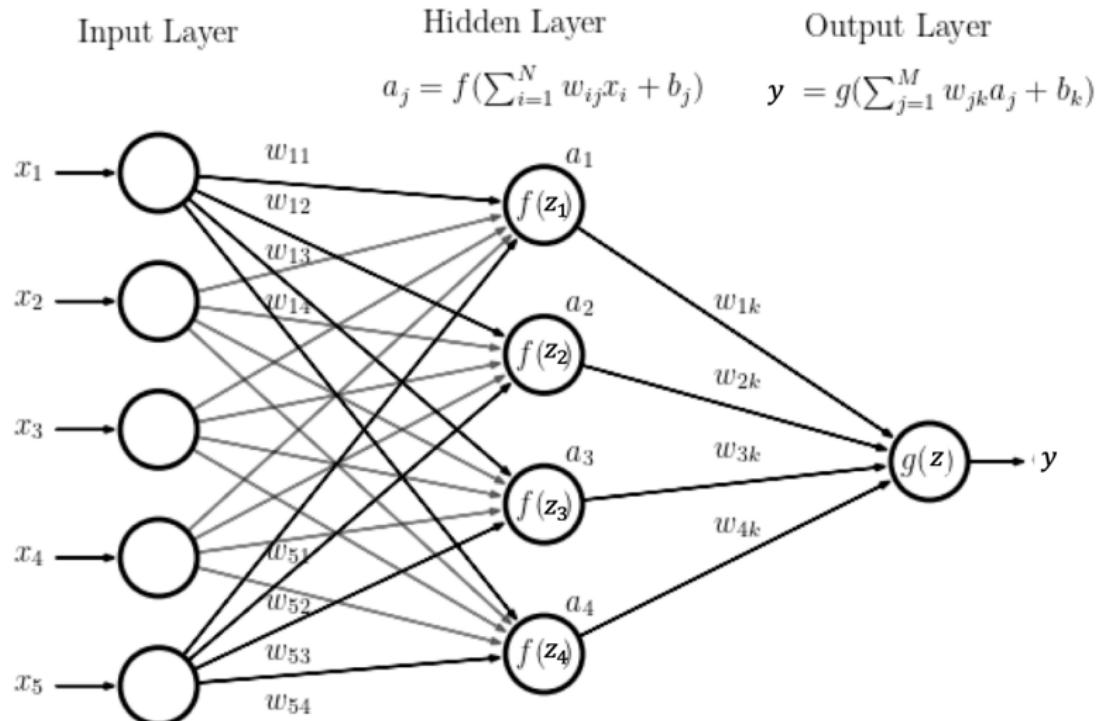
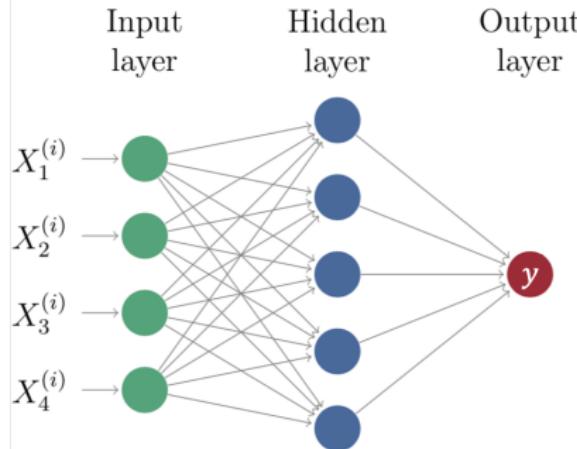


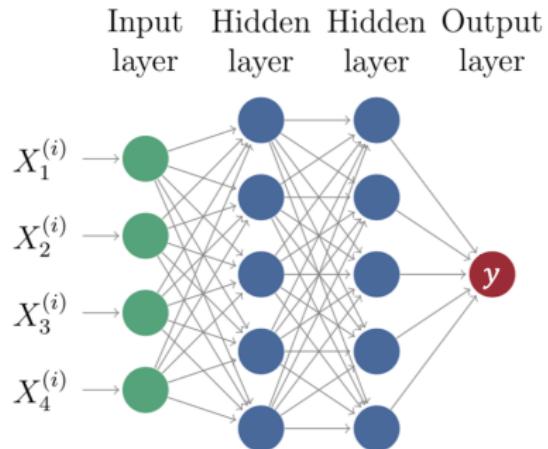
Figure: A MultiLayer Perceptrons with One Hidden Layer

(Deep) Feedforward Neural Network

(a) Single Layer



(b) Deep Neural Network



- ▶ i is observation i 's index; each one has 4 features (green circle)
- ▶ blue circle: linear and non-linear transformation ($\phi(z)$)
- ▶ red circle: linear combination to produce the final output

Activation Function

Transformation of the linear combination of inputs

$$\mathbf{a}^l = \phi(\mathbf{z}^l) = \phi(\mathbf{b}^l + \mathbf{w}^l \cdot \mathbf{a}^{l-1})$$

- ▶ $l = 1, \dots, L$ is the layer index, with $l = L$ as the output layer
- ▶ **b** is sometimes referred as “bias”
- ▶ Without nonlinear $\phi(z)$, y is just a linear combination of \mathbf{x}

