

Neural Networks Review

Arin Ghazarian
Chapman University

Kullback-Leibler (KL) divergence

- If we have two separate probability distributions $P(x)$ and $Q(x)$ over the same random variable x , we can measure how different these two distributions are using the Kullback-Leibler (KL) divergence:

$$D_{\text{KL}}(P\|Q) = \mathbb{E}_{x \sim P} \left[\log \frac{P(x)}{Q(x)} \right] = \mathbb{E}_{x \sim P} [\log P(x) - \log Q(x)].$$

Kullback-Leibler (KL) divergence

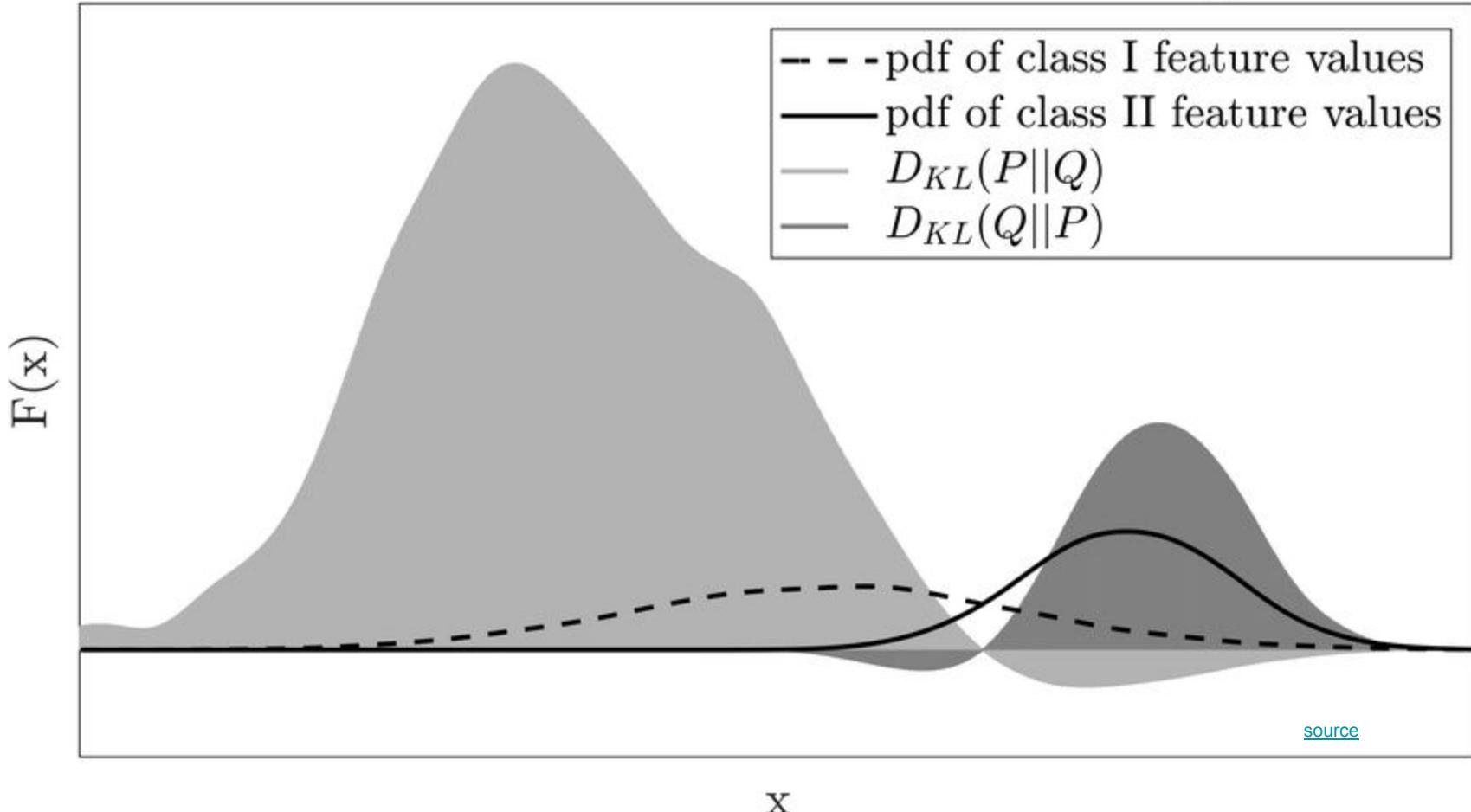
- In the case of discrete variables, it is the extra amount of information (measured in bits if we use the base 2 logarithm, but in machine learning we usually use nats and the natural logarithm) needed to send a message containing symbols drawn from probability distribution P, when we use a code that was designed to minimize the length of messages drawn from probability distribution Q.

Kullback-Leibler (KL) divergence as a Measure of Distance

- The KL divergence has many useful properties, most notably that it is nonnegative.
- The KL divergence is 0 if and only if P and Q are the same distribution in the case of discrete variables, or equal “almost everywhere” in the case of continuous variables.
- Because the KL divergence is non-negative and measures the difference between two distributions, it is often conceptualized as measuring some sort of distance between these distributions.
- However, it is not a true distance measure because it is not symmetric:

$$D_{\text{KL}}(P\|Q) \neq D_{\text{KL}}(Q\|P) \text{ for some } P \text{ and } Q$$

Visualization of the Kullback–Leibler divergence



Cross-entropy

- A quantity that is closely related to the KL divergence is the cross-entropy $H(P, Q) = H(P) + D_{KL}(PQ)$, which is similar to the KL divergence but lacking the term on the left:

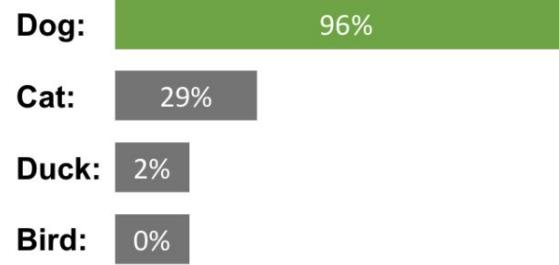
$$H(P, Q) = -\mathbb{E}_{x \sim P} \log Q(x).$$

- Minimizing the cross-entropy with respect to Q is equivalent to minimizing the KL divergence, because Q does not participate in the omitted term.

Likelihood Example

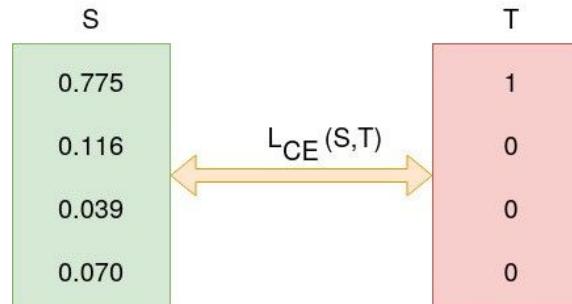
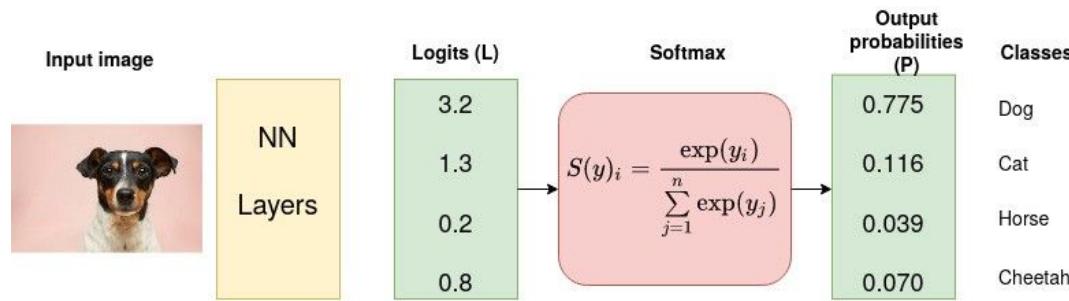


[source](#)



$$L = .96^1 \times .29^0 \times .02^0 \times .0^0 = 0.96$$

Cross Entropy Data Distribution vs Model Distribution



LogLikelihood

$$\text{loglik}(\theta) = \sum_{j=1}^M y_j \log(\bar{y}_j)$$

It is easier to minimise the negative log-likelihood function than maximising the direct likelihood

$$\text{loglik}(\theta) = - \sum_{j=1}^M y_j \log(\bar{y}_j)$$

Cross-entropy and Loglikelihood

- For a multiclass classification problem, we use cross-entropy as a loss function.
 - The average bits of information required to identify an event drawn from the estimated probability distribution $f(x)$, rather than the true distribution $g(x)$

$$H(P, Q) = -\mathbb{E}_{x \sim P} \log Q(x).$$

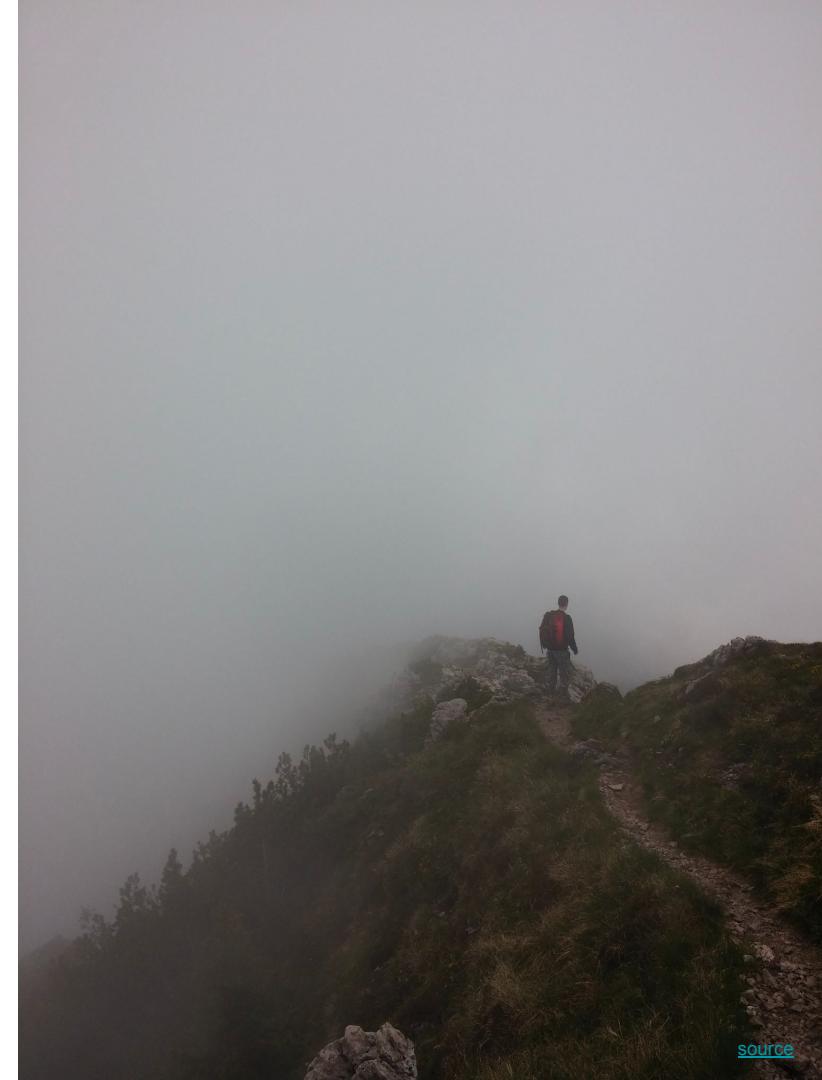
$$P=y_i \text{ and } Q=\bar{y}_i$$

$$H(P, Q) = -\sum_{j=1}^M y_j \log(\bar{y}_j)$$

Thus cross-entropy is equal to negative log-likelihood. By minimizing the negative loglikelihood, in fact we are maximizing the likelihood

Gradient descent

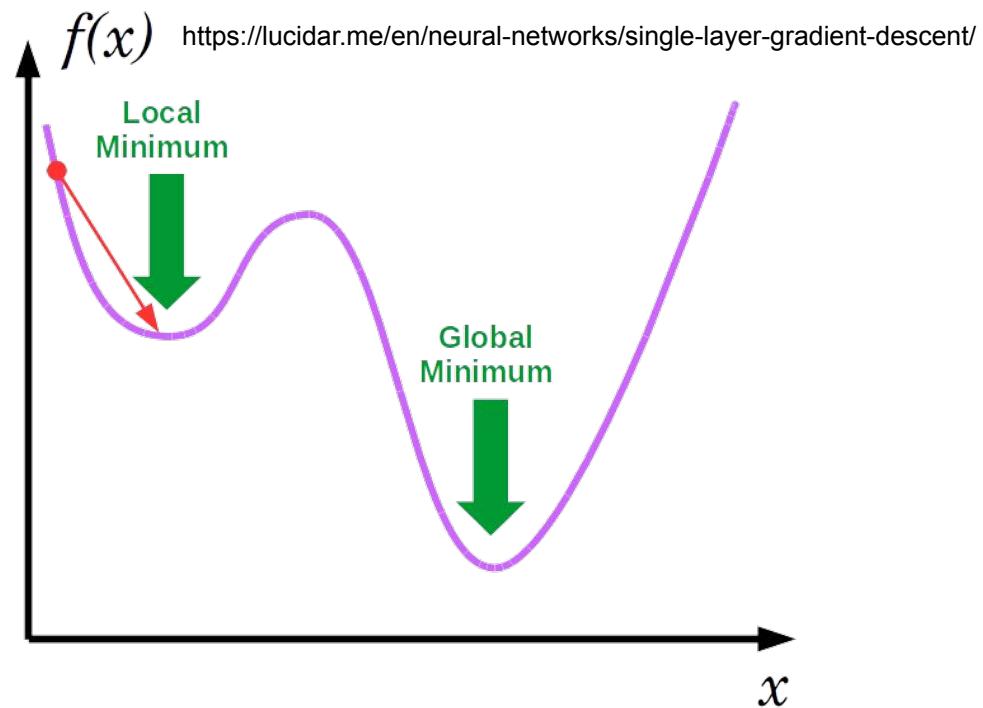
- Hiker in the fog approach
- Imagine you are hiking in the fog and need to get down a mountain.
- You cannot see where the optimal downhill path is.
- But you can see your immediate surroundings.
- So you decide to choose the steepest downhill path at any point.
- If you are lucky, this method will get you down the hill. If not, you might get stuck in some valley or ravine, still up the mountain.



[source](#)

Gradient descent

- Taking the best possible downhill path at any point is the essence of the gradient descent algorithm.



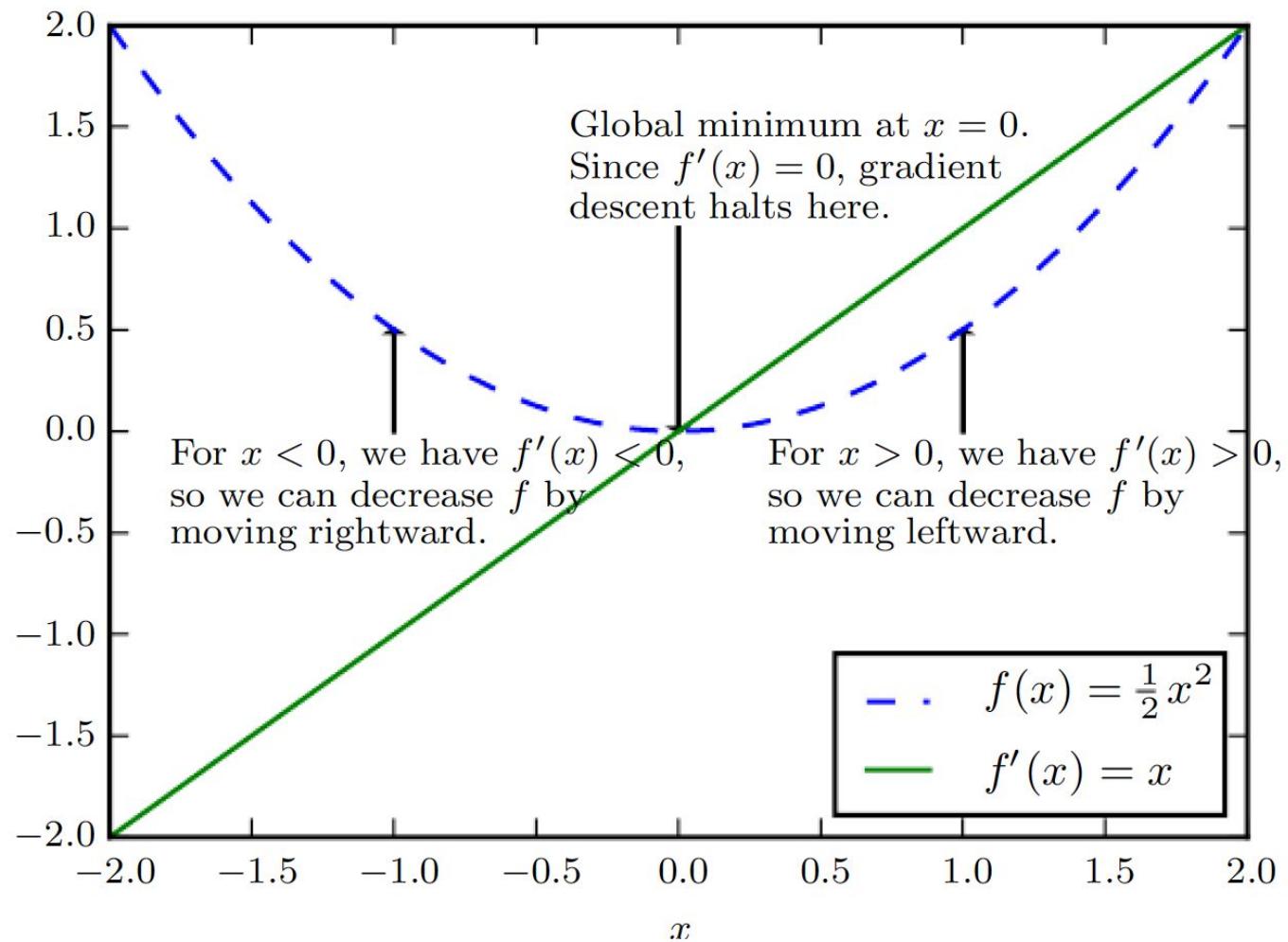
Derivative

- Using the derivative we can estimate how a small change in input will scale the output:

$$f(x + \varepsilon) \approx f(x) + \varepsilon f'(x)$$

Gradient Descent

- In gradient descent, we move in small steps **in the opposite to the derivative sign**, hoping it will reduce the error function.
- The derivative is therefore useful for minimizing a function because it tells us how to change x in order to make a small improvement in y .
- For example, we know that $f(x - \varepsilon \operatorname{sign}(f'(x)))$ is less than $f(x)$ for small enough ε . We can thus reduce $f(x)$ by moving x in small steps with opposite sign of the derivative. This technique is called gradient descent (Cauchy, 1847).



