



上海大学未来技术学院
SCHOOL OF FUTURE TECHNOLOGY, SHANGHAI UNIVERSITY

上海大学人工智能研究院
INSTITUTE OF ARTIFICIAL INTELLIGENCE, SHANGHAI UNIVERSITY

人工智能导论

——第1课：人工智能基本理论

叶林奇

未来技术学院（人工智能研究院）

2024秋季学期





课程简介

课程内容

- 第1周, 绪论: 人工智能基本理论
- 第2周, 强化学习
- 第3周, 机器人智能控制
- 第4周, 生成式AI
- 第5周, 大模型
- 第6周, 智能体
- 第7周, 具身智能
- 第8周, AI for Science
- 第9周, 项目汇报, 学生分组完成一个人工智能项目。
- 第10周, 项目汇报, 学生分组完成一个人工智能项目。

考核方式: 平时成绩30%+项目汇报70%

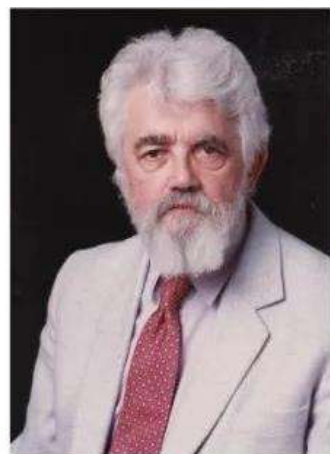


人工智能基本理论

1956年：**麦卡锡**召集哈佛大学、麻省理工学院、IBM公司、贝尔实验室的研究人员召开**达特茅斯会议**正式提出“**人工智能**”