Activation Functions

Pros

- Smooth
- Bounded output
- Cons
- Non-zero centered output
- Computational expensive
- Vanishing gradient

Pros

- Output centered at 0 where gradient is the steepest
- Cons
- Similar to sigmoid

Pros

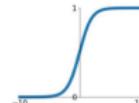
- Quick computation
- Non-saturating

Cons

- Dying neurons: both output and derivatives become 0 when inputs are negative

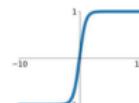
Sigmoid

$$\sigma(x) = \frac{1}{1+e^{-x}}$$



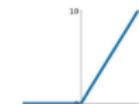
tanh

$$\tanh(x)$$



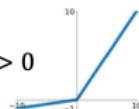
ReLU

$$\max(0, x)$$



Leaky ReLU

$$\max(ax, x), 1 \gg a > 0$$

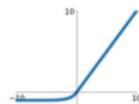


PReLU

$$\max(\alpha x, x), \alpha \text{ is learned by training}$$

ELU

$$\begin{cases} x & x \geq 0 \\ \alpha(e^x - 1) & x < 0 \end{cases}$$



Pros

- Prevent dying neurons
- Similar to ReLU
- Cons
- Another hyper-parameter α to tune

Pros

- α is learned by backpropagation

Cons

- Can over-fit

Pros

- Smooth when $\alpha = 1$
- No dying neurons/vanishing gradient
- Faster convergence
- Output closer to 0 (self-normalizing for SELU)

Cons

- Computational expensive

- In general, SELU > ELU > leaky ReLU (and its variants) > ReLU > tanh > sigmoid
- For network structure that prevents self-normalizing, ELU can outperform SELU
- Actual performance varies by problem. Better to trial-and-error.

Loss (Cost) Functions

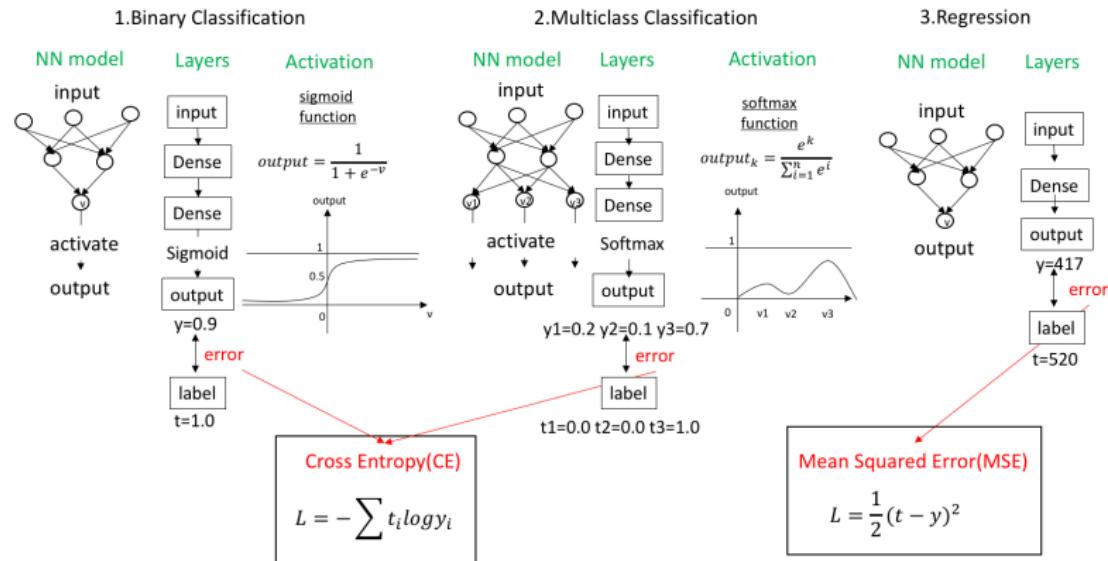


Figure: Common Loss Functions

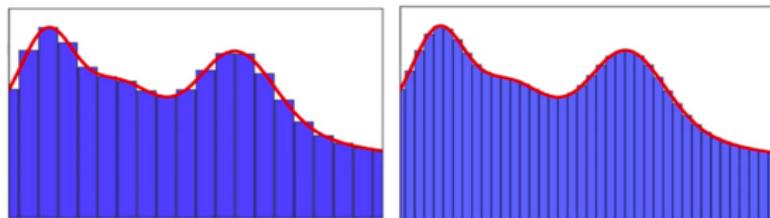
Universal Approximation Theorem

Hornik, K., Stinchcombe, M., and White, H. (1989)

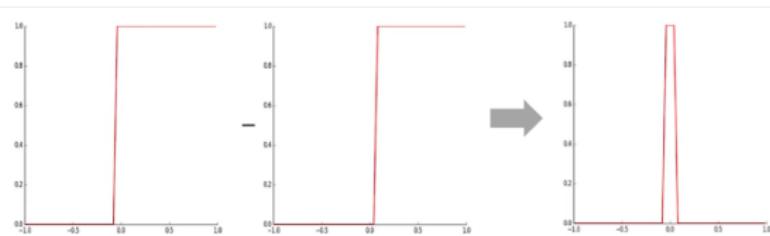
One-hidden layer feed-forward network with arbitrary squashing functions can approximate any Borel measurable function.