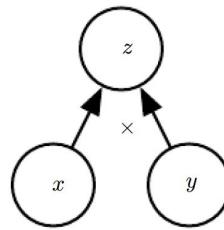
Gradient Descent

- Gradient descent proposes a new point, where epsilon is the learning rate, a positive scalar determining the size of the step:

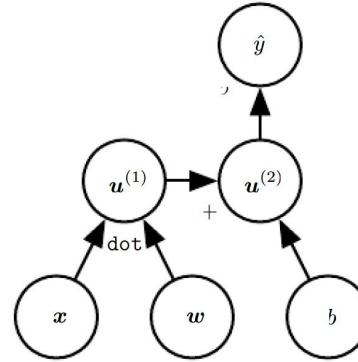
$$\mathbf{x}' = \mathbf{x} - \epsilon \nabla_{\mathbf{x}} f(\mathbf{x})$$

- Gradient descent converges when every element of the gradient is zero (or, in practice, very close to zero). In some cases, we may be able to avoid running this iterative algorithm, and just jump directly to the critical point by solving the equation $\nabla_{\mathbf{x}} f(\mathbf{x}) = 0$ for \mathbf{x} .

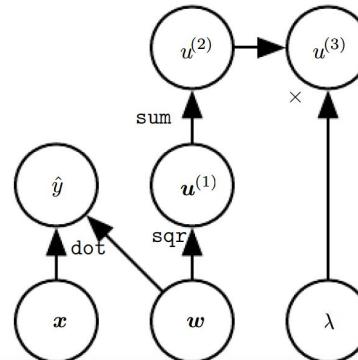
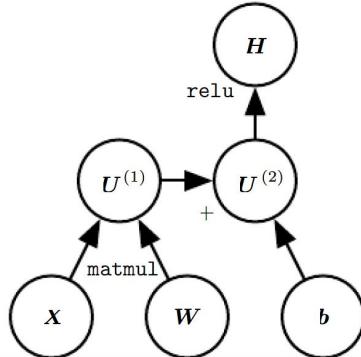
Computational Graphs



(a)

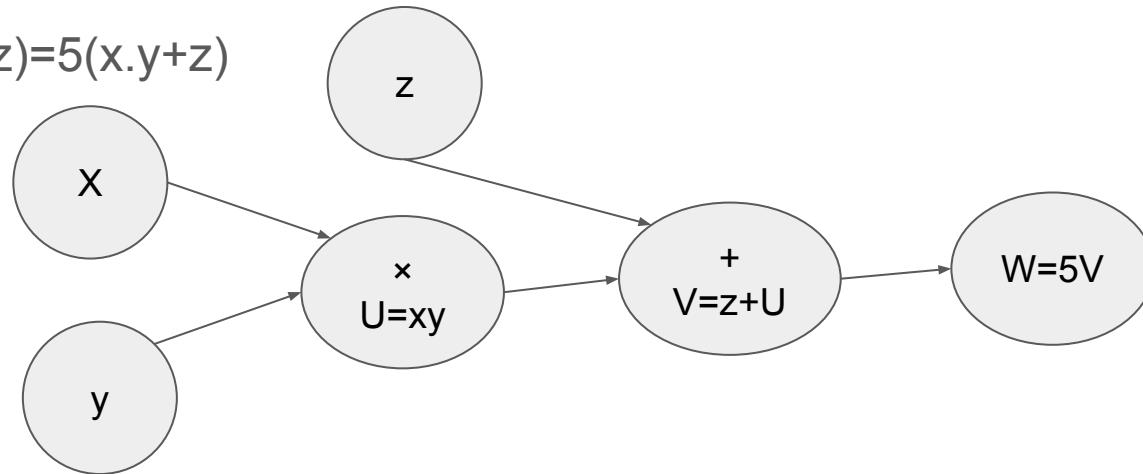


(b)



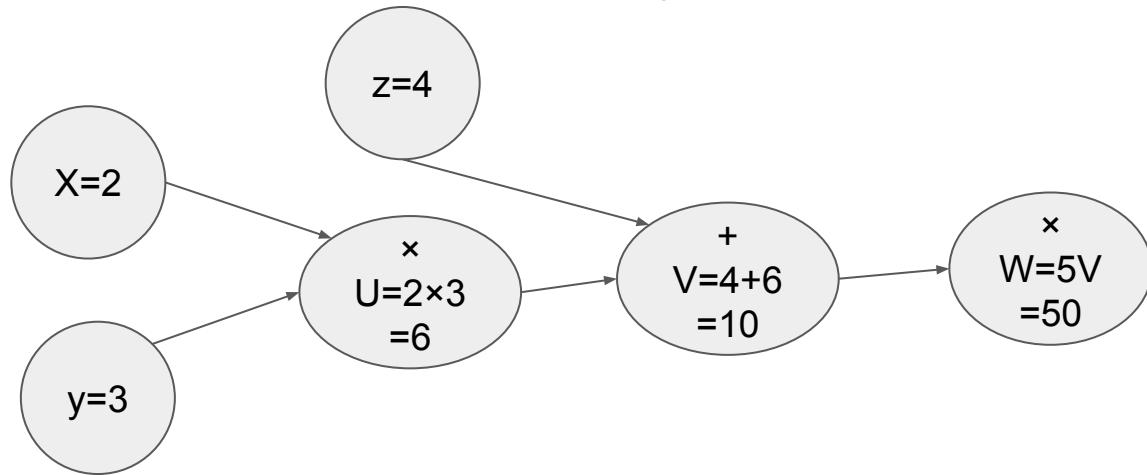
Computational Graphs

$$f(x,y,z) = 5(x \cdot y + z)$$



Computational Graphs:

Forward Propagation

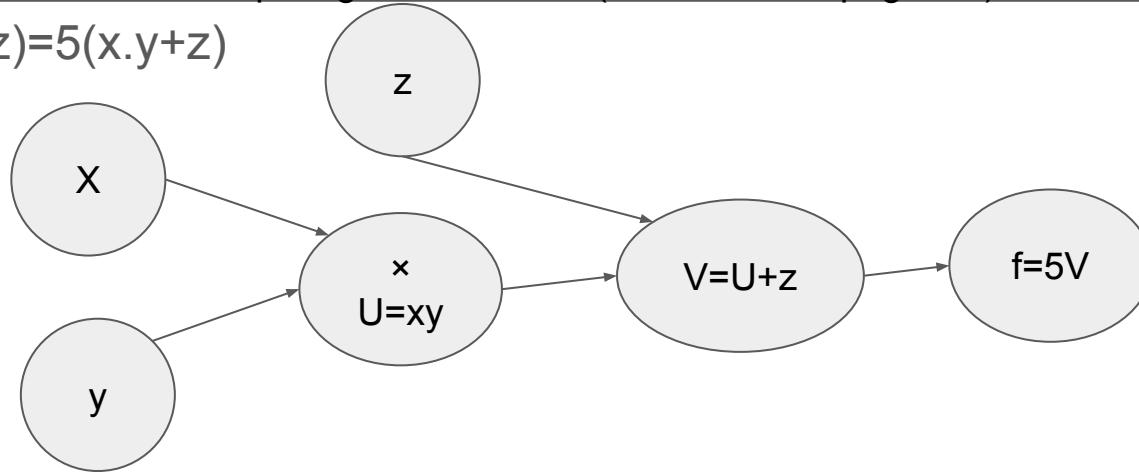


$$x=2, y=3, z=4$$

Computational Graphs:

Computing the Derivative(Backward Propagation)

$$f(x,y,z) = 5(x \cdot y + z)$$



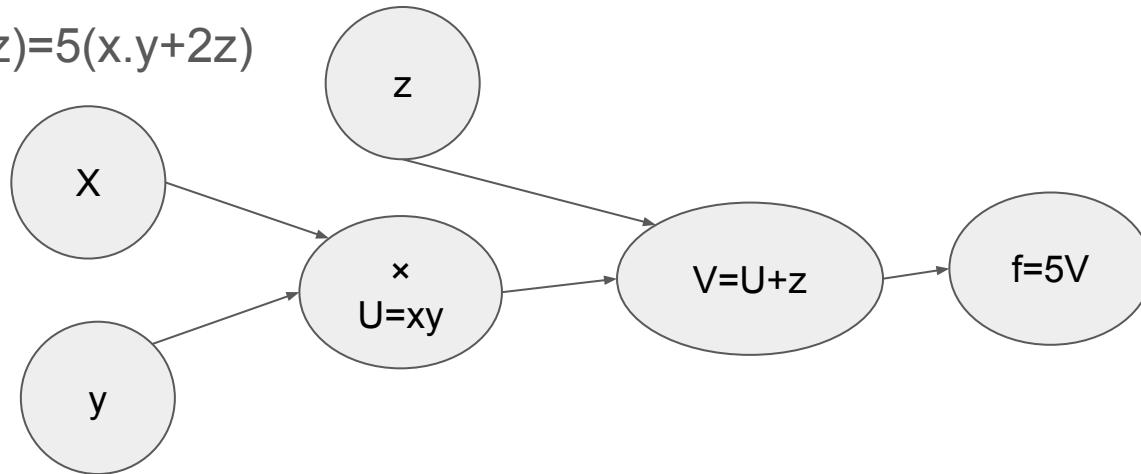
$$\partial f / \partial V = 5$$

$$\partial V / \partial z = 1$$

$$\partial f / \partial z = \partial f / \partial V \cdot \partial V / \partial z = 5 \cdot 1 = 5$$

Computational Graphs: Computing the Derivative

$$f(x,y,z) = 5(x.y + 2z)$$



$$\begin{aligned}x &= 2, y = 3 \\z &= 4 \rightarrow 4.01 \\V &= 10 \rightarrow 10.01 \\f &= 50 \rightarrow 50.05\end{aligned}$$

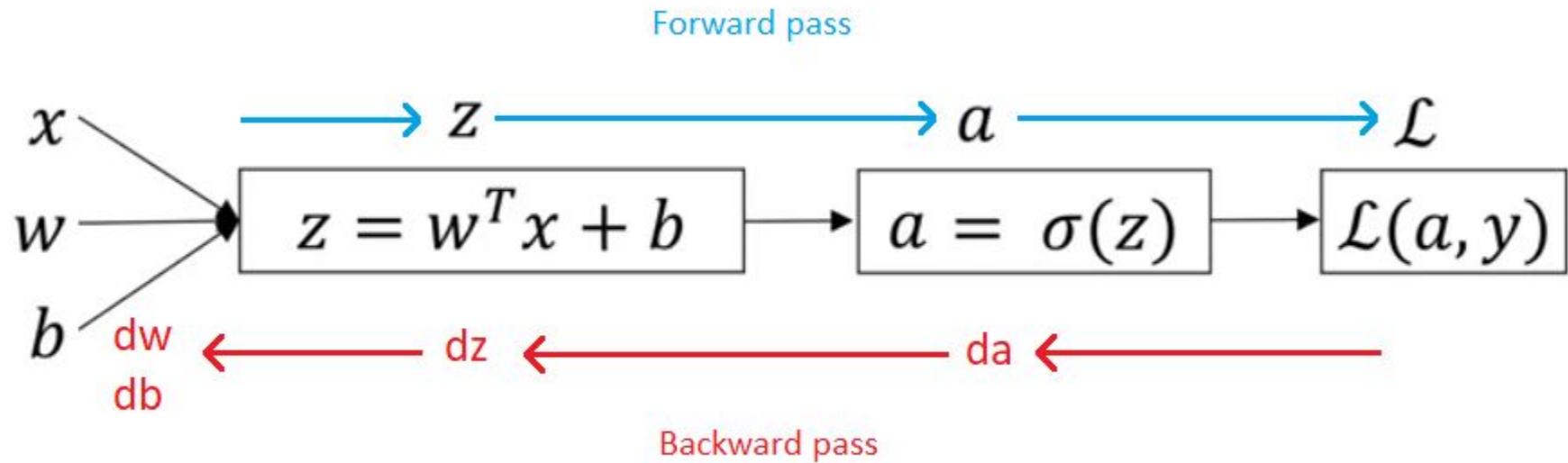
$$\frac{\partial f}{\partial V} = 5$$

$$\frac{\partial V}{\partial z} = 1$$

$$\frac{\partial f}{\partial z} = \frac{\partial f}{\partial V} \cdot \frac{\partial V}{\partial z} = 5 \cdot 2 = 10$$

Computational Graph For Logistic Regression

da means $\partial L / \partial a$ and dz means $\partial L / \partial z$



Derivation Rules: Reminder

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

$$(\ln x)' = \frac{1}{x}$$

$$(f(x) g(x))' = f'(x) g(x) + f(x) g'(x)$$

$$\sigma'(x) = \frac{d}{dx} \sigma(x) = \sigma(x)(1 - \sigma(x))$$

$$\frac{d}{dx} [\ln u] = \frac{1}{u} \cdot u' = \frac{u'}{u}$$

$$L(a,y) = -[y \log(a) + (1-y) \log(1-a)]$$

$$\partial L / \partial a = da = -[(y \cdot (1/a)) + ((1-y) \cdot (-1/(1-a)))] = -[((1-a)y + a(y-1))/(a(1-a))] = -(y-a)/(a(1-a))$$

$$\frac{dL}{dz} = \frac{dL}{da} \times \frac{da}{dz}$$

$$\frac{dL}{dz} = \frac{a-y}{a(1-a)} \times a(1-a) = a - y$$

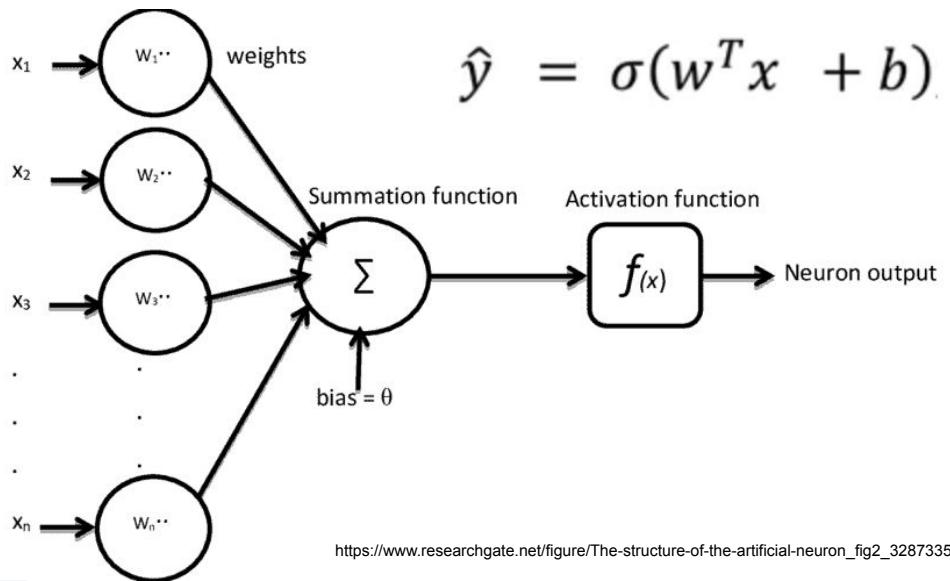
$$\partial L / \partial w = \partial L / \partial z. \quad \partial z / \partial w = (a-y)x$$

$$\partial L / \partial b = \partial L / \partial z. \quad \partial z / \partial b = (a-y)$$

Logistic Regression as a Single Neuron Neural Network

- Logistic regression is in fact a single neuron neural network model with a sigmoid activation function: $P(X) = \sigma(\beta_0 + \beta_1 x)$
- The cost function is the cross entropy between the predictions and actuals
- We used gradient descent to find the W and b which minimizes the cost

$$P = \frac{1}{1+e^{-(\beta_0 + \beta_1 X_1 + \beta_2 X_2 + \beta_3 X_3 + \dots + \beta_n X_n)}}$$



Vectorizing Logistic Regression

- Assuming we have M training samples, we need to go over each sample in a loop

$$J = 0, dw_1 = 0, dw_2 = 0, db = 0$$

for i = 1 to m:

$$z^{(i)} = w^T x^{(i)} + b$$

$$a^{(i)} = \sigma(z^{(i)})$$

$$J += -[y^{(i)} \log a^{(i)} + (1 - y^{(i)}) \log(1 - a^{(i)})]$$

$$dz^{(i)} = a^{(i)} - y^{(i)}$$

$$dw_1 += x_1^{(i)} dz^{(i)}$$

$$dw_2 += x_2^{(i)} dz^{(i)}$$

$$db += dz^{(i)}$$

$$J = J/m, dw_1 = dw_1/m, dw_2 = dw_2/m$$

$$\underline{db} = db/m$$

Vectorizing Logistic Regression:Forward Propagation

- The size of training matrix X will be (N_x by M).
 - N_x is the number of features
- W is (N_x by 1) matrix (W of the logistic regression model)

$$\begin{aligned}Z &= [z^{(1)} \ z^{(2)} \ z^{(3)} \dots \ z^{(m)}] \\&= w^T X^T + [b \ b \ b \dots b] \\&= w^T \left[\begin{array}{c|c|c|c|c} \vdots & \vdots & \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \end{array} \right] + [b \ b \ \dots \ b] \\&= \left[\begin{array}{c|c|c|c} \overbrace{w^T x^{(1)} + b}^{z^{(1)}} & \overbrace{w^T x^{(2)} + b}^{z^{(2)}} & \dots & \overbrace{w^T x^{(m)} + b}^{z^{(m)}} \end{array} \right] \quad \begin{matrix} \rightarrow (1 \times m) \\ \text{dimensions} \end{matrix} \\&\quad \begin{matrix} \rightarrow (1 \times m) \\ \text{dimensions} \end{matrix}\end{aligned}$$

$$A = [a^{(1)} \ a^{(2)} \ a^{(3)} \dots \ a^{(m)}] = \sigma(Z)$$

$$X = \begin{bmatrix} \vdots & \vdots & \vdots & \vdots & \vdots \\ x_1 & x_2 & x_3 & x_4 & \dots & x_m \end{bmatrix} \quad \begin{matrix} (n_x, m) \\ \mathbb{R}^{n_x \times m} \end{matrix}$$