2006年达特茅斯会议当事人重聚，左起：**摩尔、麦卡锡、明斯基、塞弗里奇、所罗门诺夫**。

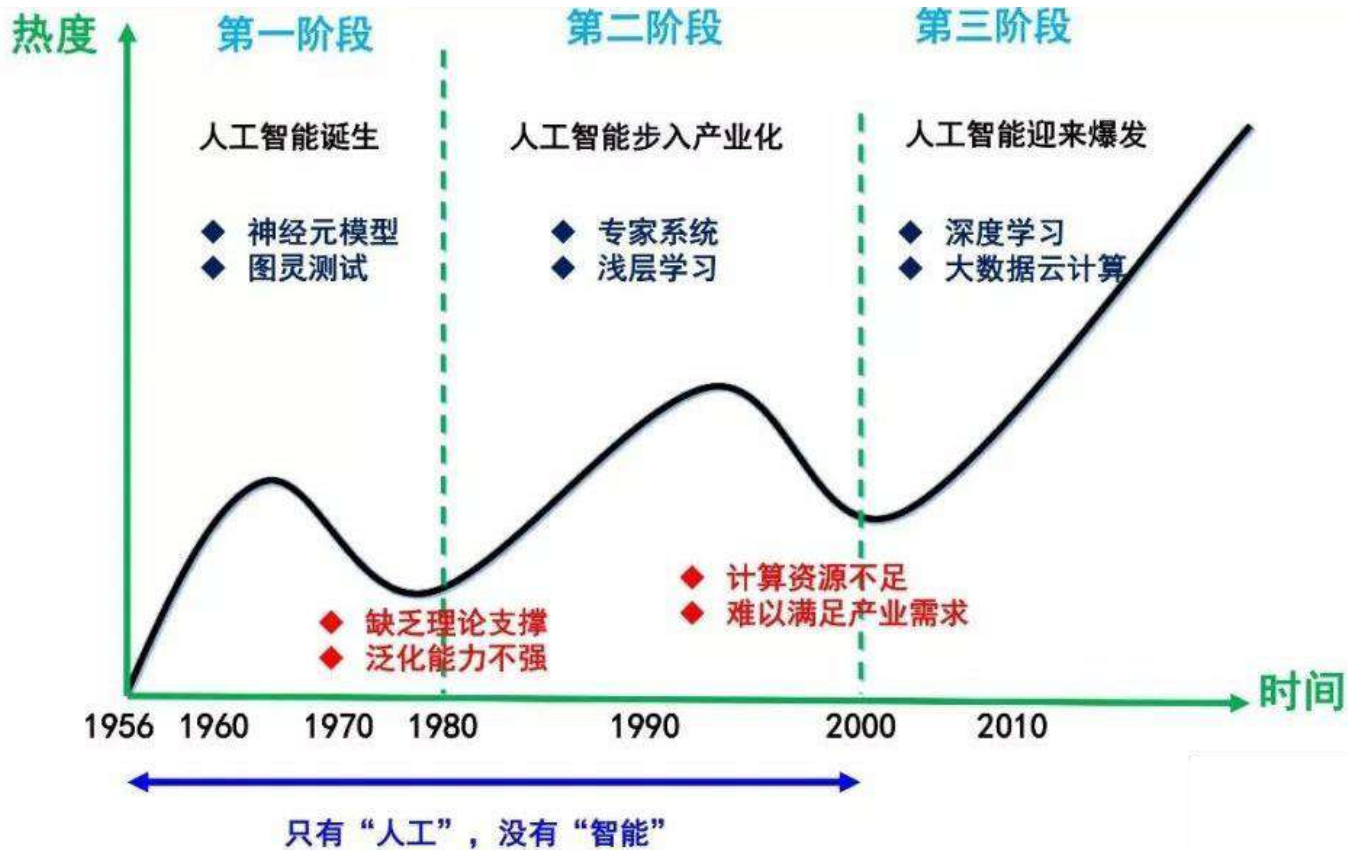


John McCarthy
人工智能之父

制造智能机器的科学和工程，特别是智能计算机程序。它与使用计算机理解人类智能的类似任务有关，但人工智能并不局限于生物学观察。



人工智能基本理论





人工智能基本理论



2006年，Hinton和他的学生Salakhutdinov在《科学》上发表了一篇文章，开启了深度学习在学术界和工业界的浪潮。

504 28 JULY 2006 VOL 313 SCIENCE
Reducing the Dimensionality of
Data with Neural Networks
G. E. Hinton* and R. R. Salakhutdinov

High-dimensional data can be converted to low-dimensional codes by training a multilayer neural network with a small central layer to reconstruct high-dimensional input vectors. Gradient descent can be used for fine-tuning the weights in such “autoencoder” networks, but this works well only if the initial weights are close to a good solution. We describe an effective way of initializing the weights that allows deep autoencoder networks to learn low-dimensional codes that work much better than principal components analysis as a tool to reduce the dimensionality of data



What is Deep Learning?

ARTIFICIAL INTELLIGENCE

Any technique that enables computers to mimic human behavior



MACHINE LEARNING

Ability to learn without explicitly being programmed



DEEP LEARNING

Extract patterns from data using neural networks

3 1 3 4 7 2
1 7 4 2 3 5

Teaching computers how to **learn a task** directly from **raw data**

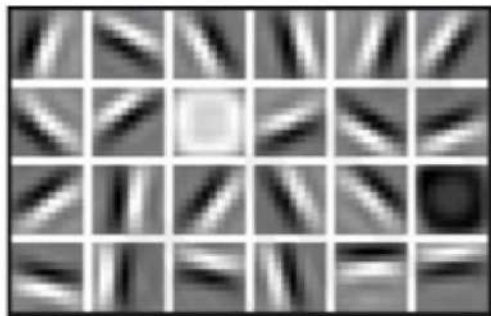


Why Deep Learning?

Hand engineered features are time consuming, brittle, and not scalable in practice

Can we learn the **underlying features** directly from data?