- ▶ It requires a sufficient amount of neurons (unbounded width)
- ▶ Extension to bounded width and Lebesgue-measurable functions: Lu. et al 2017, Hanin and Sellke 2017

Universal Approximation Theorem (Proof by Graphs*)



1. Decompose a function into a series of rectangles (bars)



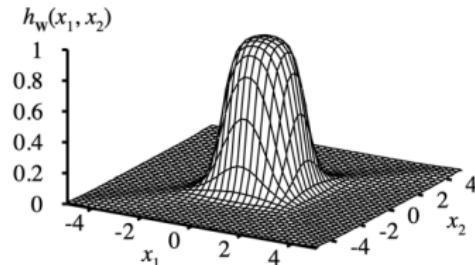
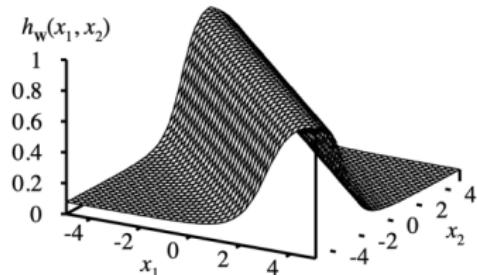
2. Each rectangle can be approximated by the difference between two sigmoid functions

A. Take a very steep sigmoid function

B. Copy the function but shift rightward a bit

C. Subtract B from A to get a rectangle

Universal Approximation Theorem (Proof by Graphs*)



Combine two opposite-facing threshold functions to make a ridge

Combine two perpendicular ridges to make a bump

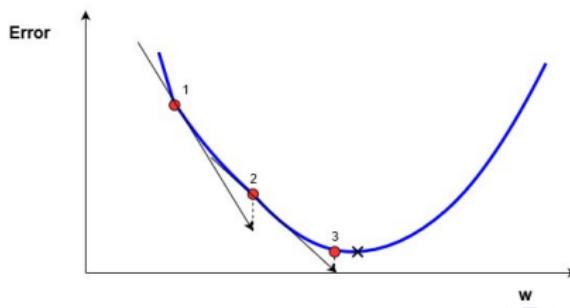
Add bumps of various sizes and locations to fit any surface

Finding the NN Parameters $\Theta = \{\mathbf{w}^l, \mathbf{b}^l\}_{l=1,\dots,L}$

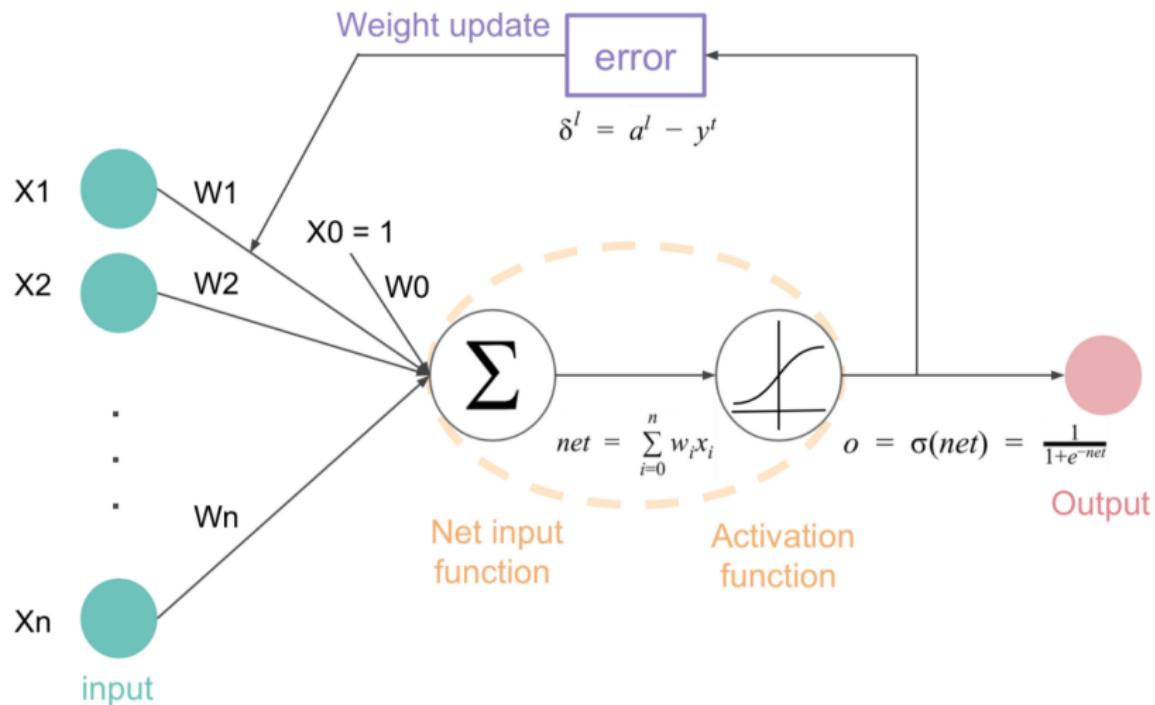
$$\Theta^* = \arg \min_{\Theta} \sum_{i=1}^N \mathbf{C}(\mathbf{Y}_i, \mathbf{X}_i; \Theta)$$