Vectorizing Logistic Regression: Backward Propagation

$$db = \frac{1}{m} \sum_{i=1}^m dz^{(i)}$$
$$= \frac{1}{m} \text{np.sum}(dz)$$

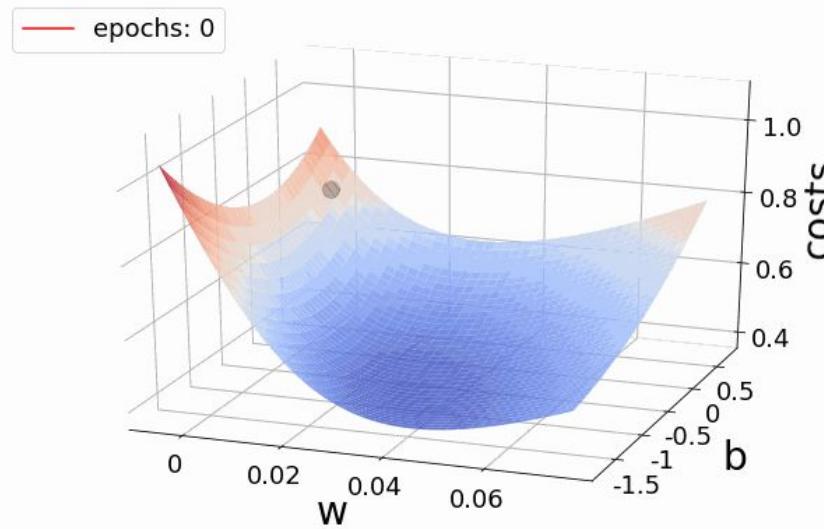
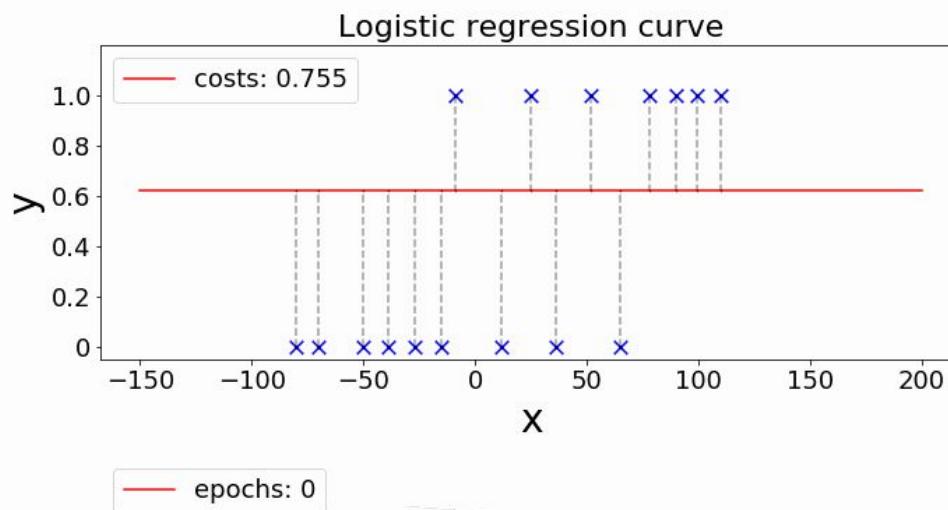
$$dw = \frac{1}{m} X dz^\top$$

$$= \frac{1}{m} \begin{bmatrix} \vdots & \vdots & \vdots & \vdots \\ x^{(1)} & x^{(2)} & x^{(3)} & \dots & x^{(m)} \\ \vdots & \vdots & \vdots & \vdots & \vdots \\ (n \times m) \end{bmatrix} \begin{bmatrix} dz^{(1)} \\ dz^{(2)} \\ \vdots \\ dz^{(m)} \\ (m \times 1) \end{bmatrix}$$

$$= \frac{1}{m} \left[x^{(1)} dz^{(1)} + x^{(2)} dz^{(2)} + \dots + x^{(m)} dz^{(m)} \right]_{(n \times 1)}$$

$$\omega := \omega - \alpha dw$$

$$b := b - \alpha db$$



<https://towardsdatascience.com/animations-of-logistic-regression-with-python-31f8c9cb420>

Chain Rule of Calculus

- Let x be a real number, and let f and g both be functions mapping from a real number to a real number. Suppose that $y = g(x)$ and $z = f(g(x)) = f(y)$. Then the chain rule states that

$$\frac{dz}{dx} = \frac{dz}{dy} \frac{dy}{dx}$$

- We can generalize this beyond the scalar case. Suppose that $x \in \mathbb{R}^m$, $y \in \mathbb{R}^n$, g maps from \mathbb{R}^m to \mathbb{R}^n , and f maps from \mathbb{R}^n to \mathbb{R} . If $y = g(x)$ and $z = f(y)$, then

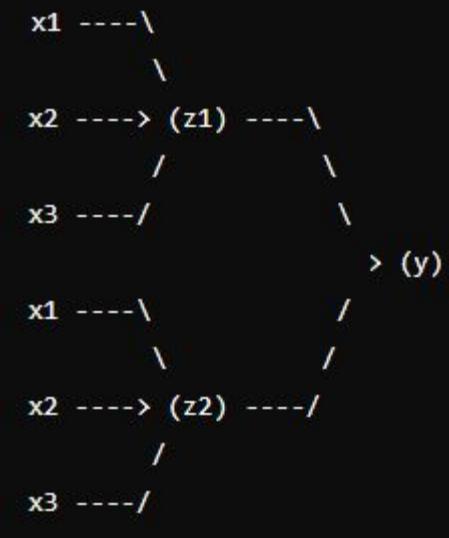
$$\frac{\partial z}{\partial x_i} = \sum_j \frac{\partial z}{\partial y_j} \frac{\partial y_j}{\partial x_i}$$

- In vector notation, this may be equivalently written as where $\partial y / \partial x$ is the $n \times m$ Jacobian matrix of g .

$$\nabla_x z = \left(\frac{\partial y}{\partial x} \right)^\top \nabla_y z$$

Jacobian-gradient Products

- The gradient of a variable x can be obtained by multiplying a Jacobian matrix $\partial y / \partial x$ by a gradient $\nabla_y z$. The back-propagation algorithm consists of performing such a Jacobian-gradient product for each operation in the graph.
- <https://suzyahyah.github.io/calculus/machine%20learning/2018/04/04/Jacobian-and-Backpropagation.html>
- Example: Given a fully connected neural net with 3 inputs, x_1, x_2, x_3 and hidden layer with two neurons z_1, z_2 , and a single scalar output y :



$$\frac{\delta y}{\delta x} = \frac{\delta y}{\delta z} \times \frac{\delta z}{\delta x} = \begin{bmatrix} \frac{\delta y}{\delta z_1} & \frac{\delta y}{\delta z_2} \end{bmatrix} \times \begin{bmatrix} \frac{\delta z_1}{\delta x_1} & \frac{\delta z_1}{\delta x_2} & \frac{\delta z_1}{\delta x_3} \\ \frac{\delta z_2}{\delta x_1} & \frac{\delta z_2}{\delta x_2} & \frac{\delta z_2}{\delta x_3} \end{bmatrix}$$

Diagram illustrating the Jacobian-gradient product. A 'Gradient' box points to the vector $\begin{bmatrix} \frac{\delta y}{\delta z_1} & \frac{\delta y}{\delta z_2} \end{bmatrix}$. A 'Jacobian' box points to the matrix $\begin{bmatrix} \frac{\delta z_1}{\delta x_1} & \frac{\delta z_1}{\delta x_2} & \frac{\delta z_1}{\delta x_3} \\ \frac{\delta z_2}{\delta x_1} & \frac{\delta z_2}{\delta x_2} & \frac{\delta z_2}{\delta x_3} \end{bmatrix}$.

Demo: Training a Logistic Regression Using Gradient Descent

- <http://localhost:8888/notebooks/jupyter-notebooks/chapman-generative-AI/NN-review-Logistic-regression.ipynb>

Automatic Differentiation

- Automatic differentiation or autodiff is a set of techniques to evaluate the derivative of a function specified by a computer program.
- AD exploits the fact that every computer program, no matter how complicated, executes a sequence of elementary arithmetic operations (addition, subtraction, multiplication, division, etc.) and elementary functions (exp, log, sin, cos, etc.). By applying the chain rule repeatedly to these operations, derivatives of arbitrary order can be computed automatically, accurately to working precision, and using at most a small constant factor more arithmetic operations than the original program.

Algorithm 6.3 Forward propagation through a typical deep neural network and the computation of the cost function. The loss $L(\hat{\mathbf{y}}, \mathbf{y})$ depends on the output $\hat{\mathbf{y}}$ and on the target \mathbf{y} (see section 6.2.1.1 for examples of loss functions). To obtain the total cost J , the loss may be added to a regularizer $\Omega(\theta)$, where θ contains all the parameters (weights and biases). Algorithm 6.4 shows how to compute gradients of J with respect to parameters \mathbf{W} and \mathbf{b} . For simplicity, this demonstration uses only a single input example \mathbf{x} . Practical applications should use a minibatch. See section 6.5.7 for a more realistic demonstration.

Require: Network depth, l

Require: $\mathbf{W}^{(i)}, i \in \{1, \dots, l\}$, the weight matrices of the model

Require: $\mathbf{b}^{(i)}, i \in \{1, \dots, l\}$, the bias parameters of the model

Require: \mathbf{x} , the input to process

Require: \mathbf{y} , the target output

$$\mathbf{h}^{(0)} = \mathbf{x}$$

for $k = 1, \dots, l$ **do**

$$\mathbf{a}^{(k)} = \mathbf{b}^{(k)} + \mathbf{W}^{(k)} \mathbf{h}^{(k-1)}$$

$$\mathbf{h}^{(k)} = f(\mathbf{a}^{(k)})$$

end for

$$\hat{\mathbf{y}} = \mathbf{h}^{(l)}$$

$$J = L(\hat{\mathbf{y}}, \mathbf{y}) + \lambda \Omega(\theta)$$

Algorithm 6.4 Backward computation for the deep neural network of algorithm 6.3, which uses in addition to the input \mathbf{x} a target \mathbf{y} . This computation yields the gradients on the activations $\mathbf{a}^{(k)}$ for each layer k , starting from the output layer and going backwards to the first hidden layer. From these gradients, which can be interpreted as an indication of how each layer's output should change to reduce error, one can obtain the gradient on the parameters of each layer. The gradients on weights and biases can be immediately used as part of a stochastic gradient update (performing the update right after the gradients have been computed) or used with other gradient-based optimization methods.

After the forward computation, compute the gradient on the output layer:

$$\mathbf{g} \leftarrow \nabla_{\hat{\mathbf{y}}} J = \nabla_{\hat{\mathbf{y}}} L(\hat{\mathbf{y}}, \mathbf{y})$$

for $k = l, l - 1, \dots, 1$ **do**

Convert the gradient on the layer's output into a gradient into the pre-nonlinearity activation (element-wise multiplication if f is element-wise):

$$\mathbf{g} \leftarrow \nabla_{\mathbf{a}^{(k)}} J = \mathbf{g} \odot f'(\mathbf{a}^{(k)})$$

Compute gradients on weights and biases (including the regularization term, where needed):

$$\nabla_{\mathbf{b}^{(k)}} J = \mathbf{g} + \lambda \nabla_{\mathbf{b}^{(k)}} \Omega(\theta)$$

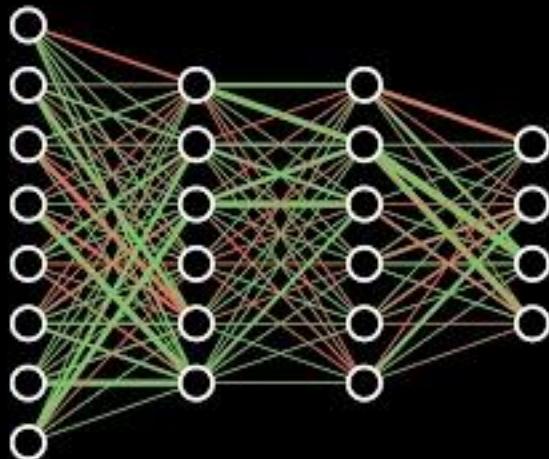
$$\nabla_{\mathbf{W}^{(k)}} J = \mathbf{g} \mathbf{h}^{(k-1)\top} + \lambda \nabla_{\mathbf{W}^{(k)}} \Omega(\theta)$$

Propagate the gradients w.r.t. the next lower-level hidden layer's activations:

$$\mathbf{g} \leftarrow \nabla_{\mathbf{h}^{(k-1)}} J = \mathbf{W}^{(k)\top} \mathbf{g}$$

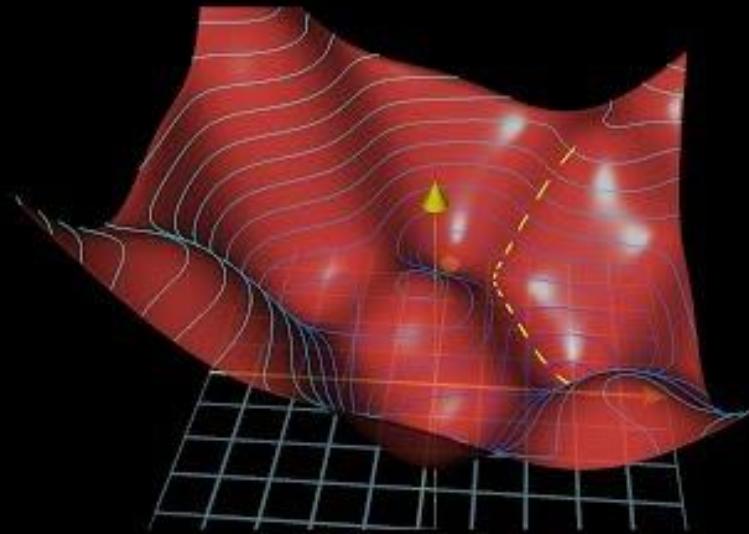
end for

Neural Networks



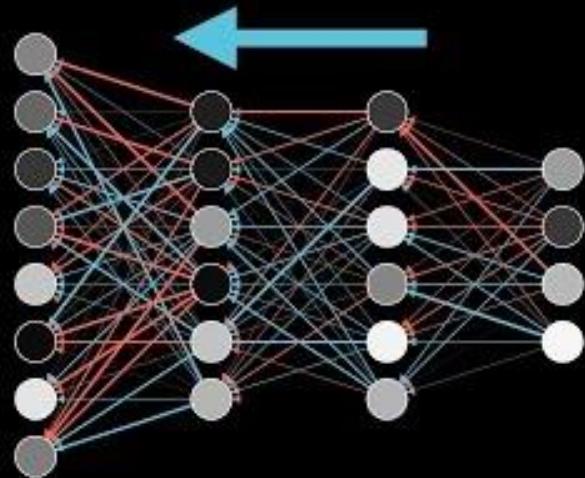
From the
ground up

How machines learn



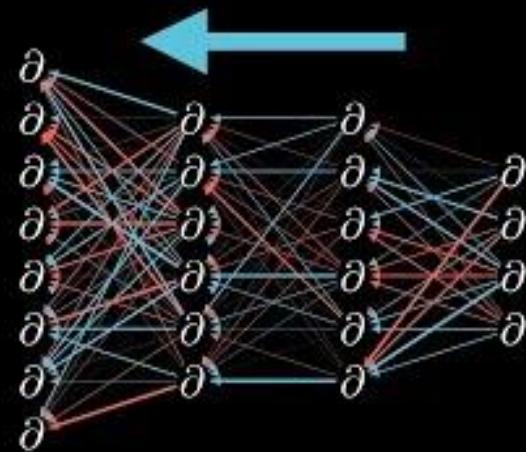
https://www.youtube.com/watch?v=IHZwWFHWa-w&list=PLZHQBObOWTQDNU6R1_67000Dx_ZCJB-3pi&index=2

Backpropagation



https://www.youtube.com/watch?v=llg3gGewQ5U&list=PLZHQQbOWTQDNU6R1_67000Dx_ZCJB-3pi&index=3

Backpropagation calculus

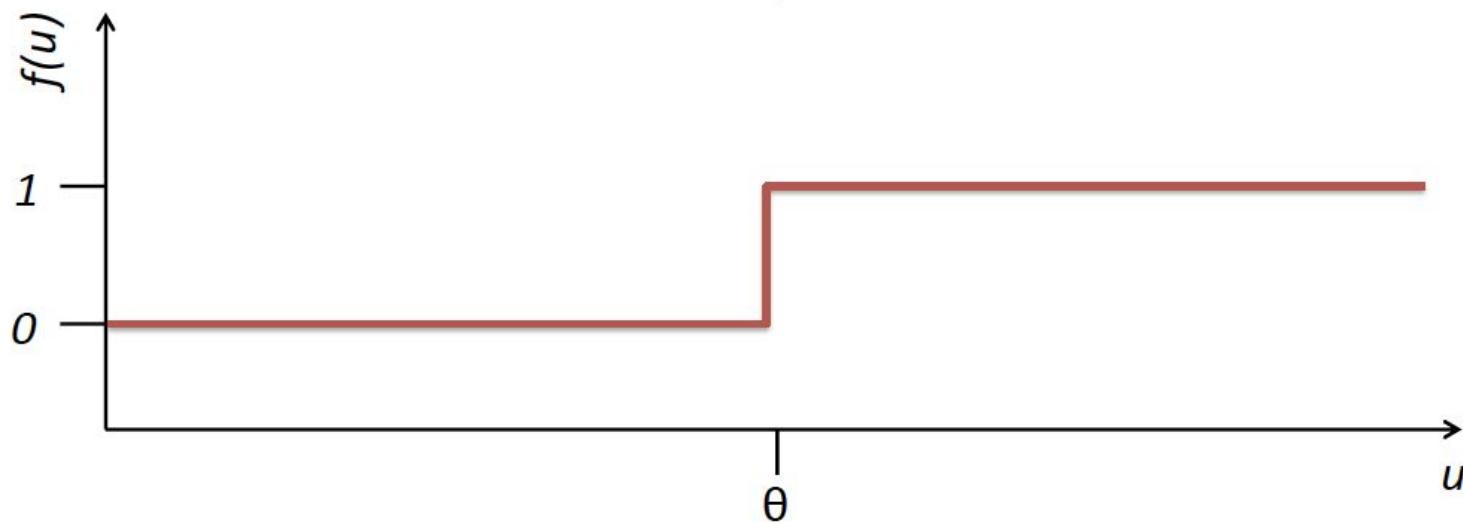


Activation Functions

Step Function

Non-continuous function:

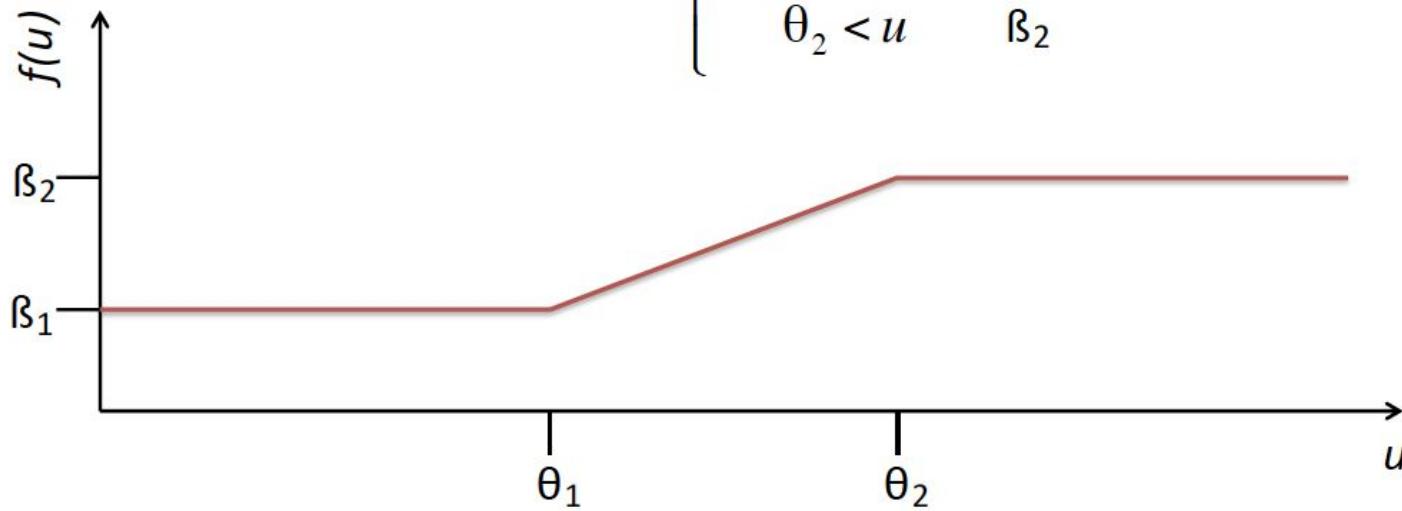
$$f(u) = \begin{cases} 0 & u < \theta \\ 1 & u \geq \theta \end{cases} \quad f : \mathfrak{R} \rightarrow \{0,1\}$$