Low Level Features



Lines & Edges

Mid Level Features



Eyes & Nose & Ears

High Level Features



Facial Structure

提纲

一、前馈网络

二、循环网络

三、卷积网络



上海大学
SHANGHAI UNIVERSITY

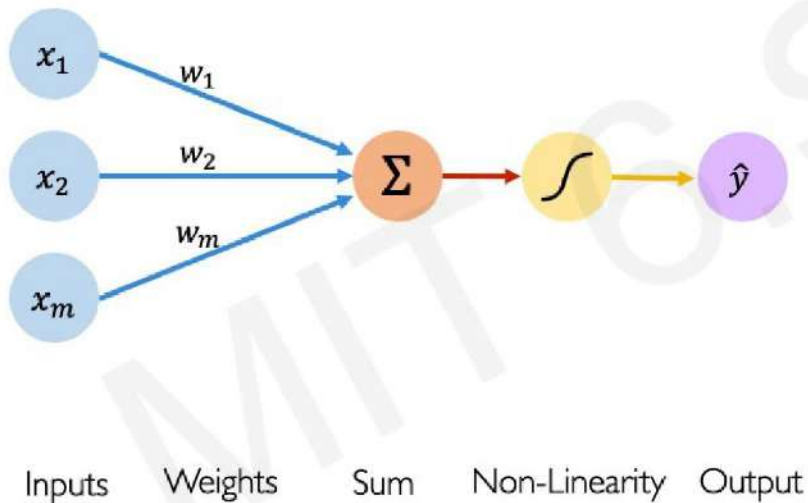


The Perceptron

The structural building block of deep learning



The Perceptron: Forward Propagation



Linear combination of inputs

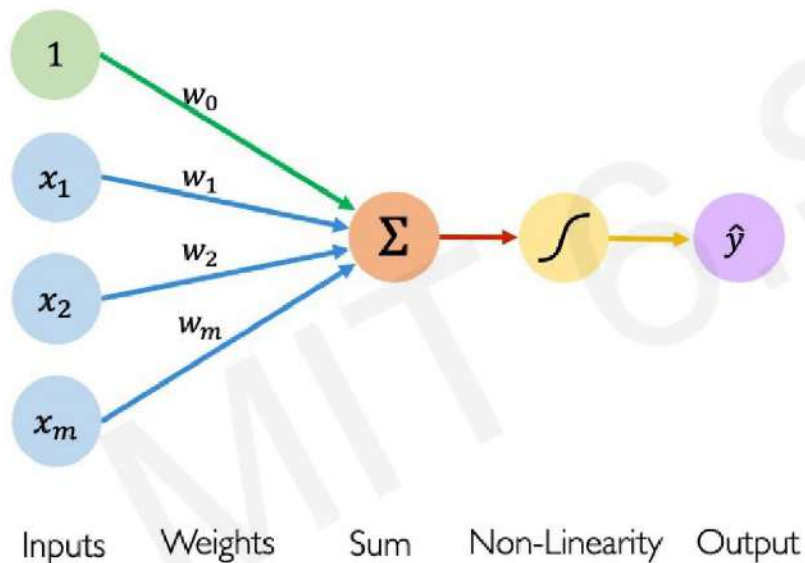
$$\hat{y} = g \left(\sum_{i=1}^m x_i w_i \right)$$

Output

Non-linear activation function



The Perceptron: Forward Propagation



Output

Linear combination of inputs

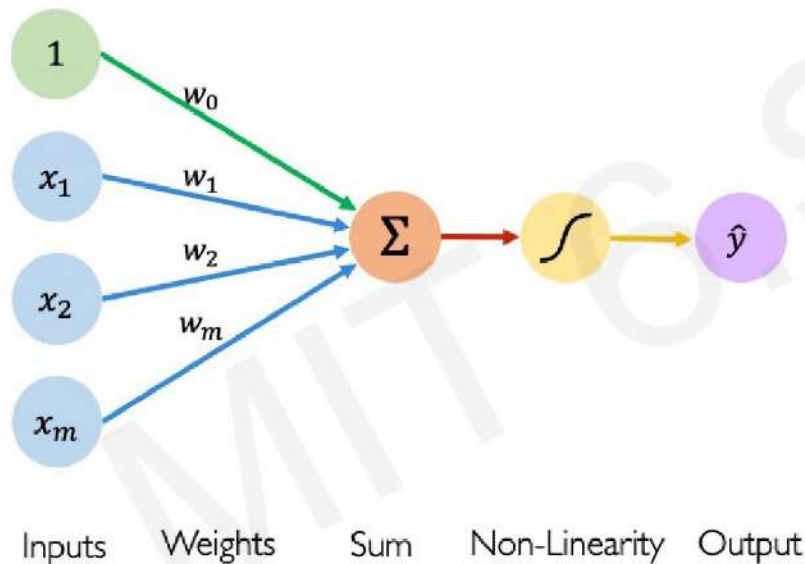
$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