- ▶ \mathbf{X}_i is the input vector for observation i , \mathbf{Y}_i is the target
- ▶ $\mathbf{C}(\cdot)$ is the loss function
- ▶ Θ is updated according to (ϵ is the learning rate)

$$\Theta' = \Theta - \epsilon \nabla_{\Theta} \sum_{i=1}^{1 \leq j \leq N} \mathbf{C}(\mathbf{Y}_i, \mathbf{X}_i; \Theta)$$



How to Calculate Gradient: Backpropagation (Chain Rule)



How to Calculate Gradient: Backpropagation (Chain Rule)

A one-unit per-layer network without bias term.

$$\text{Cost function } C = \frac{1}{2}(a^2 - y)^2$$

Forward propagation to calculate output

$$x \rightarrow \sigma(w^1 x = z^1) \rightarrow a^1 \rightarrow \sigma(w^2 a^1 = z^2) \rightarrow a^2 \rightarrow \frac{1}{2}(a^2 - y)^2 = C$$

Backward propagation to calculate derivative

$$\delta^2 \equiv \frac{\partial C}{\partial z^2} = \nabla C \cdot \sigma'(z^2) = \frac{\partial C}{\partial a^2} \frac{\partial a^2}{\partial z^2} = (a^2 - y)\sigma'(z^2)$$

$$\delta^1 \equiv \frac{\partial C}{\partial z^1} = \frac{\partial C}{\partial z^2} \frac{\partial z^2}{\partial a^1} \frac{\partial a^1}{\partial z^1} = \delta^2 w^2 \sigma'(z^1)$$

$$\frac{\partial C}{\partial w^2} = \frac{\partial C}{\partial z^2} \frac{\partial z^2}{\partial w^2} = \delta^2 a^1$$

$$\frac{\partial C}{\partial w^1} = \frac{\partial C}{\partial z^1} \frac{\partial z^1}{\partial w^1} = \delta^1 x$$

Summary: the equations of backpropagation

$$\delta^L = \nabla_a C \odot \sigma'(z^L) \quad (\text{BP1})$$

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l) \quad (\text{BP2})$$

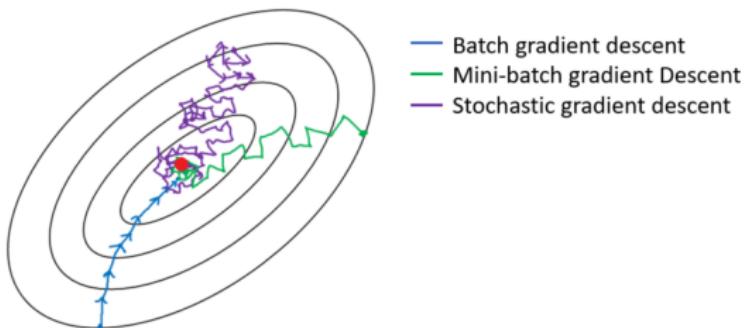
$$\frac{\partial C}{\partial b_j^l} = \delta_j^l \quad (\text{BP3})$$

$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l \quad (\text{BP4})$$

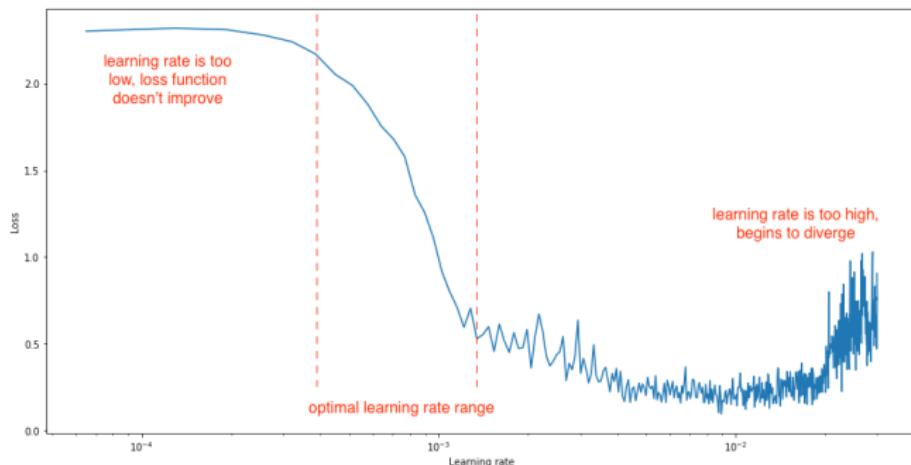
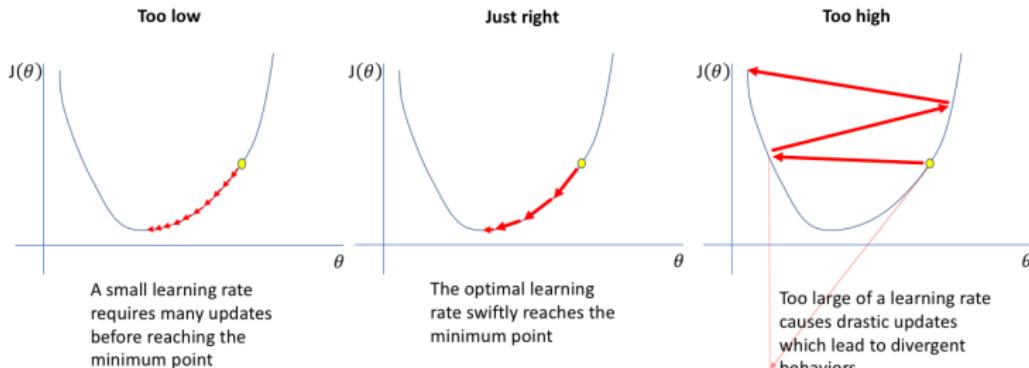
Source: Neural Networks and Deep Learning