Rectified linear (ReLU)

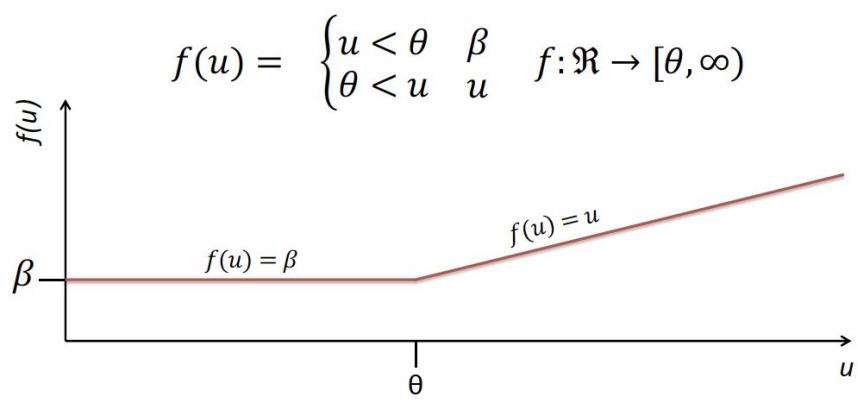
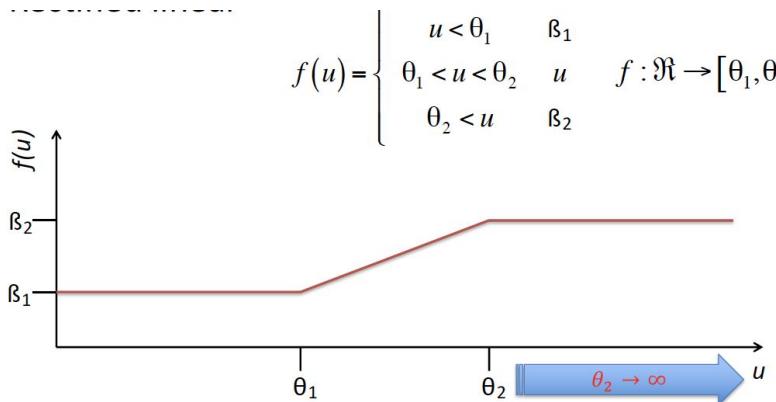
Non-continuously differentiable function

$$f(u) = \begin{cases} \beta_1 & u < \theta_1 \\ u & \theta_1 < u < \theta_2 \\ \beta_2 & \theta_2 < u \end{cases} \quad f : \mathfrak{N} \rightarrow [\theta_1, \theta_2]$$



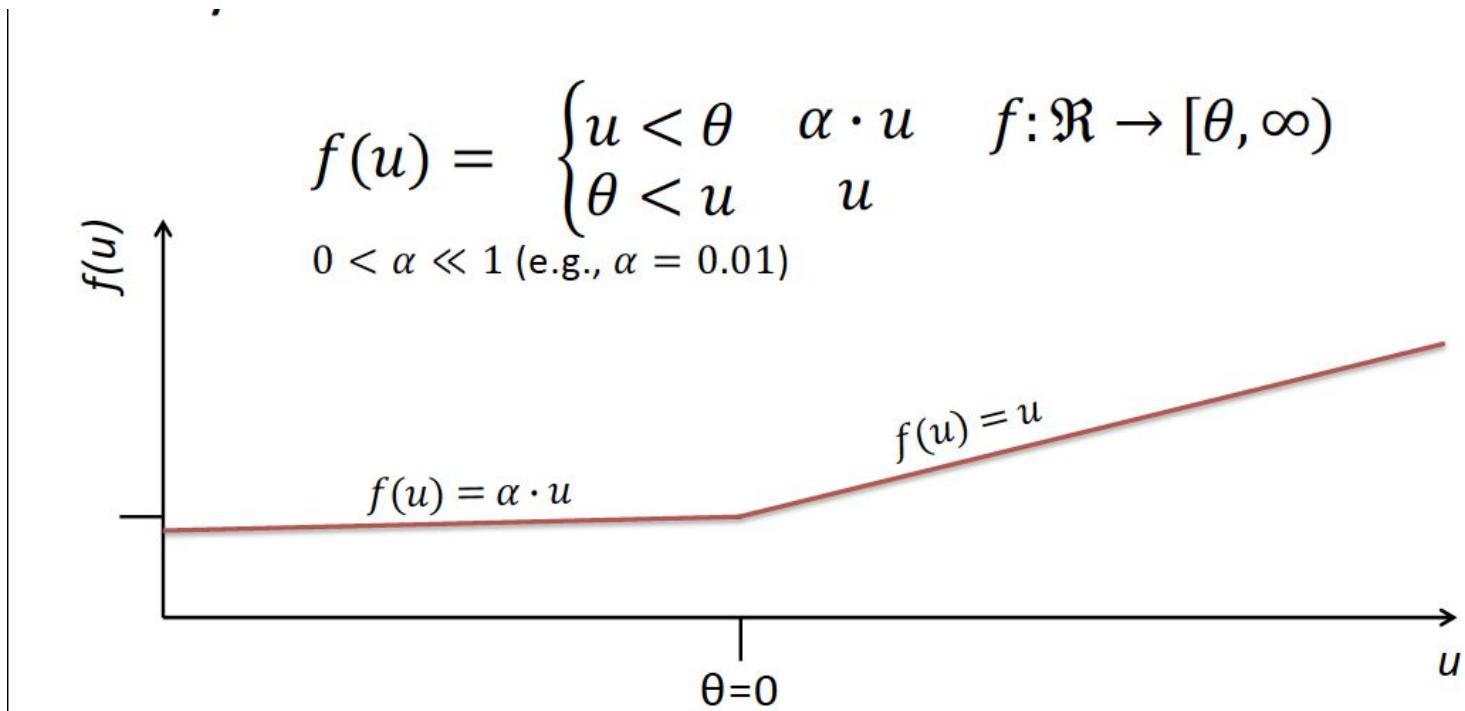
Rectified linear

Non-continuously differentiable functions



Leaky rectified linear

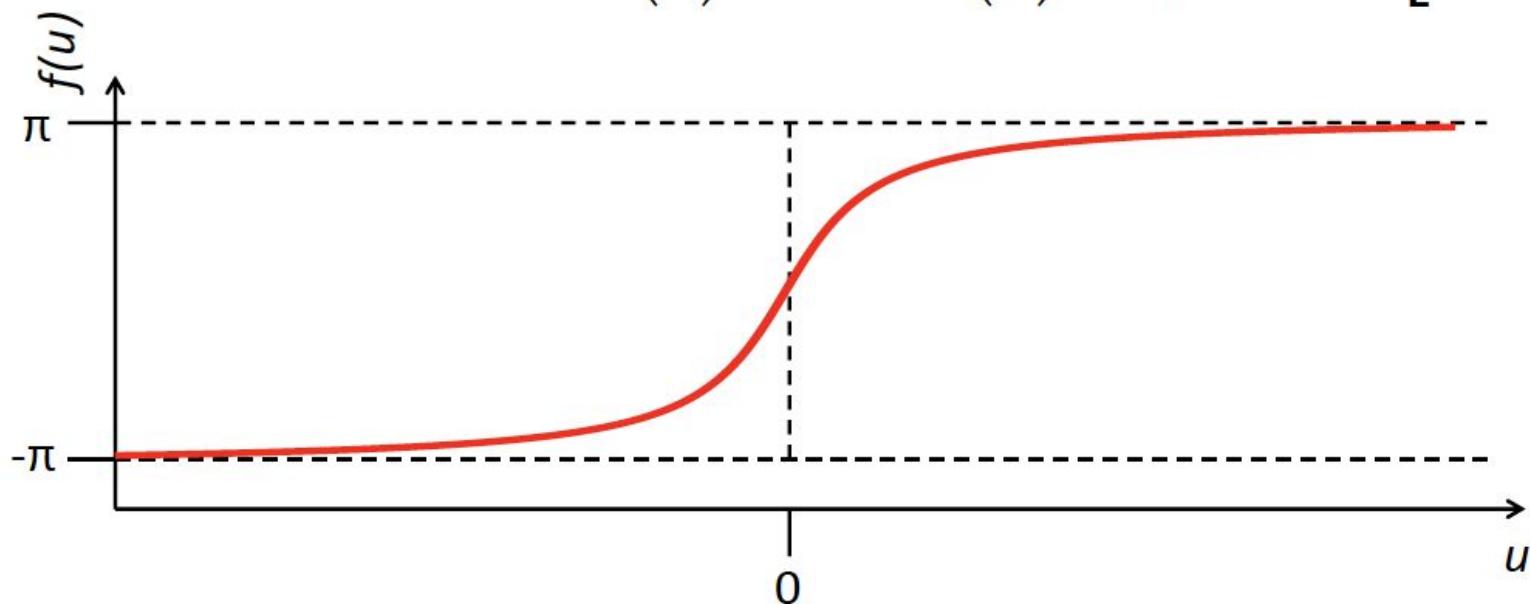
Non-continuously differentiable functions



Sigmoid functions: Arctan

continuously differentiable

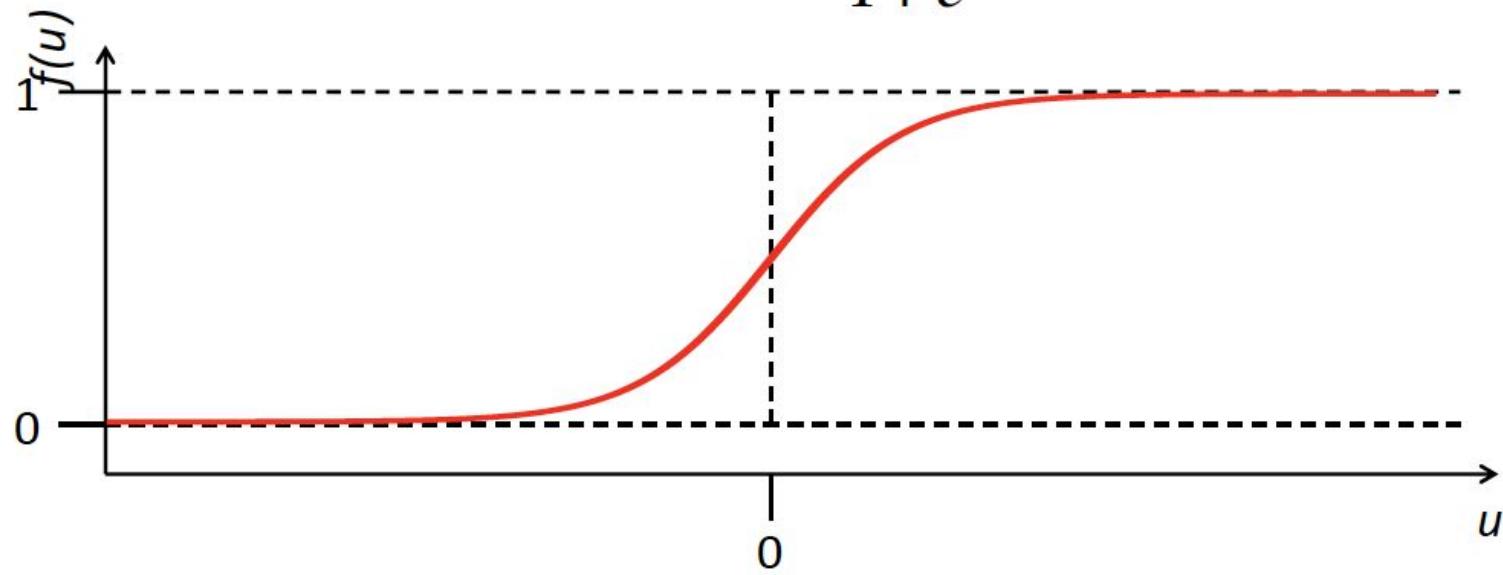
$$f(u) = \arctan(u) \quad f : \mathfrak{R} \rightarrow [-\pi, \pi]$$



Sigmoid functions: Logistic function

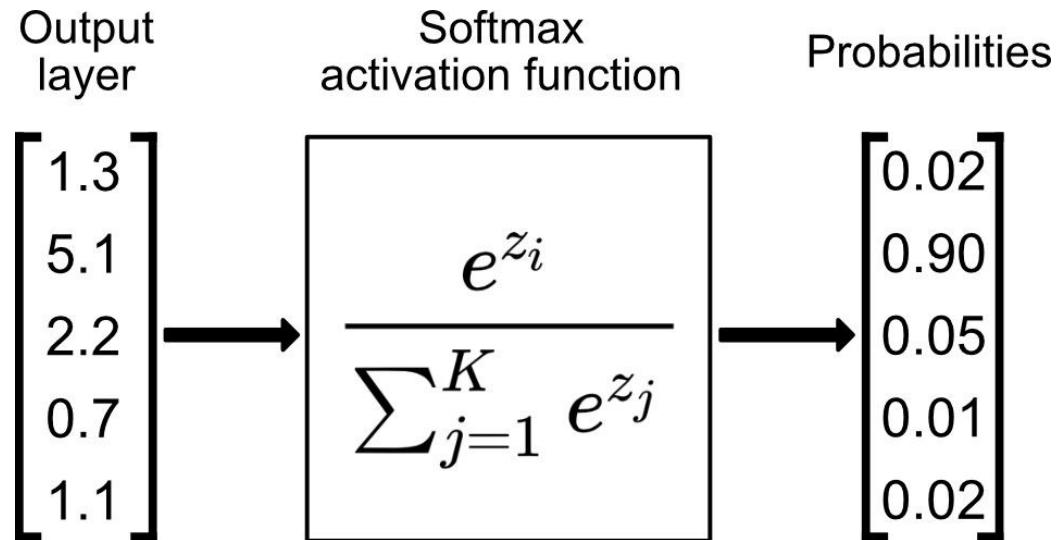
continuously differentiable

$$f(u) = \frac{1}{1 + e^{-\sigma u}} \quad f : \mathfrak{N} \rightarrow [0,1]$$



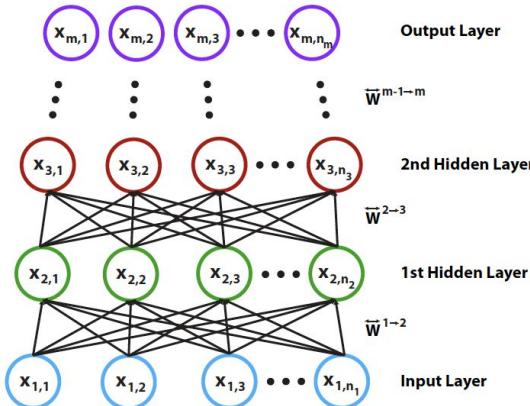
Softmax function

$$p(\vec{x}^m) = p(x_{m,1}, x_{m,2}, \dots, x_{m,n_m}) = \left(\frac{e^{x_{m,1}}}{\sum_{j=1}^{n_m} e^{x_{m,j}}}, \frac{e^{x_{m,2}}}{\sum_{j=1}^{n_m} e^{x_{m,j}}}, \dots, \frac{e^{x_{m,n_m}}}{\sum_{j=1}^{n_m} e^{x_{m,j}}} \right) = \frac{e^{\vec{x}^m}}{\sum_{j=1}^{n_m} e^{x_{m,j}}}$$



Feed Forward Neural Networks

- Information only flows forward, There are no feedback connections in which outputs of the model are fed back into itself
- Usually the number of neurons in consecutive layers decreases(lower level layers have more neurons than higher)
- For regression
 - The number of neurons at the final layers = 1
 - The single neuron at the output layer returns a continuous value
- For classification:
 - For classification into k classes, we set the number of units in the output layer to k
 - The neuron i in the output layer returns the probability that the input should be classified as class i.
 - Probabilities are computed using the softmax function



$$\vec{x}^i = (x_{i,1}, x_{i,2}, \dots, x_{i,n_i})^T$$

$$\vec{W}^{i \rightarrow i+1} = \begin{bmatrix} w_{1,1}^{i \rightarrow i+1} & w_{1,2}^{i \rightarrow i+1} & \dots & w_{1,n_i}^{i \rightarrow i+1} \\ w_{2,1}^{i \rightarrow i+1} & w_{2,2}^{i \rightarrow i+1} & \dots & w_{2,n_i}^{i \rightarrow i+1} \\ \vdots & \vdots & \dots & \vdots \\ w_{n_{i-1},1}^{i \rightarrow i+1} & w_{n_{i-1},2}^{i \rightarrow i+1} & \dots & w_{n_{i-1},n_i}^{i \rightarrow i+1} \end{bmatrix}$$

Deep Learning Programming

Keras

- It doesn't handle low-level computations itself; instead, it passes it to another library called the Backend. After version 2.4 only TensorFlow is supported.
- In the newer versions, Keras acts as an interface for the TensorFlow
- Supports distributed training of deep-learning models on clusters of GPUs and TPUs
- `input_shape=(784,)` is equal to `batch_input_shape=(None, 784)` is equal to `input_dim=784`
- `input_shape=(10, 64)` is equal to `batch_input_shape=(None, 10, 64)` is equal to `input_dim=64`

Sequential API

```
model = Sequential()  
  
model.add(Dense(10,activation='relu')) # You don't need to specify input  
model.add(Dense(5,activation='relu'))  
model.add(Dense(1))
```

Functional API

- You define an input layer and specify the input size
- You define each layer by specifying a name, the number of neurons, activation function, etc.
- Put the previous layer in parentheses at the end of the expression

```
input_layer = Input(shape=(3,))

Layer_1 = Dense(4, activation="relu")(input_layer)

Layer_2 = Dense(4, activation="relu")(Layer_1)

output_layer= Dense(1, activation="linear")(Layer_2)

##Defining the model by specifying the input and output layers

model = Model(inputs=input_layer, outputs=output_layer)
```

Example: Keras

```
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)
```

Example: CNN Model Using Keras

```
model = Sequential()  
model.add(Conv2D(32, (3, 3), activation='relu', input_shape=(32, 32,  
3)))  
model.add(MaxPool2D())  
model.add(Conv2D(16, (3, 3), activation='relu'))  
model.add(MaxPool2D())  
model.add(Flatten())  
model.add(Dense(10, activation='softmax'))
```