Non-linear activation function

Bias



The Perceptron: Forward Propagation



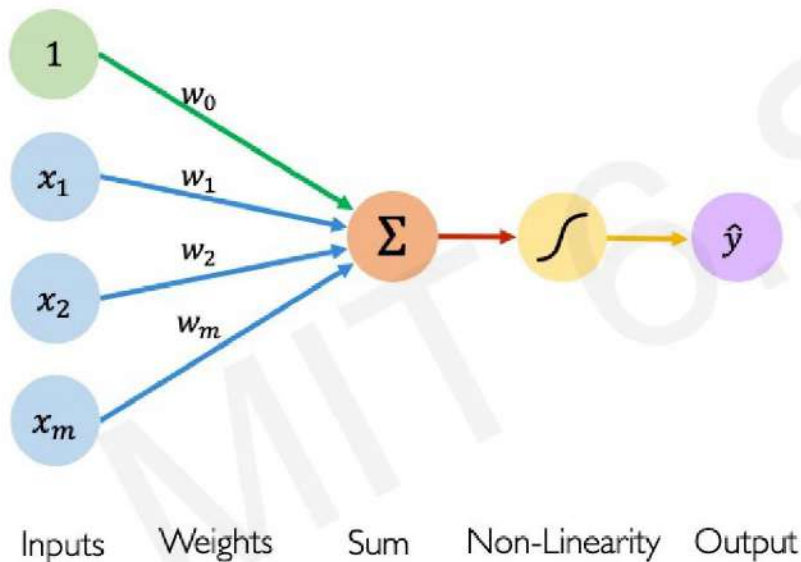
$$\hat{y} = g \left(w_0 + \sum_{i=1}^m x_i w_i \right)$$

$$\hat{y} = g (w_0 + \mathbf{X}^T \mathbf{W})$$

where: $\mathbf{X} = \begin{bmatrix} x_1 \\ \vdots \\ x_m \end{bmatrix}$ and $\mathbf{W} = \begin{bmatrix} w_1 \\ \vdots \\ w_m \end{bmatrix}$



The Perceptron: Forward Propagation

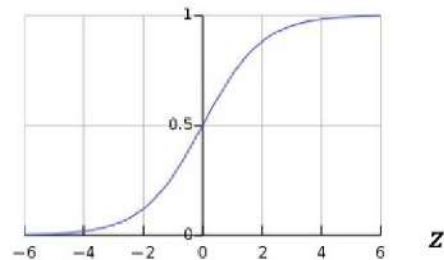


Activation Functions

$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$

- Example: sigmoid function

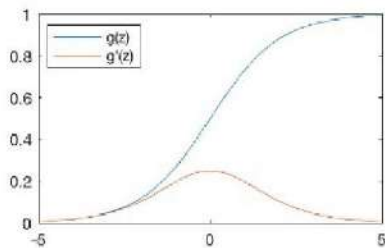
$$g(z) = \sigma(z) = \frac{1}{1 + e^{-z}}$$






Common Activation Functions

Sigmoid Function

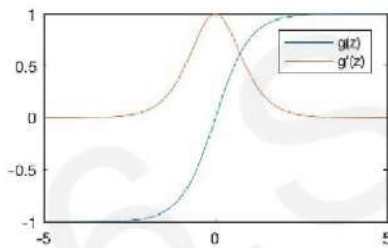


$$g(z) = \frac{1}{1 + e^{-z}}$$

$$g'(z) = g(z)(1 - g(z))$$


 `tf.math.sigmoid(z)`

Hyperbolic Tangent

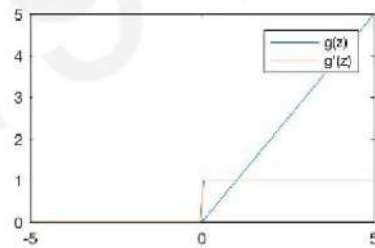


$$g(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$$

$$g'(z) = 1 - g(z)^2$$

 `tf.math.tanh(z)`

Rectified Linear Unit (ReLU)