Gradient Descent (GD)

- ▶ Batch GD $J = N$ (all observations): high computation cost when N is large; can be trapped in local minimum
- ▶ Stochastic GD $J = 1$ (randomly choose one obs): noisy and continue to wander around minimum (slow convergence)
- ▶ Mini-batch GD $J \approx [2^5, 2^9]$ (batch sampling): most popular



Source and reference: Adventures in Machine Learning; Andrew Ng; Imad Dabbura



Gradient Descent with Momentum

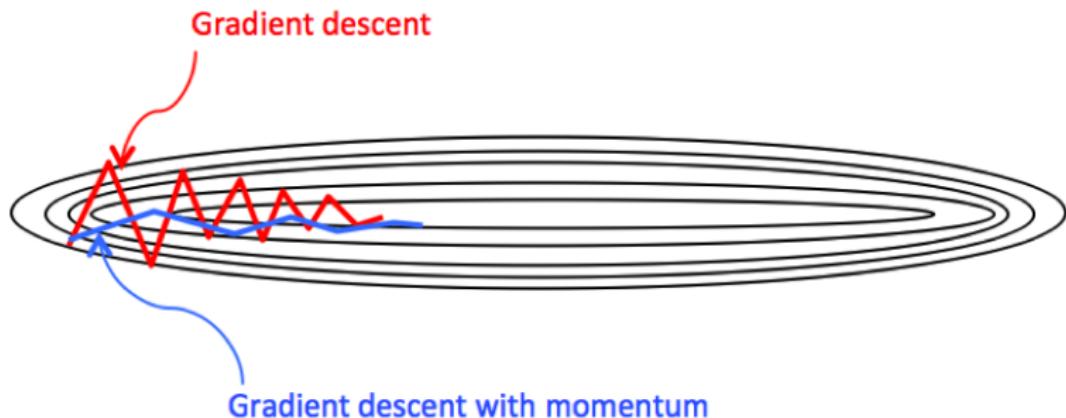
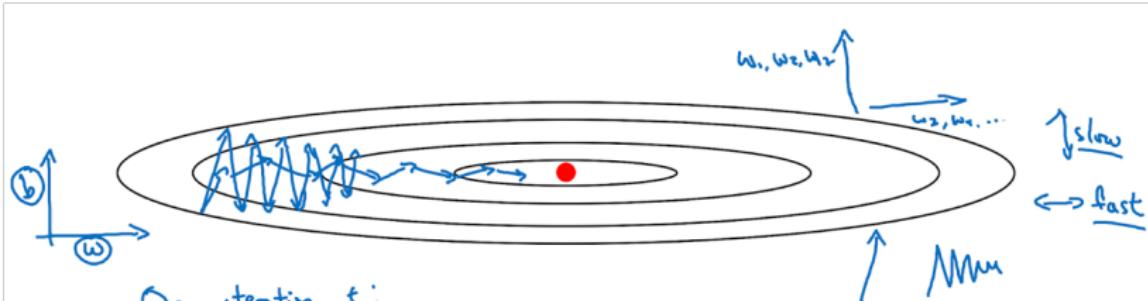


Figure: Update by the exponential moving average of gradients

- ▶ $\mathbf{m} \leftarrow \beta\mathbf{m} + (1 - \beta)\nabla_{\Theta} C$
- ▶ $\Theta \leftarrow \Theta - \epsilon\mathbf{m}$
- ▶ Cancel out the gradient in directions where it oscillates; speed up in directions where it accumulates

Gradient Descent with RMSProp



On iteration t :