R Keras Package

```
library(keras)

model <- keras_model_sequential()

model %>% layer_dense(units = 256, activation = 'relu',
input_shape = c(784)) %>%

  layer_dropout(rate = 0.4) %>%
  
  layer_dense(units = 128, activation = 'relu') %>%
  
  layer_dropout(rate = 0.3) %>%
  
  layer_dense(units = 10, activation = 'softmax')

model %>% compile(
  loss = 'categorical_crossentropy',
  optimizer = optimizer_rmsprop(),
  metrics = c('accuracy')
)

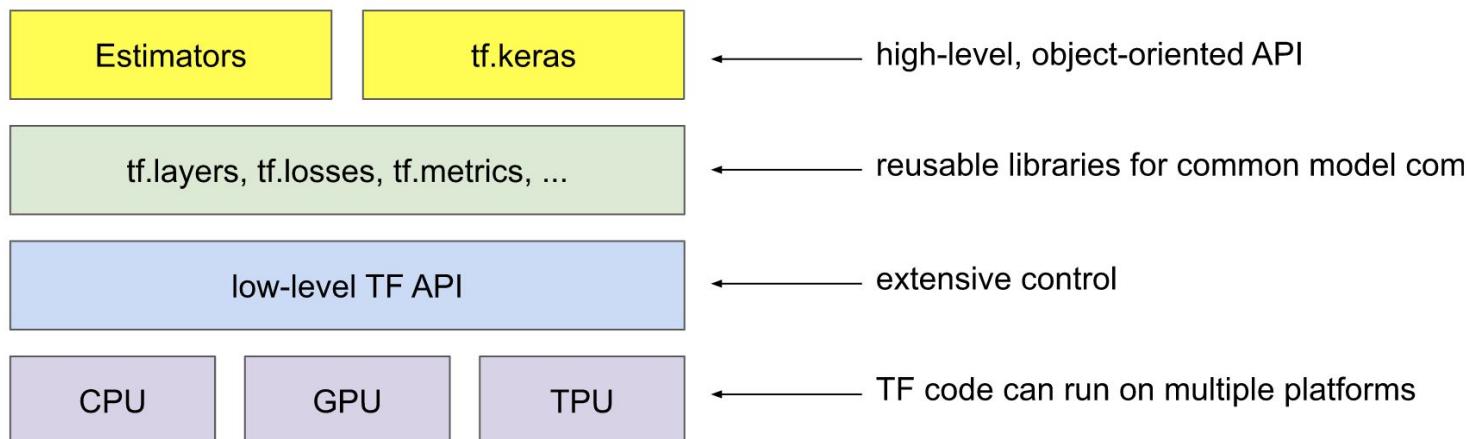
history <- model %>% fit(
  x_train, y_train,
  epochs = 30, batch_size = 128,
  validation_split = 0.2
)
```

Example code:

<https://cran.r-project.org/web/packages/keras/vignettes/index.html>

Tensorflow High-Level APIs

- Tensorflow 2.0 has two high-level deep learning APIs:
 - tf.keras
 - Tf.estimator (not recommended for new code, use keras instead)



TensorFlow 2

- TensorFlow 2.0 introduced many changes
 - Changed the automatic differentiation scheme (automatically calculating the gradient vector of a model with respect to each of its parameters) from the static computational graph to the Define-by-Run
 - TensorFlow includes an “eager execution” mode (operations are evaluated immediately as opposed to being added to a computational graph which is executed later).
 - Code executed eagerly can be examined step-by step-through a debugger, since data is augmented at each line of code rather than later in a computational graph
 - The session and placeholder objects do not exist anymore

Eager Execution

- In Tensorflow 1.0 the operations were performed inside a session
- Tensorflow 2.0 supports eager execution (no need for placeholders)

```
import tensorflow as tf

var1 = tf.constant(8)
var2 = tf.constant(9)
result_tf_1 = var1 + var2
sess = tf.Session()
print("The result for Tensorflow 1.0 is:",sess.run(result_tf_1))
```

```
import tensorflow as tf

var3 = tf.constant(5)
var4 = tf.constant(6)
result_tf_2 = var3 + var4
print("The result of tensorflow 2 is:",result_tf_2)
```

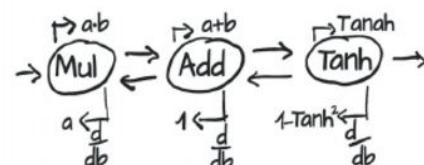
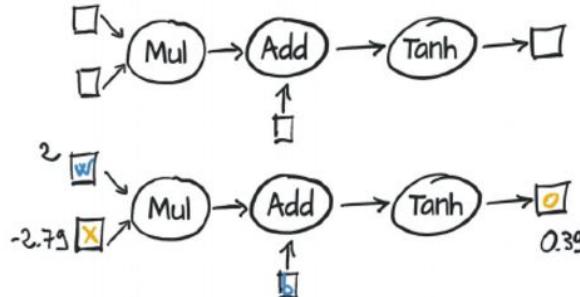
Define-and-run vs static computational graph

- **Define-and-run** (Static computational graph): A model is constructed, before execution, as a computational graph structure, which later gets executed with different inputs while remaining fixed (Theano).
- **Define-by-run** (dynamic computational graph): general-purpose AD capability available where a model is a regular program in the host programming language, whose execution dynamically constructs a computational graph on-the-fly that can freely change in each iteration (PyTorch)

Define-and-run vs Define-by-run

STATIC GRAPH

$$o = \tanh(wx + b)$$



$$\frac{d_o}{dw} = \frac{dTanh}{dwx+b} \frac{dwx+b}{dw} = (1 - \tanh^2(wx+b))x$$

COMPILE SYMBOLIC GRAPH
↓
EVALUATE

AUTOMATIC DIFFERENTIATION

$$o = \tanh(wx + b)$$

$$x = -2.79$$

$$x_1 = wx = 2 \times (-2.79) = -5.58$$

$$x_2 = x_1 + b = -5.58 + 6 = 0.42$$

$$o = \tanh x_2 = \tanh 0.42 = 0.3969\dots$$

GREEDY EVALUATION
(NO GRAPH)

STATEMENTS

$$x_1 = wx$$



$$x_1 \approx -5.58$$

$$x_2 = x_1 + b$$



$$x_2 \approx 0.42$$

$$o = \tanh x_2$$



DYNAMIC GRAPH
"DEFINE BY RUN"



$$o \approx 0.3969\dots$$

BACKWARD

Linear Regression: TensorFlow 2

```
import tensorflow as tf
import numpy as np

x_data = np.random.rand(100).astype(np.float32)
y_data = x_data * 0.1 + 0.3

W = tf.Variable(tf.random.normal([1]))
b = tf.Variable(tf.zeros([1]))

# mean squared error between y_data and computed y
def mse_loss():
    y = W * x_data + b
    loss = tf.reduce_mean(tf.square(y - y_data))
    return loss

optimizer = tf.keras.optimizers.Adam()
for step in range(5000):
    optimizer.minimize(mse_loss, var_list=[W,b])
    if step % 500 == 0:
        print(step, W.numpy(), b.numpy())
```

PyTorch

- Good debugging Support
- PyTorch provides two APIs:
 - Low-level Tensor computing (similar to NumPy)
 - High-level API
- PyTorch has a class called Tensor (`torch.Tensor`) to store and operate on homogeneous multidimensional rectangular arrays of numbers, similar to NumPy Arrays
 - PyTorch Tenors can be operated on a CUDA-capable NVIDIA GPU

Example: Low-level Tensor API

```
import torch

dtype = torch.float

device = torch.device("cpu") # This executes all calculations on the CPU
# device = torch.device("cuda:0") # This executes all calculations on the GPU

# Creation of a tensor and filling of a tensor with random numbers
a = torch.randn(2, 3, device=device, dtype=dtype)

b = torch.randn(2, 3, device=device, dtype=dtype)

print(a*b)
```

High Level API Example: Standard Neural Net

```
import torch
from torch import nn

class NeuralNetwork(nn.Module):
    def __init__(self): # super(NeuralNetwork, self).__init__()
        self.flatten = nn.Flatten()
        self.linear_relu_stack = nn.Sequential(
            nn.Linear(28*28, 512), # Linear Layers have an input and output shape
            nn.ReLU(),
            nn.Linear(512, 512),
            nn.ReLU(),
            nn.Linear(512, 10),
    )
    def forward(self, x):
        x = self.flatten(x)
        logits = self.linear_relu_stack(x)
        return logits
```

Neural networks are defined as classes, inherit from nn.Module class

Layers and variables are defined in the __init__ method

This function defines the forward propagation logic which implies the structure of the network

Example: CNN Using PyTorch

```
class Net(nn.Module):
    def __init__(self):
        super(Net, self).__init__()
        self.conv1 = nn.Conv2d(3, 32, 3)
        self.conv2 = nn.Conv2d(32, 16, 3)
        self.fc1 = nn.Linear(16 * 6 * 6, 10)
        self.pool = nn.MaxPool2d(2, 2)

    def forward(self, x):
        x = self.pool(F.relu(self.conv1(x)))
        x = self.pool(F.relu(self.conv2(x)))

        # view returns a new tensor with the same data as the self tensor but of a different shape.
        x = x.view(-1, 16 * 6 * 6)
        x = F.log_softmax(self.fc1(x), dim=-1)
        return x
```

Deep Learning Regularization

Weight Decay/Ridge Regression/Tikhonov regularization

- This regularization strategy drives the weights closer to the origin by adding a regularization term $\Omega(\boldsymbol{\theta}) = \frac{1}{2} \|\boldsymbol{w}\|_2^2$

$$\tilde{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = \frac{\alpha}{2} \boldsymbol{w}^\top \boldsymbol{w} + J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}),$$

$$\nabla_{\boldsymbol{w}} \tilde{J}(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}) = \alpha \boldsymbol{w} + \nabla_{\boldsymbol{w}} J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y}).$$

- To take a single gradient step to update the weights, we perform this update:

$$\boldsymbol{w} \leftarrow \boldsymbol{w} - \epsilon (\alpha \boldsymbol{w} + \nabla_{\boldsymbol{w}} J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y})).$$

$$\boldsymbol{w} \leftarrow (1 - \epsilon \alpha) \boldsymbol{w} - \epsilon \nabla_{\boldsymbol{w}} J(\boldsymbol{w}; \boldsymbol{X}, \boldsymbol{y})$$

Keras/Pytorch Weight Decay

Keras Weight Decay (L2), L1 and L1_L2 Regularizers:

```
keras.regularizers.l1(0.01)
```

```
keras.regularizers.l2(0.01)
```

```
keras.regularizers.l1_l2(l1=0.01, l2=0.01)
```

```
model.add(Dense(32, kernel_regularizer=l2(0.01), bias_regularizer=l2(0.01)))
```

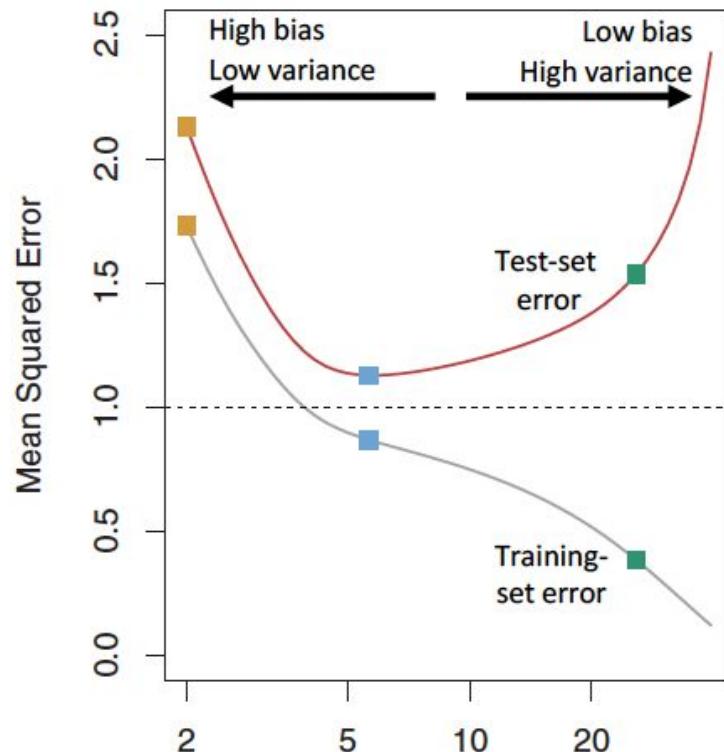
Pytorch Regularizers:

```
optimizer = torch.optim.SGD(model.parameters(), lr=1e-3, weight_decay=1e-4)
```

```
optimizer = torch.optim.Adam(model.parameters(), lr=1e-3, weight_decay=1e-4)
```

Early Stopping

- Early stopping may be used either alone or in conjunction with other regularization strategies.
- Early stopping requires a validation set, which means some training data is not fed to the model. To best exploit this extra data, one can perform extra training after the initial training with early stopping has completed.
 - One strategy is to initialize the model again and retrain on all of the data. In this second training pass, we train for the same number of steps as the early stopping procedure determined was optimal in the first pass.



Early Stopping: Keras

```
callback =  
[EarlyStopping(monitor='val_loss',patience=2),ModelCheckpoint(filepath='best_model.h5',  
monitor='val_loss', save_best_only=True)]  
  
model = tf.keras.models.Sequential([tf.keras.layers.Dense(10)])  
  
model.compile(tf.keras.optimizers.SGD(), loss='mse')  
  
history = model.fit(X,Y, epochs=10, batch_size=1, callbacks=[callback])
```

Early Stopping: Pytorch

```
def traindata(device, model, epochs, optimizer, loss_function, train_loader, valid_loader):
```

```
...
```

Early stopping

```
    current_loss = validation(model, device, valid_loader, loss_function)  
    print('The Current Loss:', current_loss)
```

```
    if current_loss > last_loss:
```

```
        trigger_times += 1
```

```
        print('Trigger Times:', trigger_times)
```

```
        if trigger_times >= patience:
```

```
            print('Early stopping!\nStart to test process.')
```

```
            return model
```

```
    else:
```

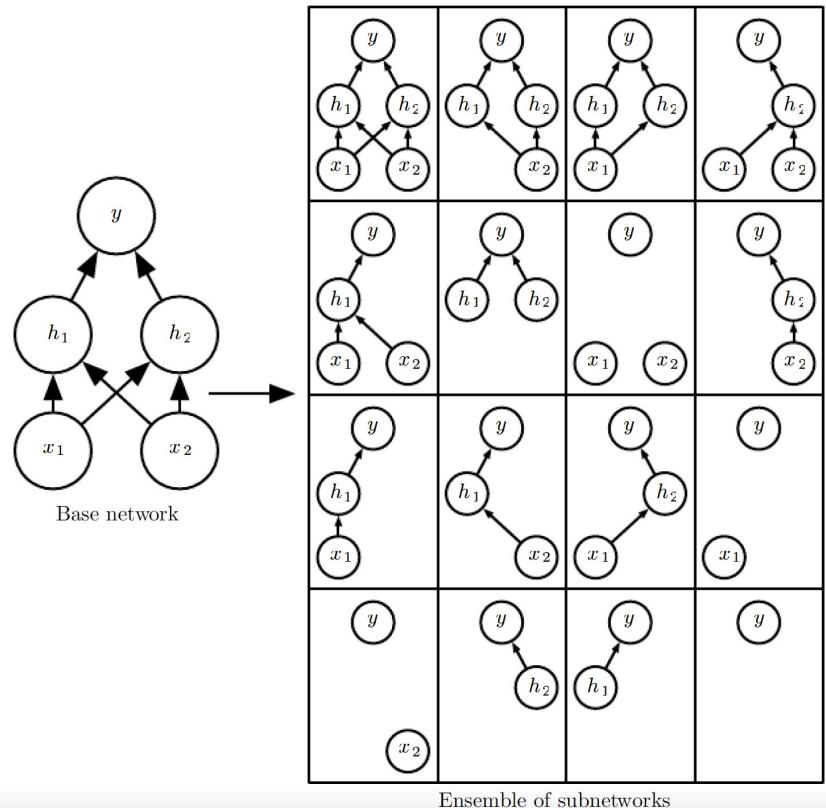
```
        print('trigger times: 0')
```

```
        trigger_times = 0
```

```
    last_loss = current_loss
```

Dropout

- Dropout provides an inexpensive approximation to training and evaluating a **bagged ensemble** of exponentially many neural networks.
- Specifically, dropout trains the ensemble consisting of all sub-networks that can be formed by removing non-output units from an underlying base network
- In most modern neural networks, based on a series of affine transformations and nonlinearities, we can effectively remove a unit from a network by multiplying its output value by zero.



Dropout: Keras

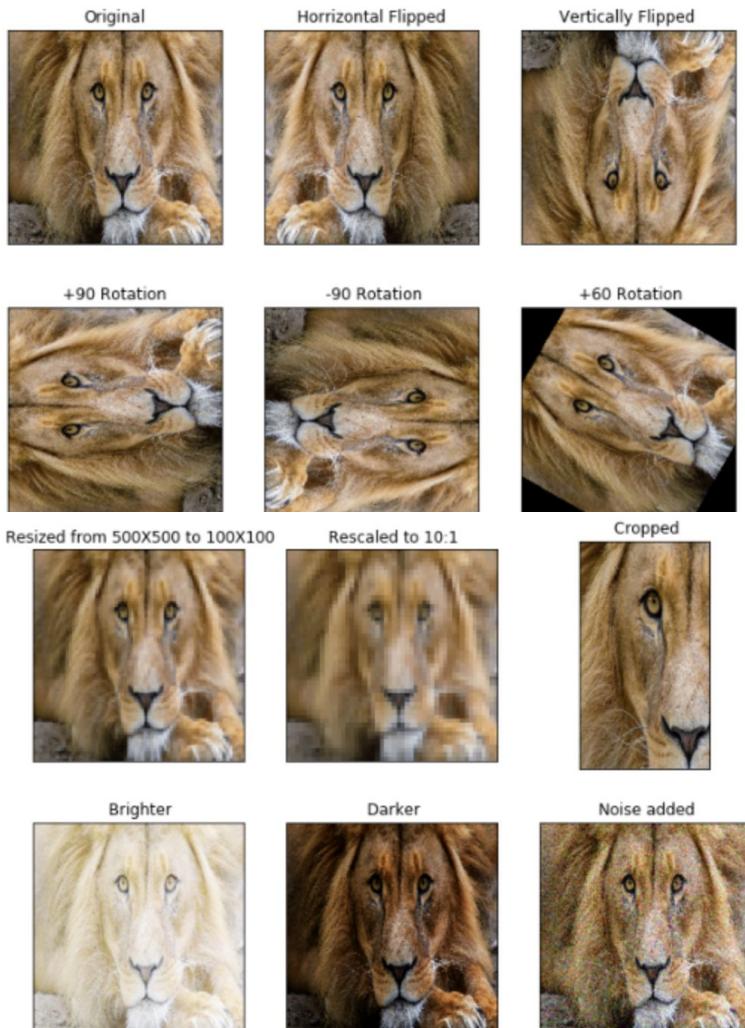
```
model = Sequential()
model.add(Dense(60, input_shape=(60,), activation='relu',
kernel_constraint=MaxNorm(3)))
model.add(Dropout(0.2))
model.add(Dense(30, activation='relu',
kernel_constraint=MaxNorm(3)))
model.add(Dropout(0.2))
model.add(Dense(1, activation='sigmoid'))
```