$$g(z) = \max(0, z)$$

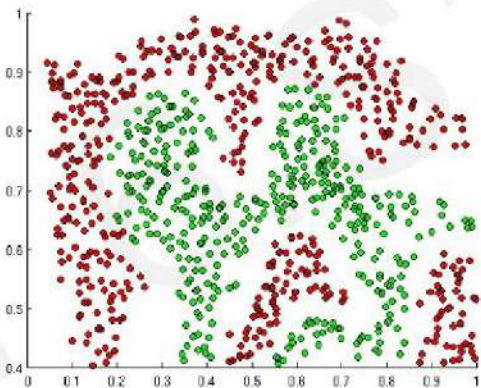
$$g'(z) = \begin{cases} 1, & z > 0 \\ 0, & \text{otherwise} \end{cases}$$

 `tf.nn.relu(z)`



Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network

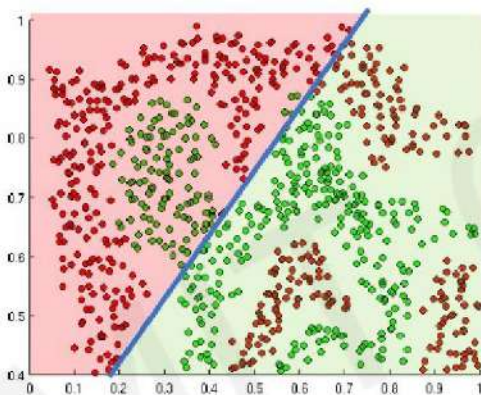


What if we wanted to build a neural network to distinguish green vs red points?



Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network

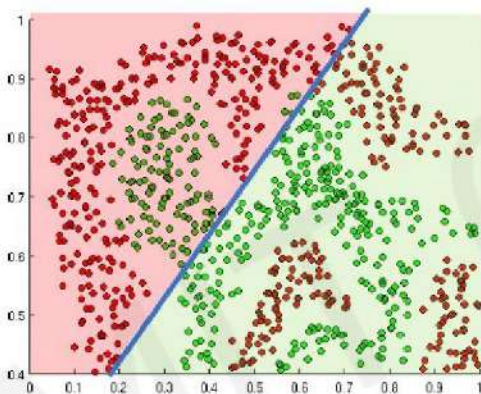


Linear activation functions produce linear decisions no matter the network size

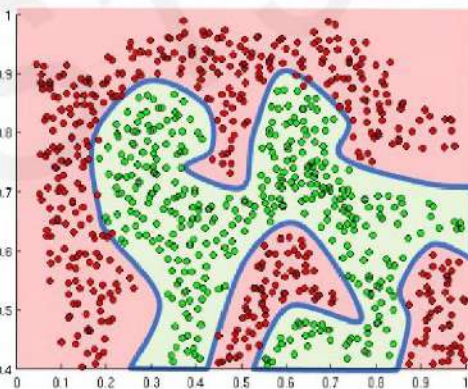


Importance of Activation Functions

The purpose of activation functions is to **introduce non-linearities** into the network



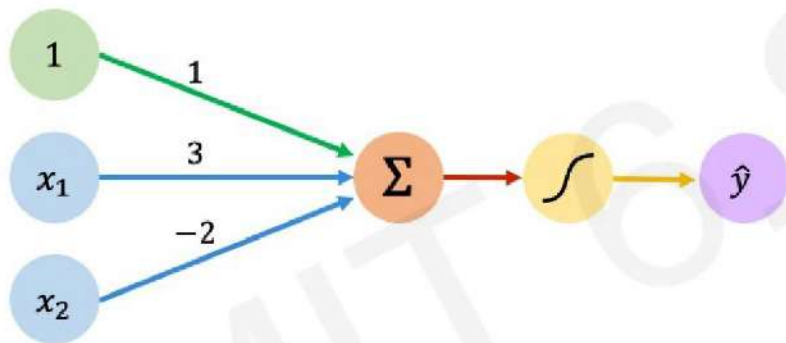
Linear activation functions produce linear decisions no matter the network size