Compute dW, db on current mini-batch

$$S_{dW} = \beta_2 S_{dW} + (1-\beta_2) dW^2 \quad \text{element-wise} \quad \leftarrow \text{small}$$

$$\rightarrow S_{db} = \beta_2 S_{db} + (1-\beta_2) d b^2 \quad \leftarrow \text{large}$$

$$w := w - \frac{\alpha}{\sqrt{S_{dW} + \epsilon}} dW \quad \leftarrow \quad b := b - \frac{\alpha}{\sqrt{S_{db} + \epsilon}} d b \quad \leftarrow$$

$$\epsilon = 10^{-8}$$

Figure: Adaptive Learning Rate

Other Considerations

- ▶ Weight Initialization (LeCun, Glorot/Xavier, He Initialization)
- ▶ Optimization: Adagrad, RMSprop, Adam, Nesterov Accelerated Gradient, Nadam ...
- ▶ Batch normalization (initializing inputs with zero mean and unit variance)
- ▶ Gradient Clipping
- ▶ Learning rate scheduling
- ▶ Regularization: l_1 and l_2 Norm, Dropout/Monte Carlo Dropout, Early Stopping
- ▶ Bagged and Boosted NN, Bayesian NN

See HOML Chapter 11 for reference.

Outline

Overview of (Feedforward) Neural Networks

Structure of Neural Networks

Components of Neural Networks

Train Neural Networks

Neural Networks In Action

Build a DNN from Scratch

Software Packages

Limitation of Neural Networks

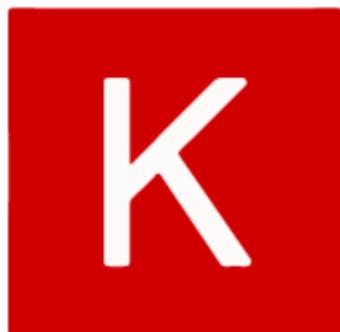
Hyper-parameter Tuning

Build and Train A Neural Network From Scratch (link)

```
94 # Transfer neuron activation
95 def transfer(activation):
96     return 1.0 / (1.0 + exp(-activation))
97
98 # Forward propagate input to a network output
99 def forward_propagate(network, row):
100    inputs = row
101    for layer in network:
102        new_inputs = []
103        for neuron in layer:
104            activation = activate(neuron['weights'], inputs)
105            neuron['output'] = transfer(activation)
106            new_inputs.append(neuron['output'])
107        inputs = new_inputs
108    return inputs
109
110 # Calculate the derivative of an neuron output
111 def transfer_derivative(output):
112     return output * (1.0 - output)
113
114 # Backpropagate error and store in neurons
115 def backward_propagate_error(network, expected):
116     for i in reversed(range(len(network))):
117         layer = network[i]
118         errors = []
119         if i != len(network)-1:
120             for j in range(len(layer)):
121                 error = 0.0
122                 for neuron in network[i+1]:
123                     error += (neuron['weights'][j] * neuron['delta'])
124                 errors.append(error)
125         else:
126             for j in range(len(layer)):
127                 neuron = layer[j]
128                 errors.append(expected[j] - neuron['output'])
129         for j in range(len(layer)):
130             neuron = layer[j]
131             neuron['delta'] = errors[j] * transfer_derivative(neuron['output'])
132
133 # Update network weights with error
134 def update_weights(network, row, l_rate):
135     for i in range(len(network)):
136         inputs = row[:-1]
137         if i != 0:
138             inputs = [neuron['output'] for neuron in network[i-1]]
139         for neuron in network[i]:
140             for j in range(len(inputs)):
141                 neuron['weights'][j] += l_rate * neuron['delta'] * inputs[j]
142                 neuron['weights'][-1] -= l_rate * neuron['delta']
143
144 # Train a network for a fixed number of epochs
145 def train_network(network, train, l_rate, n_epoch, n_outputs):
146     for epoch in range(n_epoch):
147         for row in train:
148             outputs = forward_propagate(network, row)
149             expected = [0 for i in range(n_outputs)]
150             expected[row[-1]] = 1
151             backward_propagate_error(network, expected)
152             update_weights(network, row, l_rate)
```

Machine Learning Packages

Keras



TensorFlow



PyTorch

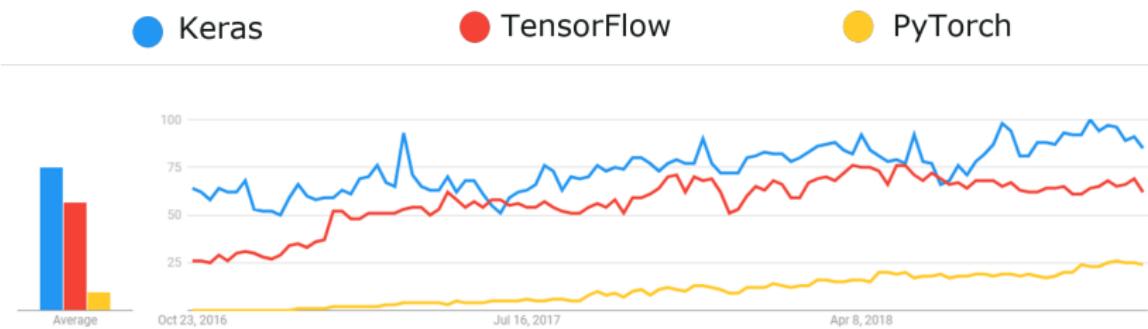


- ▶ High-level API Standard
- ▶ Very easy to learn

- ▶ Use Keras as its core API
- ▶ Easy to deploy into production

- ▶ More Pythonic
- ▶ Popular in research community

Popularity of Keras, TensorFlow and PyTorch



Source: edureka!