Dropout: PyTorch

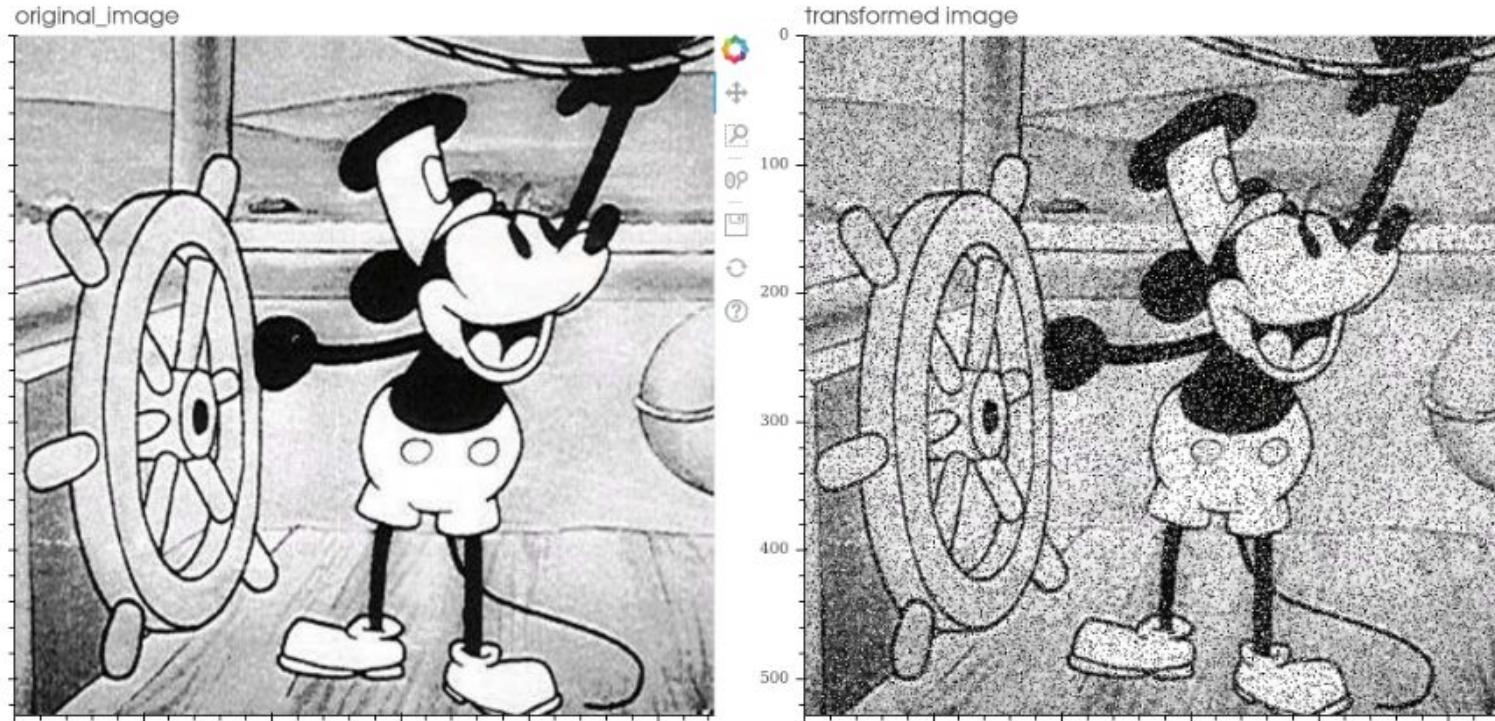
```
def __init__(self):  
    super(Net1, self).__init__()  
  
    self.hid1 = T.nn.Linear(4, 8)  
  
    self.drop1 = T.nn.Dropout(0.50)  
  
    self.hid2 = T.nn.Linear(8, 8)  
  
    self.drop2 = T.nn.Dropout(0.25)  
  
    self.outp = T.nn.Linear(8, 1)
```

Dataset Augmentation: Image

- One must be careful not to apply transformations that would change the correct class. For example, optical character recognition tasks require recognizing the difference between 'b' and 'd' and the difference between '6' and '9', so horizontal flips and 180° rotations are not appropriate ways of augmenting datasets for these tasks.



Dataset Augmentation: Noise Injection Example



Data Augmentation in Keras

```
data_augmentation = tf.keras.Sequential([  
    layers.experimental.preprocessing.RandomFlip("horizontal_and_vertical"),  
    layers.experimental.preprocessing.RandomRotation(0.2)])
```

```
datagen = ImageDataGenerator(rotation_range=90)
```

Data augmentation in PyTorch

- Transforms library from the torchvision package has popular datasets, model architectures, and common image transformations for Computer Vision tasks

```
from torchvision import transforms as tr

from torchvision.transforms import Compose

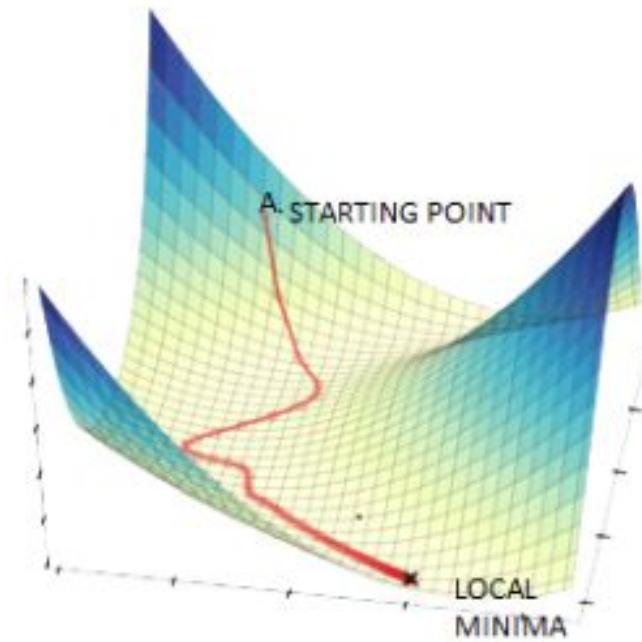
pipeline = Compose(
    [tr.RandomRotation(degrees = 90),
     tr.RandomRotation(degrees = 270)])

augmented_image = pipeline(img = img)
```

Deep Learning Optimization

Optimization Goal

- We want to converge faster and to a better minima
- Machine learning usually acts indirectly. In most machine learning scenarios, we care about some performance measure P , that is defined with respect to the test set and may also be intractable. We therefore optimize P only indirectly. We reduce a different cost function $J(\theta)$ in the hope that doing so will improve P . This is in contrast to pure optimization, where minimizing J is a goal in and of itself.
 - If we knew the true distribution $p_{\text{data}}(x, y)$, risk minimization would be an optimization task solvable by an optimization algorithm.
 - However, when we do not know $p_{\text{data}}(x, y)$ but only have a training set of samples, we have a machine learning problem.



Optimizers Cheat Sheet

Gradient Descent:

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta)$$

Stochastic Gradient Descent

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta; \text{sample})$$

Mini-Batch Gradient Descent

$$\theta = \theta - \alpha \cdot \nabla_{\theta} J(\theta; N \text{ samples})$$

SGD + Momentum

$$v = \gamma \cdot v + \eta \cdot \nabla_{\theta} J(\theta)$$

$$\theta = \theta - \alpha \cdot v$$

SGD + Momentum + Acceleration

$$v = \gamma \cdot v + \eta \cdot \nabla_{\theta} J(\theta - \gamma \cdot v)$$

$$\theta = \theta - \alpha \cdot v$$

Adagrad

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{G_{t,ii}} + \epsilon} \nabla_{\theta_{t,i}} J(\theta_{t,i})$$

Adadelta

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{E[G_{t,ii}]} + \epsilon} \nabla_{\theta_{t,i}} J(\theta_{t,i})$$

Adam

$$\theta_{t+1,i} = \theta_{t,i} - \frac{\eta}{\sqrt{E[G_{t,ii}]} + \epsilon} \times E[g_{t,i}]$$

Keras Example

```
keras.optimizers.RMSprop(lr=0.001, rho=0.9, epsilon=None, decay=0.0)
```

```
keras.optimizers.Adagrad(lr=0.01, epsilon=None, decay=0.0)
```

```
keras.optimizers.Adadelta(lr=1.0, rho=0.95, epsilon=None, decay=0.0)
```

```
keras.optimizers.Adam(lr=0.001, beta_1=0.9, beta_2=0.999, epsilon=None, decay=0.0, amsgrad=False)
```

Pytorch Example

```
optimizer = optim.SGD(model.parameters(), lr=0.01, momentum=0.9)
```

```
optimizer = optim.Adam(model.parameters(), lr=0.0001)
```

```
optimizer= optim.RMSprop(model.parameters(), lr=0.0005)
```

Batch normalization

- Batch normalization is a method used to make training of artificial neural networks faster and more stable through normalization of the layers' inputs by re-centering and re-scaling (Let H be a minibatch of activations of the layer to normalize):

$$H' = \frac{H - \mu}{\sigma},$$

$$\mu = \frac{1}{m} \sum_i H_{i,:}$$

$$\sigma = \sqrt{\delta + \frac{1}{m} \sum_i (H - \mu)_i^2},$$

Where δ is a small positive value such as 10^{-8} , imposed to avoid encountering the undefined gradient of \sqrt{z} at $z=0$.

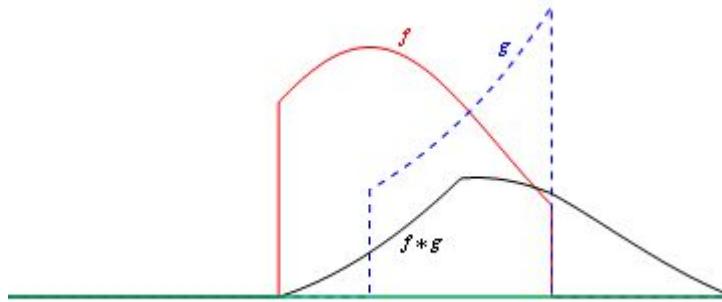
Keras Example

```
model = Sequential()  
  
model.add(Dense(10, activation="relu", input_shape=(X_train.shape[1],)))  
  
model.add(BatchNormalization())  
  
model.add(Dense(8, activation="relu"))  
  
model.add(BatchNormalization())  
  
model.add(Dense(1, activation="sigmoid"))
```

PyTorch Example

```
class MLP(nn.Module):
    def __init__(self):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Flatten(),
            nn.Linear(32 * 32 * 3, 64),
            nn.BatchNorm1d(64),
            nn.ReLU(),
            nn.Linear(64, 32),
            nn.BatchNorm1d(32),
            nn.ReLU(),
            nn.Linear(32, 10)
        )
```

Convolutional Neural Networks



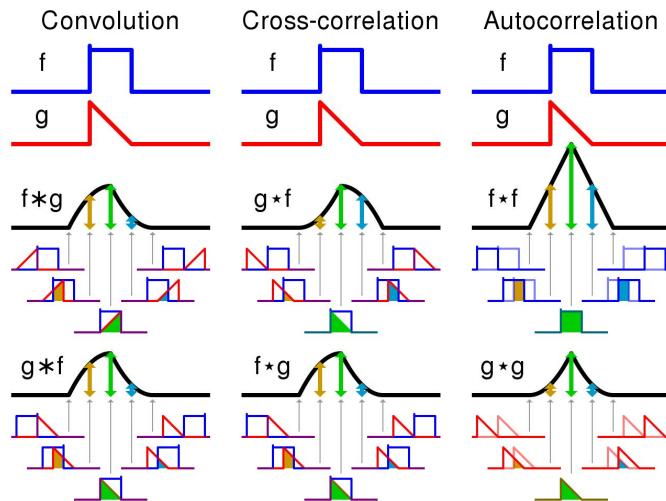
$$s(t) = \int x(a)w(t-a)da.$$

$$s(t) = (x * w)(t).$$

$$S(i, j) = (I * K)(i, j) = \sum \sum I(m, n)K(i - m, j - n).$$

- many neural network libraries implement a related function called the cross-correlation, which is the same as convolution but without flipping the kernel

$$S(i, j) = (K * I)(i, j) = \sum_m \sum_n I(i + m, j + n)K(m, n).$$



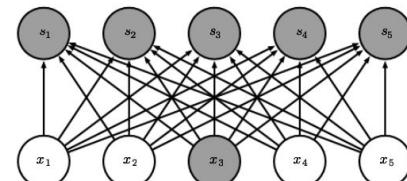
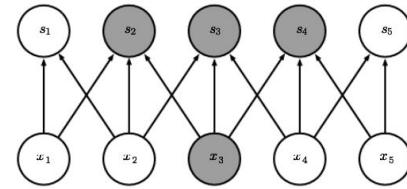
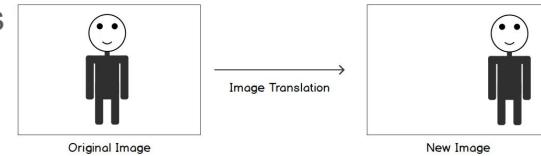
The diagram shows the computation of a convolution step using a 3x3 kernel. The input image is a 5x5 grid with values [0, 25, 75, 80, 80; 0, 75, 80, 80, 80; 0, 75, 80, 80, 80; 0, 70, 75, 80, 80; 0, 0, 0, 0, 0]. The kernel is a 3x3 matrix with values [-1, 0, 1; -2, 0, 2; -1, 0, 1]. The result of the multiplication and summation is shown in the bottom right corner, resulting in a value of 235.

0	25	75	80	80
0	75	80	80	80
0	75	80	80	80
0	70	75	80	80
0	0	0	0	0

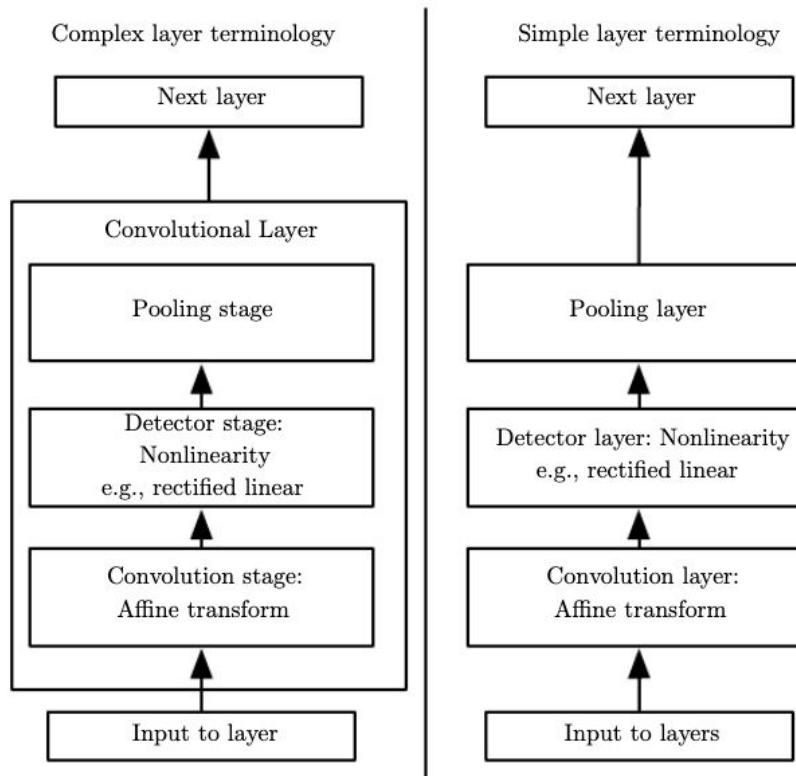
$$\begin{matrix} & \begin{matrix} -1 & 0 & 1 \\ -2 & 0 & 2 \\ -1 & 0 & 1 \end{matrix} \end{matrix} = \begin{matrix} & \begin{matrix} 0 & 0 & 80 \\ 0 & 0 & 80 \\ 0 & 0 & 80 \end{matrix} \end{matrix} \Sigma \begin{matrix} & 235 \end{matrix}$$

CNN

- Specialized kind of neural network for processing data that has a known grid-like topology, Examples include
 - Time-series data, which can be thought of as a 1-D grid taking samples at regular time intervals
 - Image data, which can be thought of as a 2-D grid of pixels.
- Convolution leverages three important ideas that can help improve a machine learning system
 - Sparse interactions
 - Parameter sharing
 - Parameter sharing refers to using the same parameter for more than one function in a model. In a traditional neural net, each element of the weight matrix is used exactly once when computing the output of a layer. It is multiplied by one element of the input and then never revisited.
 - Equivariant representations

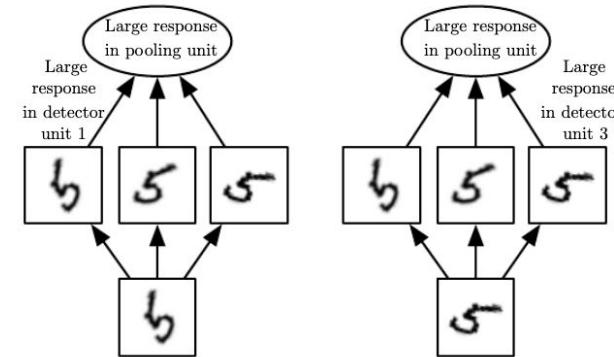
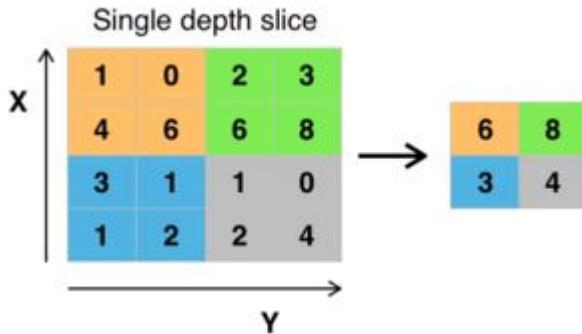


CNN Blocks



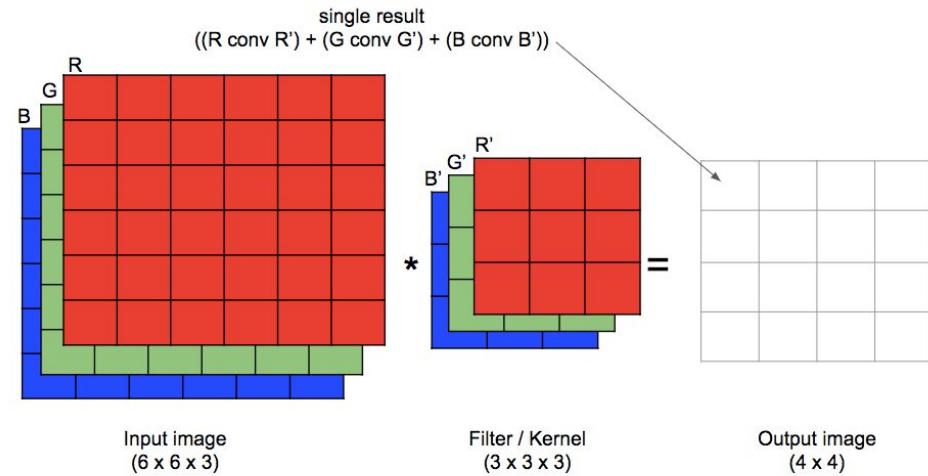
pooling

- pooling helps to make the representation approximately invariant to small translations of the input. Invariance to translation means that if we translate the input by a small amount, the values of most of the pooled outputs do not change.