Non-linearities allow us to approximate arbitrarily complex functions



The Perceptron: Example



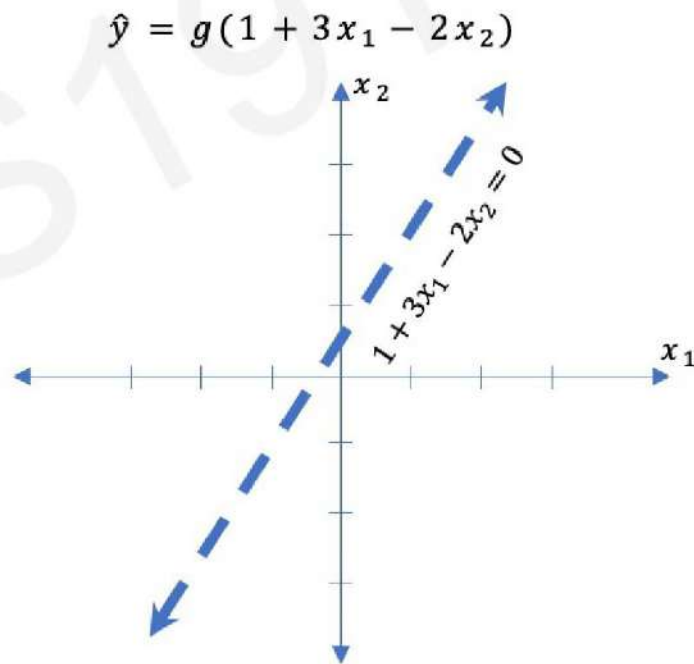
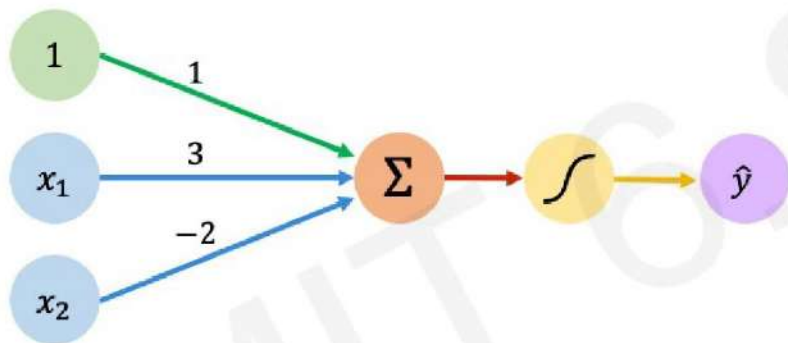
We have: $w_0 = 1$ and $\mathbf{W} = \begin{bmatrix} 3 \\ -2 \end{bmatrix}$

$$\begin{aligned}\hat{y} &= g(w_0 + \mathbf{X}^T \mathbf{W}) \\ &= g\left(1 + \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^T \begin{bmatrix} 3 \\ -2 \end{bmatrix}\right) \\ \hat{y} &= g(1 + 3x_1 - 2x_2)\end{aligned}$$

This is just a line in 2D!