Keras is very easy to use!

```
import numpy as np
from keras.models import Sequential
from keras.layers import Dense, Dropout

# Generate dummy data
x_train = np.random.random((1000, 20))
y_train = np.random.randint(2, size=(1000, 1))
x_test = np.random.random((100, 20))
y_test = np.random.randint(2, size=(100, 1))

model = Sequential()
model.add(Dense(64, input_dim=20, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(64, activation='relu'))
model.add(Dropout(0.5))
model.add(Dense(1, activation='sigmoid'))

model.compile(loss='binary_crossentropy',
              optimizer='rmsprop',
              metrics=['accuracy'])

model.fit(x_train, y_train,
          epochs=20,
          batch_size=128)
score = model.evaluate(x_test, y_test, batch_size=128)
```

Figure: MLP for binary classification

If NN is so easy and so powerful, then ...

- ▶ Why don't we see it used everywhere?
- ▶ Why is gradient boosting more often the winner of Kaggle contests?

Neural Network Is Not An "Off-the-Shelf" Procedure

- ▶ Need to encode (or embed) categorical variables
- ▶ Not scale invariant and sensitive to outliers
- ▶ Too many hyper-parameters to tune (architecture, optimizer, initialization, regularization ...).
- ▶ Need domain specific priors (and many years of research) to design a good network

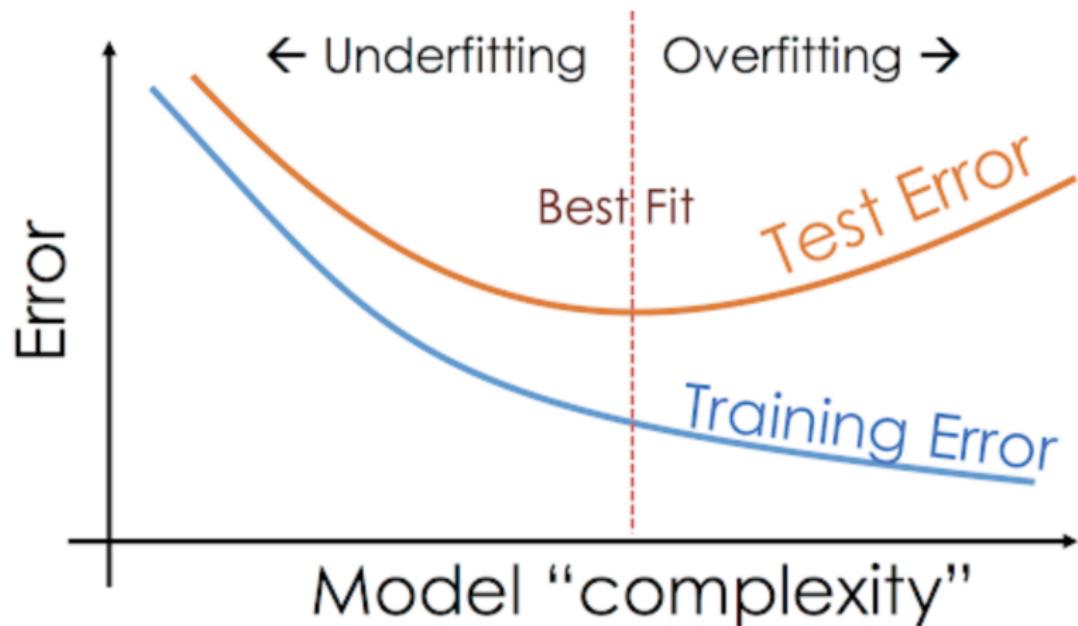
Therefore, for tasks with short timelines and no prior filed knowledge (like Kaggle), gradient boosting (and other tree-based methods) usually wins.

However, for a tasks with

- ▶ very high complexity and sufficient data (computer vision, NLP, reinforcement learning ...)
 - ▶ the right architecture (convolution NN, recurrent NN ...)
- NN (deep learning) blows any other method out of the water