From One Convolutional Layer to The Next

- Assume, we have a color image (RGB) of size 224×224 as input.
- The input to the first conv will be $224 \times 224 \times 3$
- Assume we have 32 kernels of size 3×3 at the first conv layer. if we have “same” padding, then the output of this layer will be of dimensions $224 \times 224 \times 32$.
- Note that our kernels will in fact be $3 \times 3 \times 3$ because of the three channels in the input
 - In fact the 3 channels of the input (RGB) are being merged/combined into one matrix (3d dot operator)

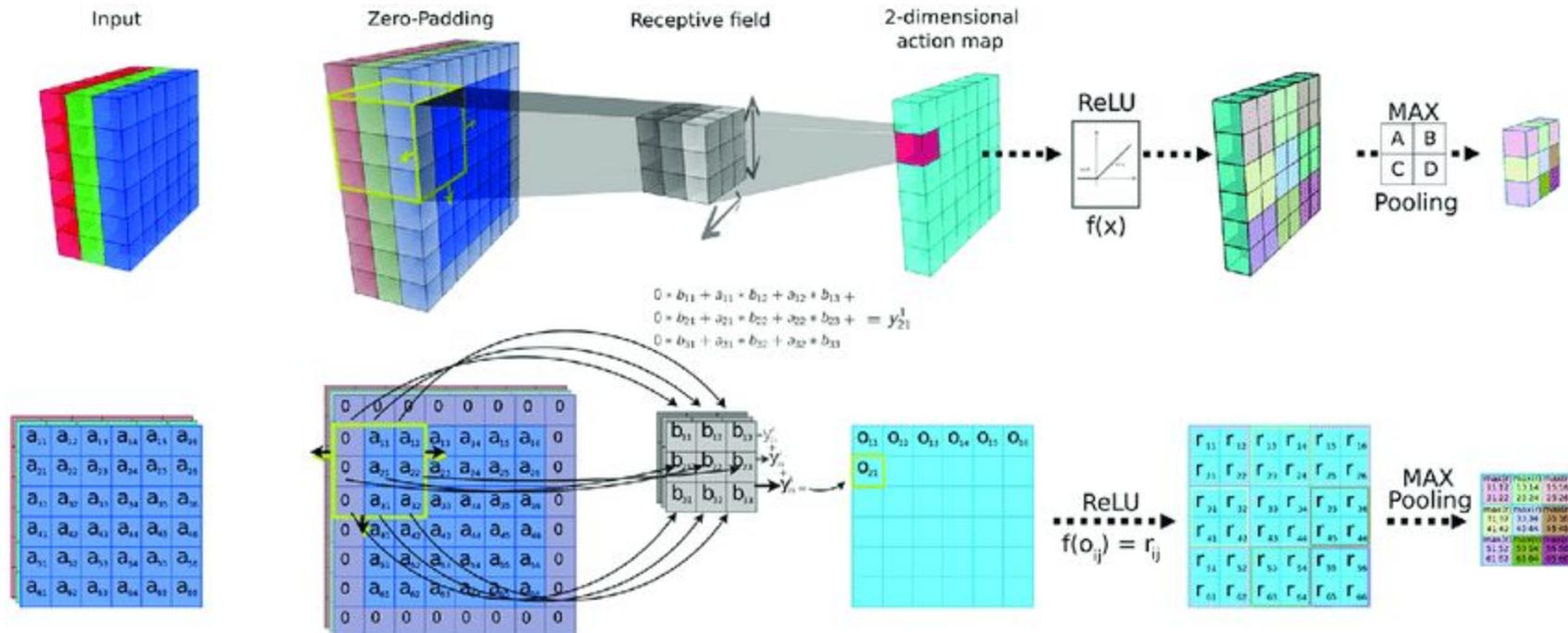


<https://blog.devgenius.io/deep-learning-in-gradient-descent-style-part-2-e159e2cf8a99>

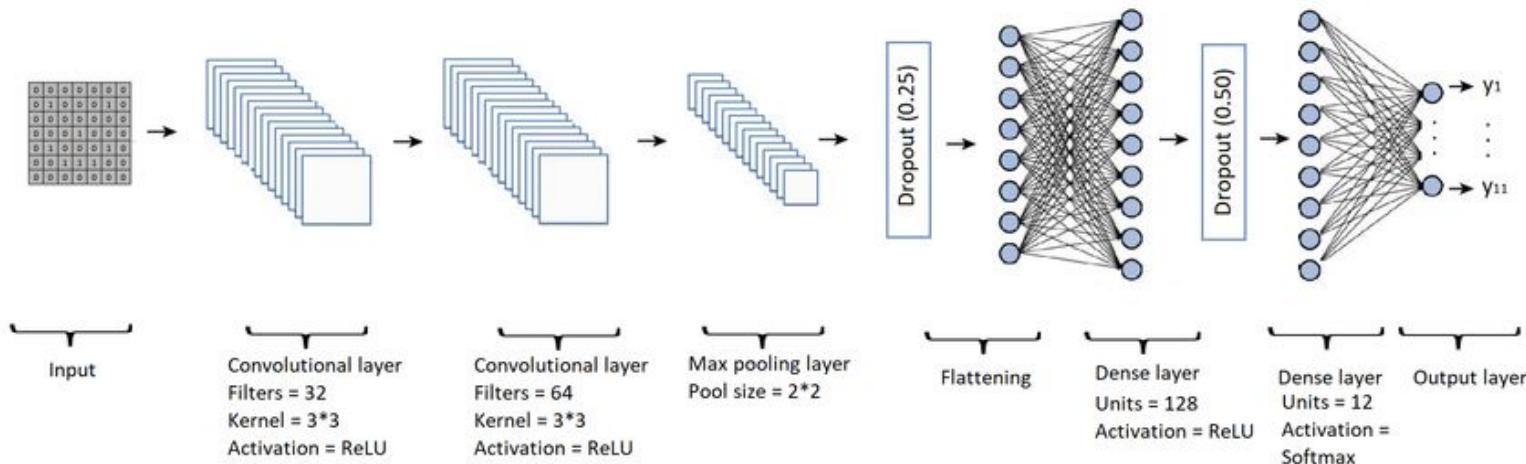
From One Convolutional Layer to The Next

- Now assume that in the second convolutional layer, we have 64 kernels of 3×3 with a “same” padding. Thus we will have 64 outputs of size 224×224 or the output will be of size $224 \times 224 \times 64$. Note that our kernels are in fact of size $3 \times 3 \times 32$, because the output of the previous layer has 32 channels.
- The number of kernels of the previous layers become the number of channels for the next layer.
 - If we would generate a separate output for each channel we would have 32 outputs for the first layer and 32×64 for the second layer

CNN Architectures

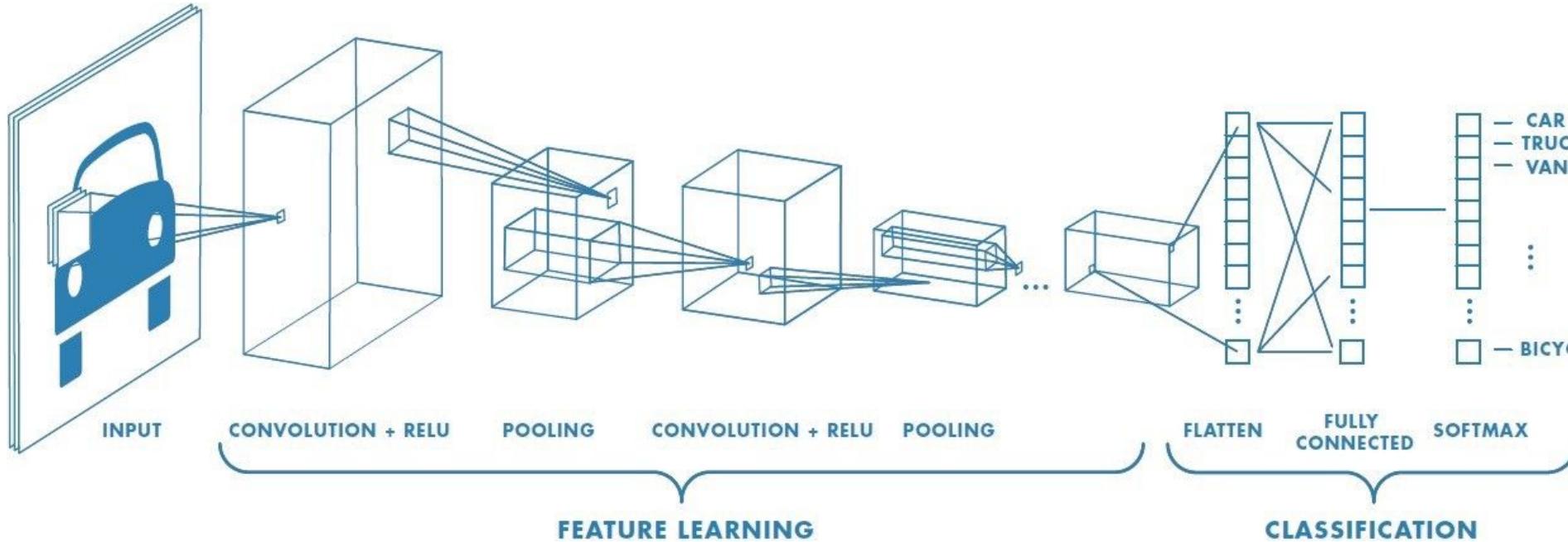


CNN Architectures



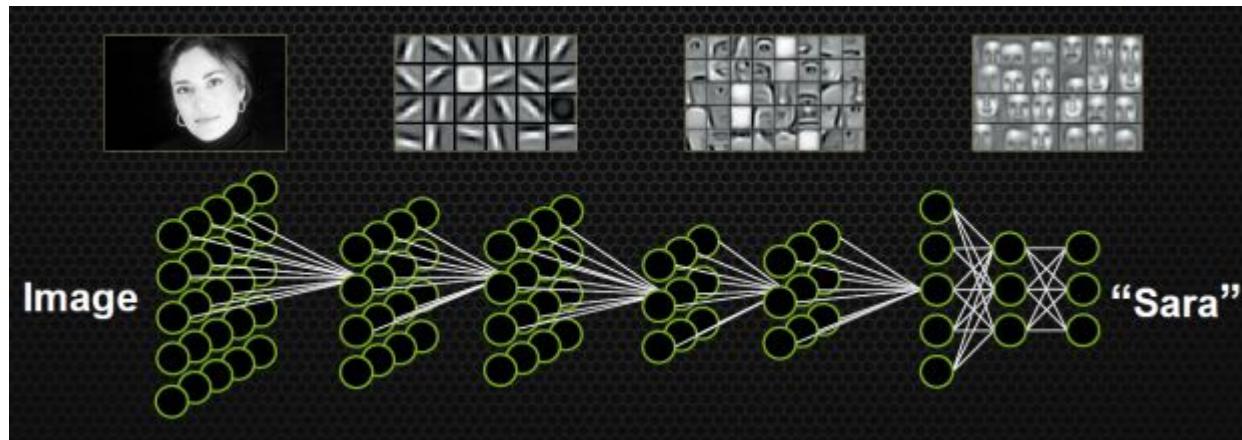
CNN Parameters

- Kernel size
- Number of kernels
 - Earlier layers (near the input) tend to have fewer filters while higher layers have more. This is done because the feature map size decreases with depth To equalize computation at each layer, the product of feature values with pixel position is kept roughly constant across layers. :
 $32 * (128*128) = 64 * (90*90)$
 - Activation Number: number of feature maps times number of pixel positions
- Padding
 - For example, for the convolutional kernels of size 3x3 would receive a 2-pixel pad (1 pixel on each side of the image)
 - Without this feature, the width of the representation shrinks by one pixel less than the kernel width at each layer.
 - Without zero padding, we are forced to choose between shrinking the spatial extent of the network rapidly and using small kernels—both scenarios that significantly limit the expressive power of the network.
- Stride
- Pooling type: max, average
- Pooling size:
 - Max pooling downsamples image greatly. We usually use a 2×2 max-pooling For larger images, we can use 4×4 pooling in the lower layers
- Dilation
 - It tries to reduce the processing/memory without significant signal loss.
 - A dilation of 2 on a 3x3 kernel expands the kernel to a 5x5 window, while still processing 9 pixels (evenly spaced).



Visualizing CNN Feature Maps

- To visualize feature maps, we apply the learned kernel filters to a sample input image



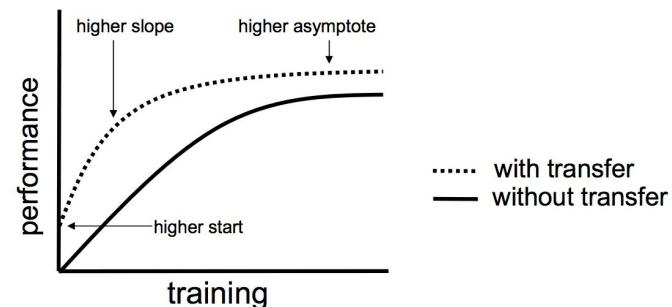
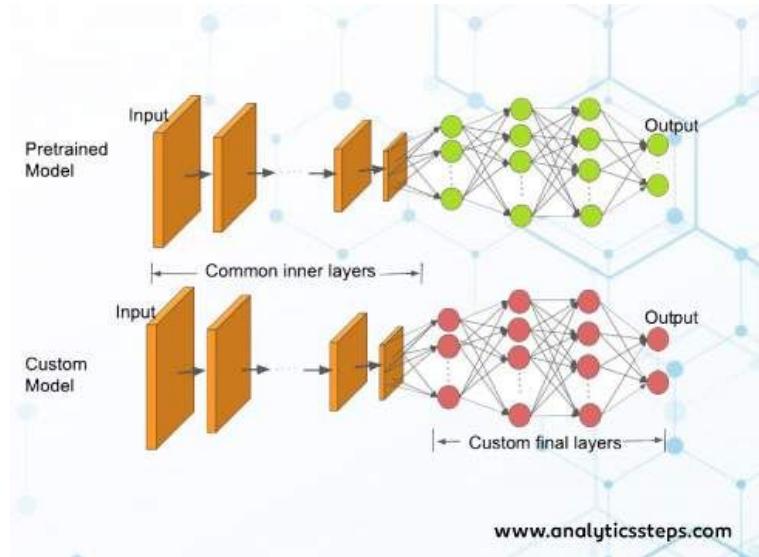
Demo: CNNs on MNIST Fashion

<http://localhost:8888/notebooks/jupyter-notebooks/chapman-generative-AI/NN-review-CNN-MNIST-Fashion.ipynb>

Pretrained Models and Transfer Learning

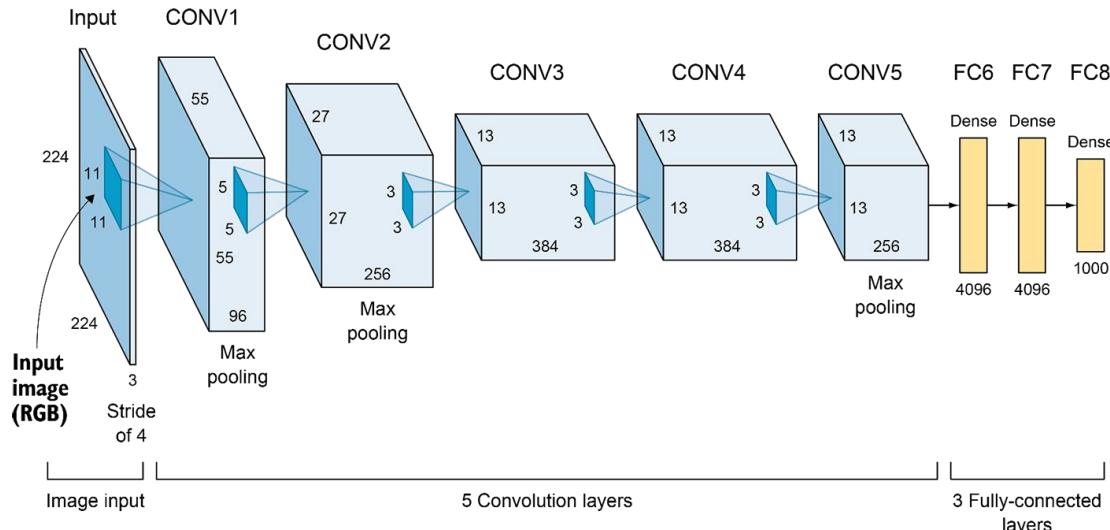
Transfer Learning

- We can also add a few new layers to the end of the existing model Using the data available for the new task, we tune the existing weights of the model to adapt for the new task and/or train the new added layers
- We might want to freeze some or all layers borrowed from the existing pre-trained model to prevent those weights from being updated
- Fine-tuning of the Whole Model

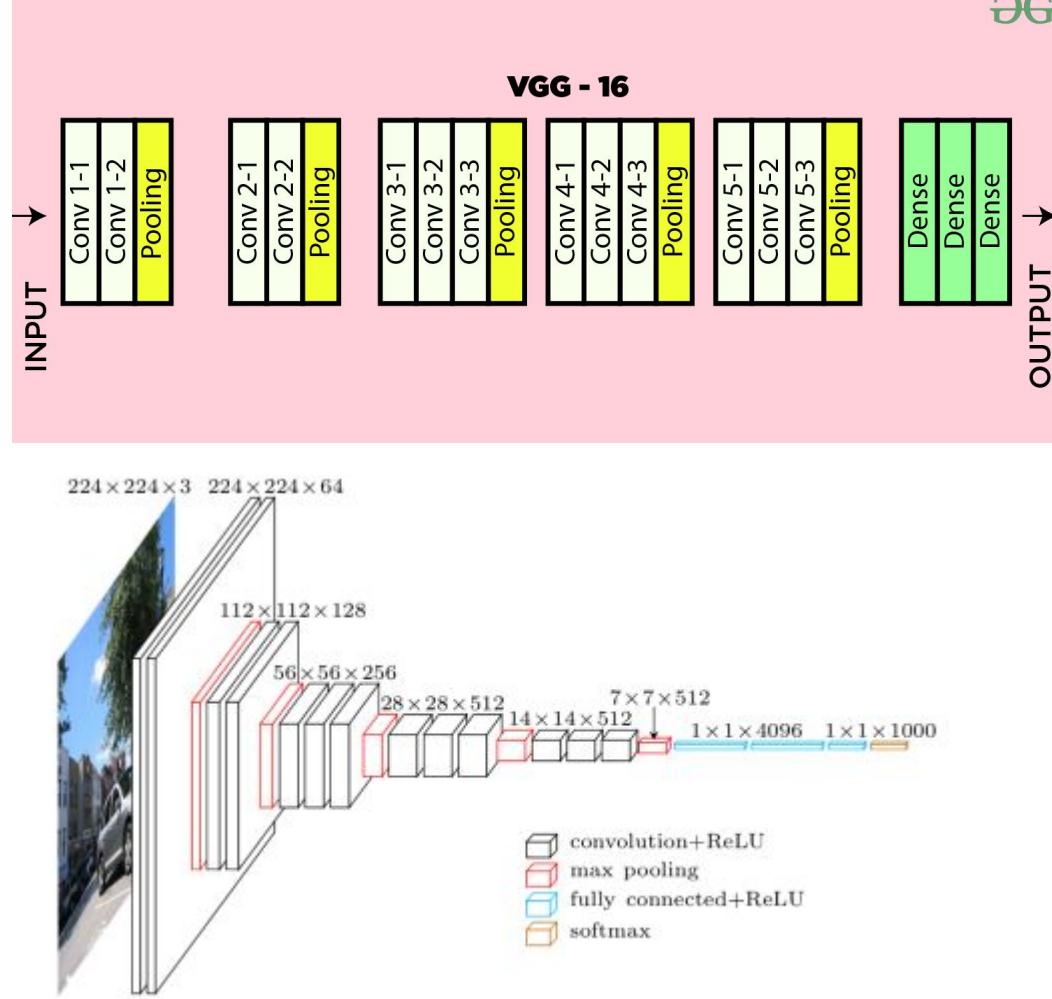


AlexNet

- It achieved a top-5 error of 15.3% in ImageNet
- AlexNet architecture consists of 5 convolutional layers, 3 max-pooling layers, 2 normalization layers, 2 fully connected layers, and 1 softmax layer.
- ReLU activation function is used in each convolutional layer.
- The input size is 224x224x3
- Has around 60 million parameters.