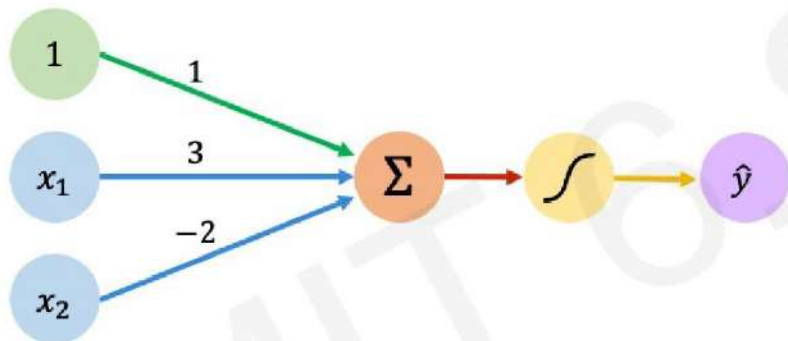


The Perceptron: Example



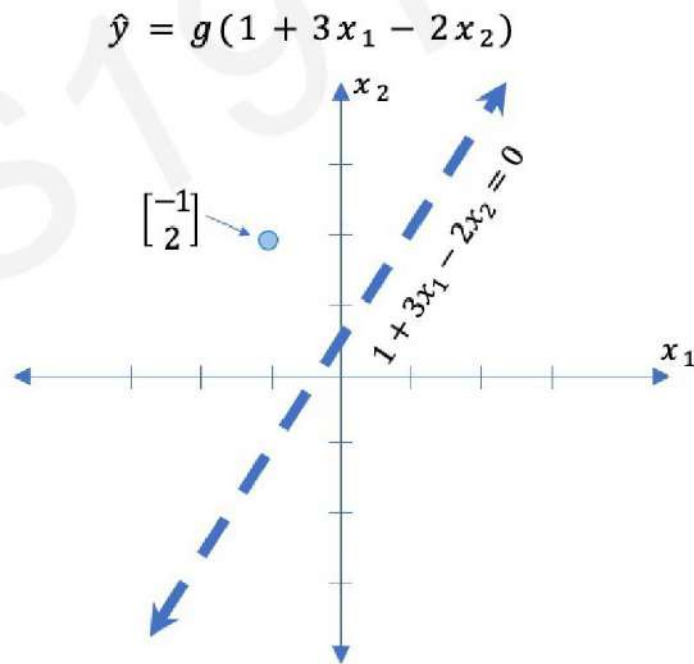


The Perceptron: Example



Assume we have input: $\mathbf{X} = \begin{bmatrix} -1 \\ 2 \end{bmatrix}$

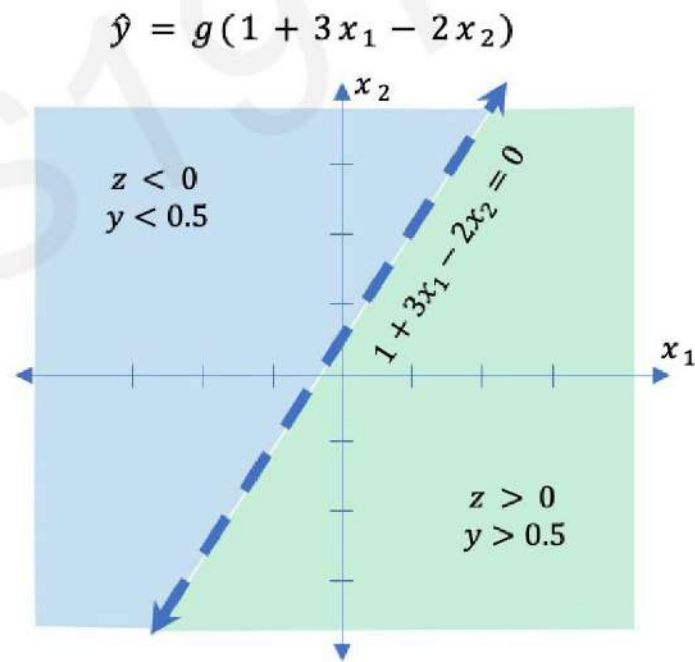
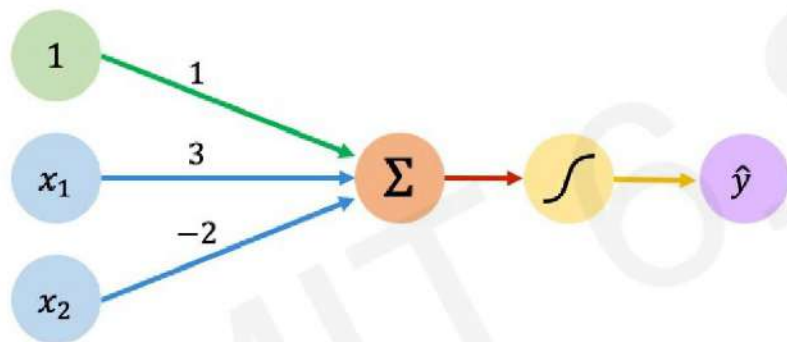
$$\begin{aligned}\hat{y} &= g(1 + (3 * -1) - (2 * 2)) \\ &= g(-6) \approx 0.002\end{aligned}$$