How to Tune Hyper-parameters



Minimize Test Error Given Data and Resources

$$\text{Test Error} = \text{Training Error} + \text{Generalization Gap}$$

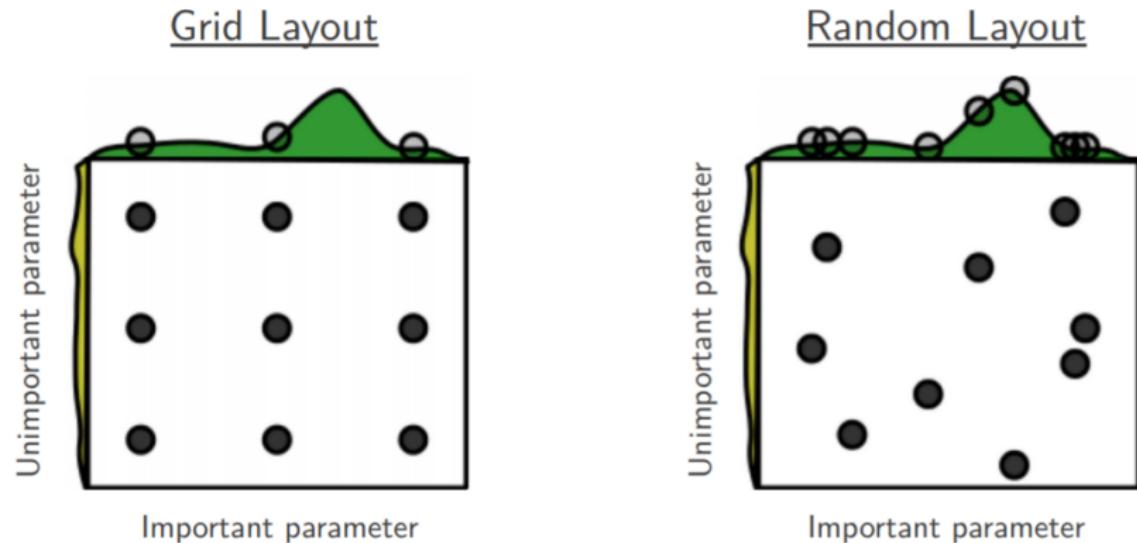
Stretch-pants and Orthogonalization

- ▶ "Stretch-pants": Increase model capacity to match the complexity of the task. First try to decrease training error and then use regularization to decrease generalization gap
- ▶ Orthogonal: focus on one thing at a time – training error or generalization

Some Best Practices in Hyper-parameter Tuning

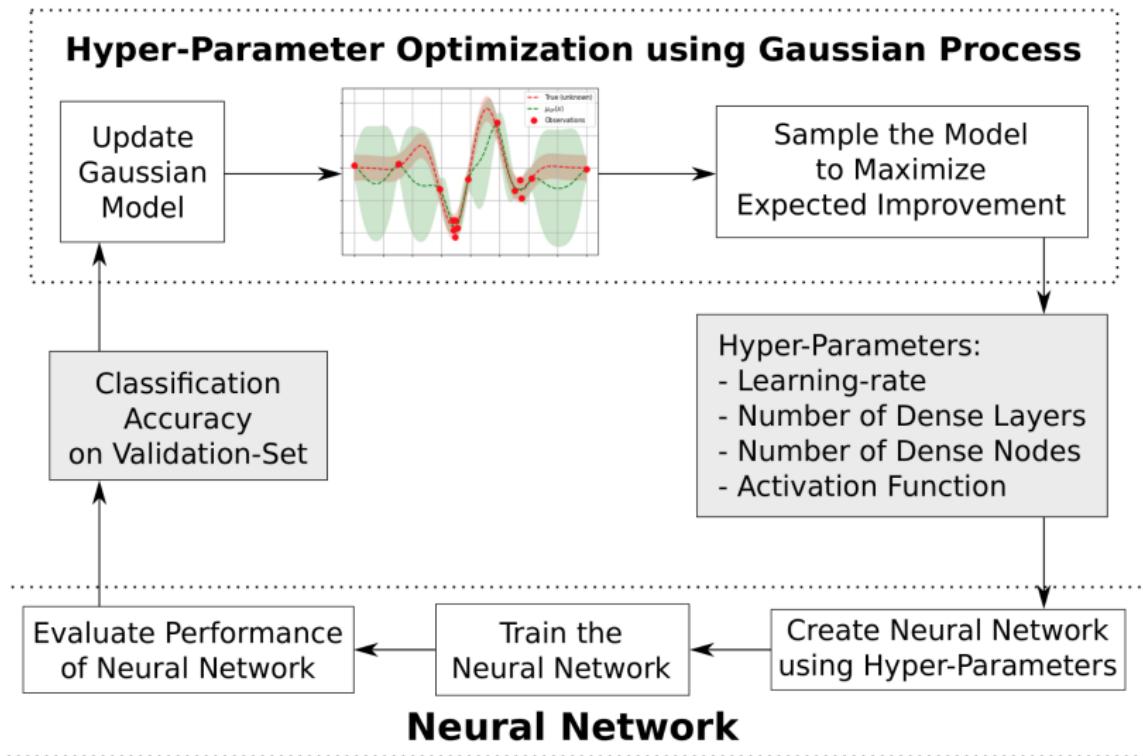
- ▶ Always starts with EDA (how many types of data, missing data, irrelevant data ...)
- ▶ Debug using a bare-bone model (Is the cost function defined correctly? Does the input tensor have the right dimension? Does the prediction make sense?)
- ▶ Trial-and-error and get your hands dirty. (Don't tune it yet; try some configuration and see if you have any insights)
- ▶ Define a default baseline model (It doesn't have to be neural network! Try boosting!)
- ▶ Can you get more data? (Plot the error against training size on a log scale)

Hands-off Automatic Tuning: Grid and Random Search



- ▶ Most of the time, random search is better if you have more than one hyper-parameter to tune.

Hands-off Automatic Tuning: Bayesian Optimization



- ▶ This can be very time consuming, and not necessarily better.