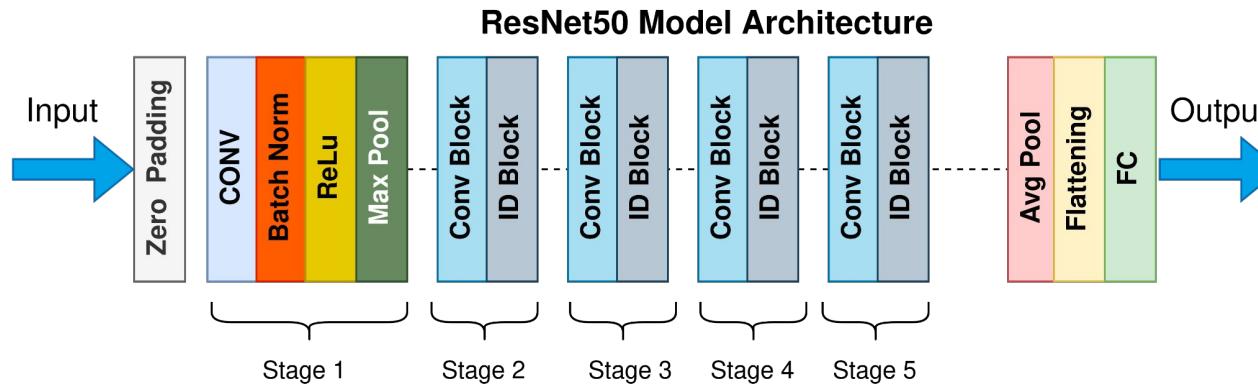
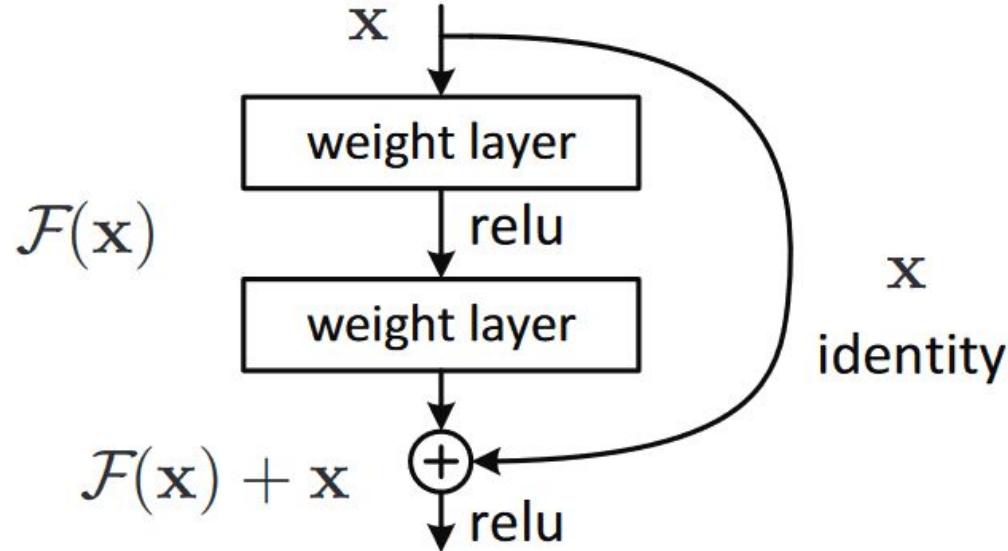


VGG16



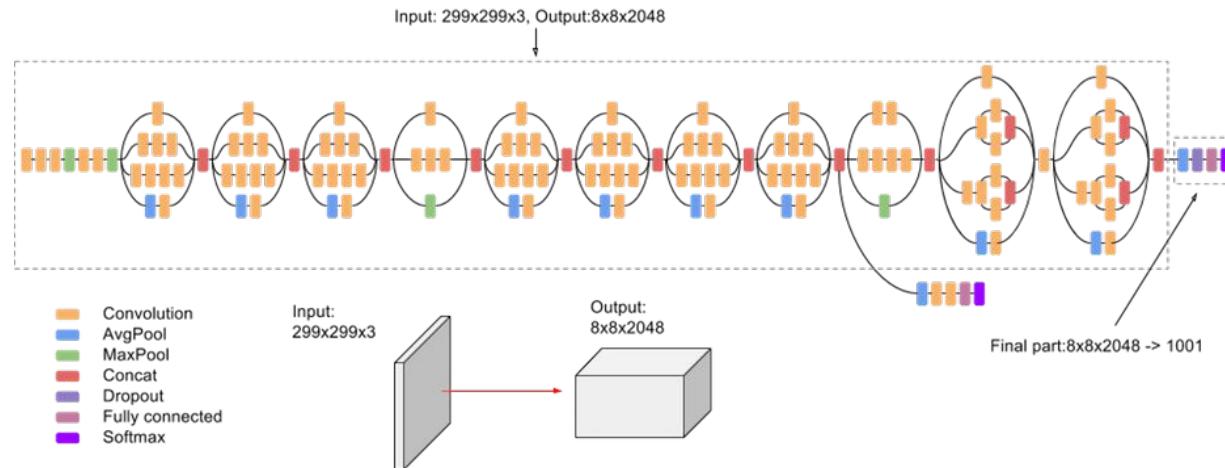
Skip Connection

- layers feed into the next layer and also directly into the layers 1, 2, or 3 or more hops away.
- These are called residual blocks
- ResNets are build using these residual blocks



Inception V3

- The main idea behind the Inception model is to calculate multiple different conversions parallelly and concatenate them as a single output to be used by the next layer.
 - For instance, convolutions of sizes 1×1 , 3×3 , 5×5 and also a 3×3 max-pooling are done in parallel on the input, and the outputs from all of these operations are stacked together to generate the output of the model.



Demo: Transfer Learning in Computer Vision

<http://localhost:8888/notebooks/jupyter-notebooks/chapman-generative-AI/NN-review-Transfer-Learning.ipynb>

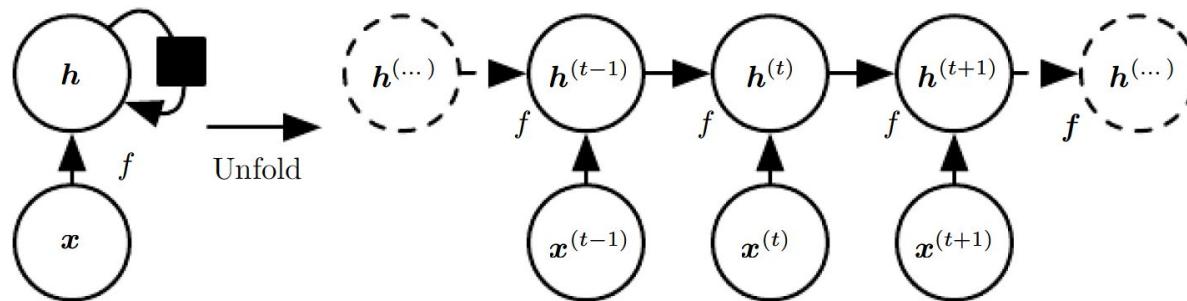
RNN

RNN

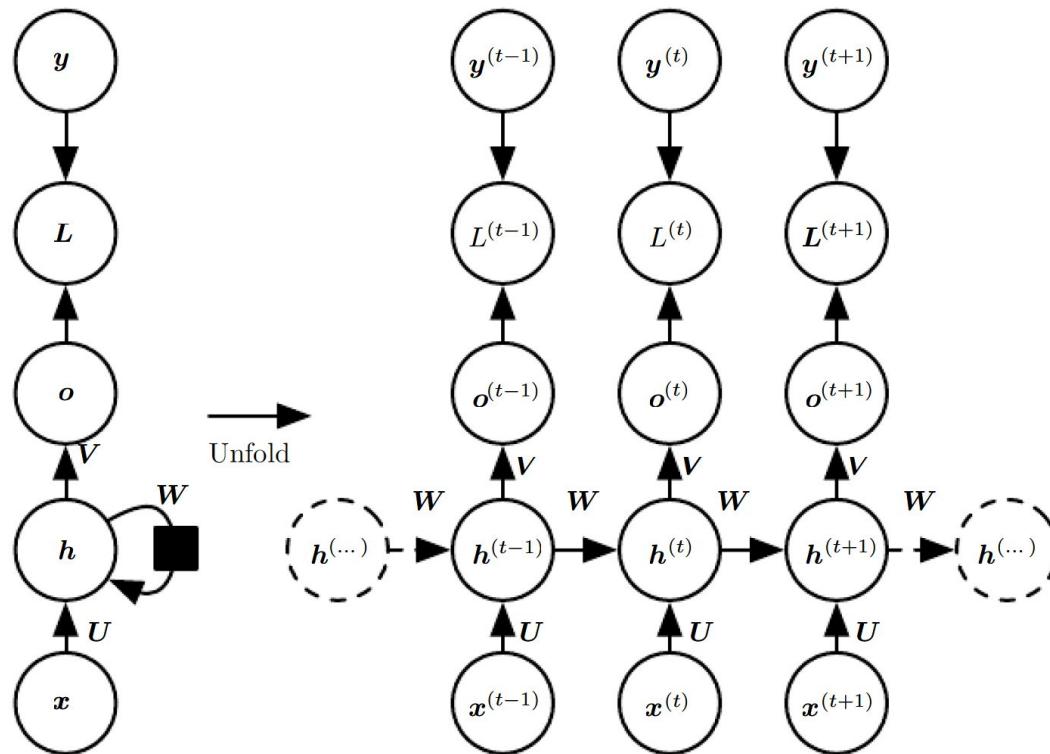
- Specialized for processing a sequence of values x_1, \dots, x_n such as words in a text
- RNNs have “memory” (store information from the previous data points in the input)
- Parameter sharing makes it possible to extend and apply the model to samples of different lengths
 - An RNN shares the same weights across several time steps
- Applications
 - Text autofill
 - Music generation
 - Speech recognition
 - Text analysis
 - Text summarization
 - RNNs can be combined with CNNs for tasks such as image to text

Unfolding the Graph

- Two ways to show a Recurrent relationship

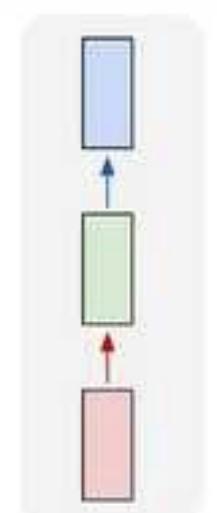


- The black square indicates a delay of a single time step.

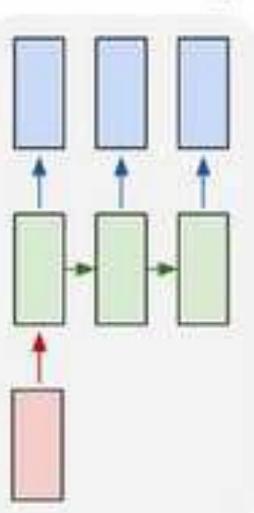


RNN Types Based on the Input and Output

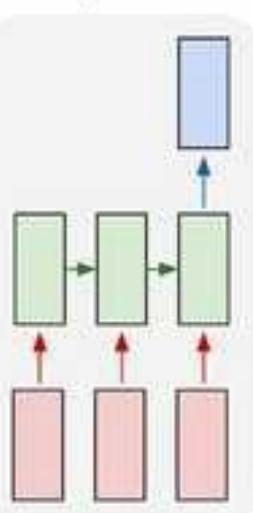
one to one



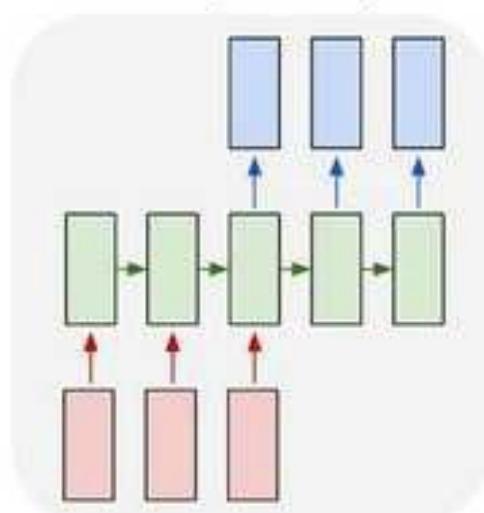
one to many



many to one



many to many



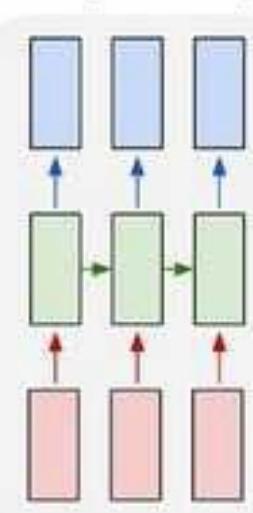
Traditional neural networks

music generation

sentiment analysis

machine translation

many to many



LSTM

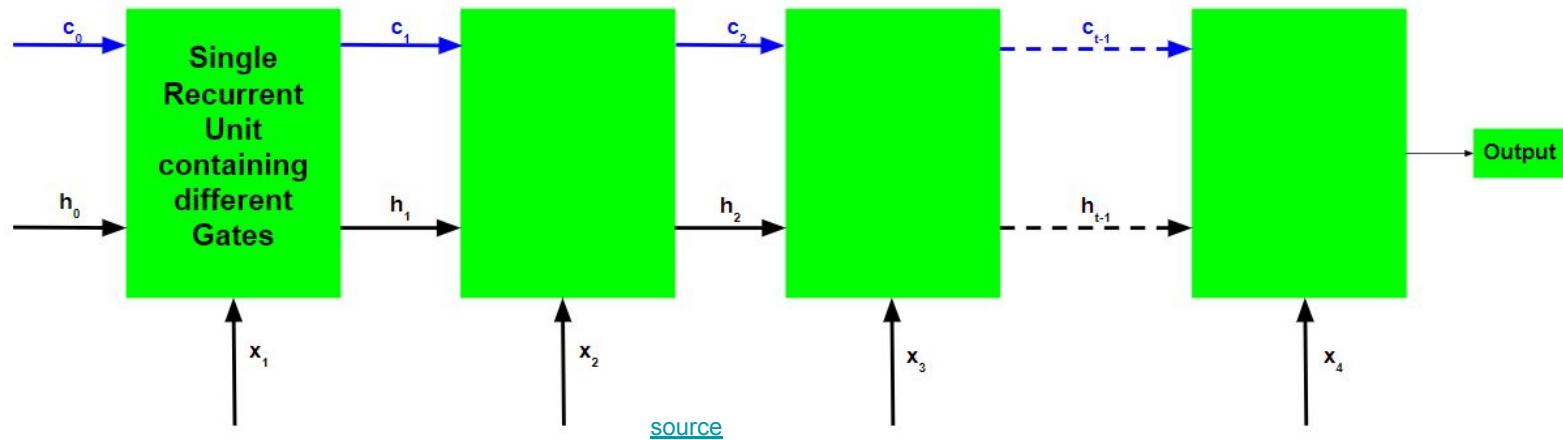
The Vanishing/Exploding Gradient Problem

- In a deep network, the gradients are propagated back in time all the way to the initial layer
- Since, we apply the chain rule in gradient calculation, gradients coming from deeper layers go through multiple matrix multiplications.
 - Thus, if they have small values (gradient<1) they diminish exponentially till they vanish (vanishing gradient)
 - On the other hand, if they have large values (gradient>1), then they grow exponentially large (exploding gradient)
- The long-term gradients in RNN which are back-propagated can "vanish" because of the computations involved in the process, which use finite-precision numbers
- Exploding/Vanishing gradient causes troubles in model training

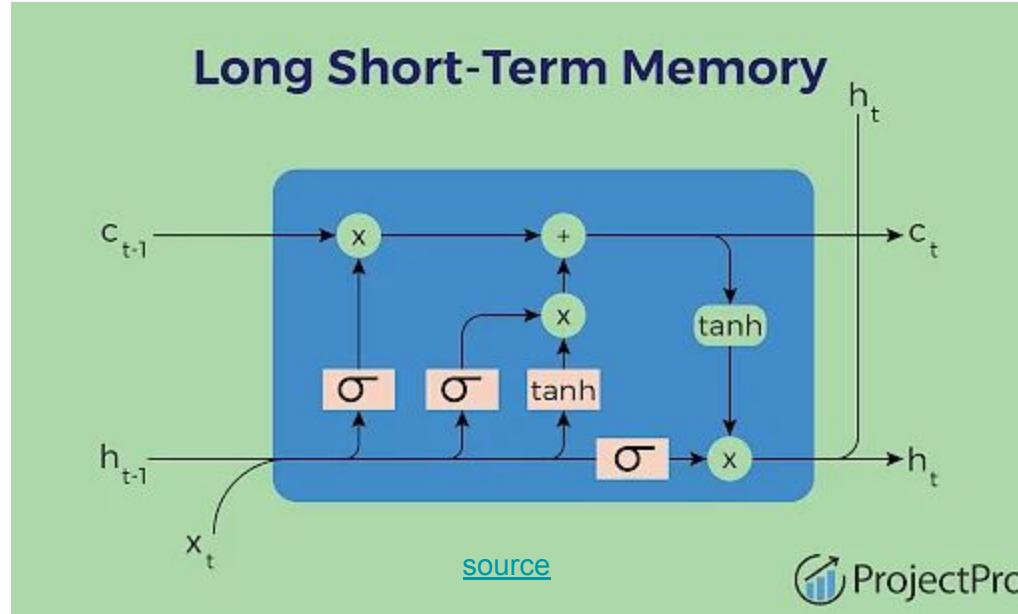
Long Short Term Memory (LSTM)

- Addresses the vanishing gradient problem in RNNs
- Use three types of gates: input, output, and forget
 - These gates decide which information to retain similar to GRU
- The recurrent units in an LSTM tries to **remember** all the past knowledge that the network has seen so far and to **forget** irrelevant data.
- LSTMs are the most popular variation of RNNs.
- Each LSTM recurrent unit has an additional vector called *Internal Cell State* which holds the information chosen to be remembered by the previous LSTM recurrent unit.
 - The internal cell state is also passed forward in addition to the hidden state

- Each unit has three inputs previous cell state (c_{t-1}), previous hidden state (h_{t-1}), and the current input (x_t)
- In addition to the unit's possible output (o_t), each unit has two outputs c_t and h_t



One LSTM Unit



Demo: Spam Detection Using LSTM

<http://localhost:8888/notebooks/jupyter-notebooks/chapman-generative-AI/NN-review-RNN-LSTM.ipynb>