$$\hat{y} = g(1 + 3x_1 - 2x_2)$$



The Perceptron: Example





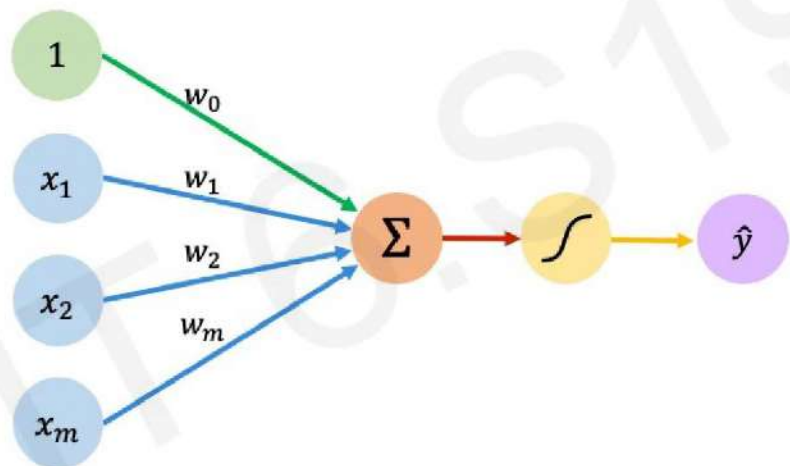
人工智能基本理论

Building Neural Networks with Perceptrons



The Perceptron: Simplified

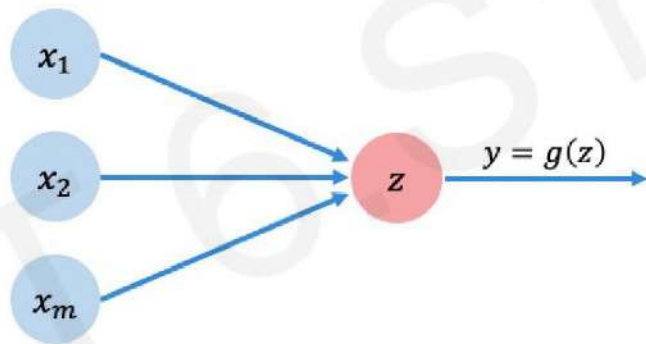
$$\hat{y} = g(w_0 + \mathbf{X}^T \mathbf{W})$$



Inputs Weights Sum Non-Linearity Output



The Perceptron: Simplified

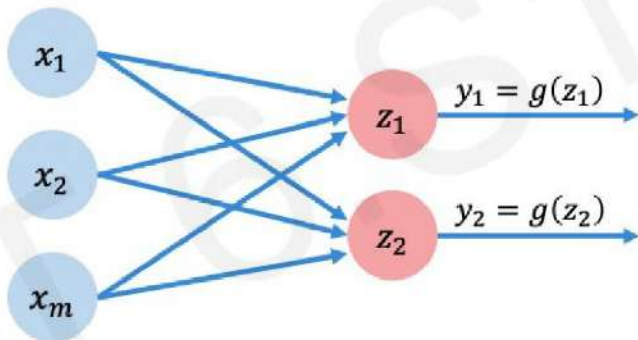


$$z = w_0 + \sum_{j=1}^m x_j w_j$$



Multi Output Perceptron

Because all inputs are densely connected to all outputs, these layers are called **Dense** layers

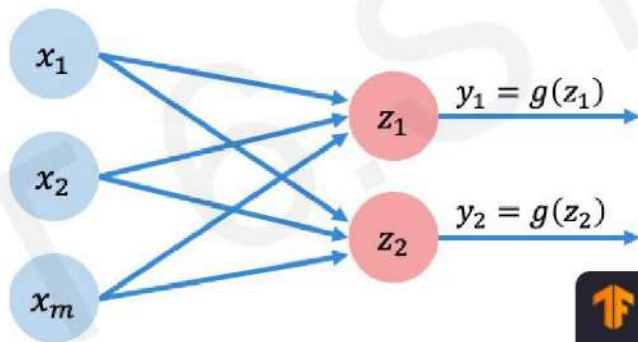


$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$



Multi Output Perceptron

Because all inputs are densely connected to all outputs, these layers are called **Dense** layers

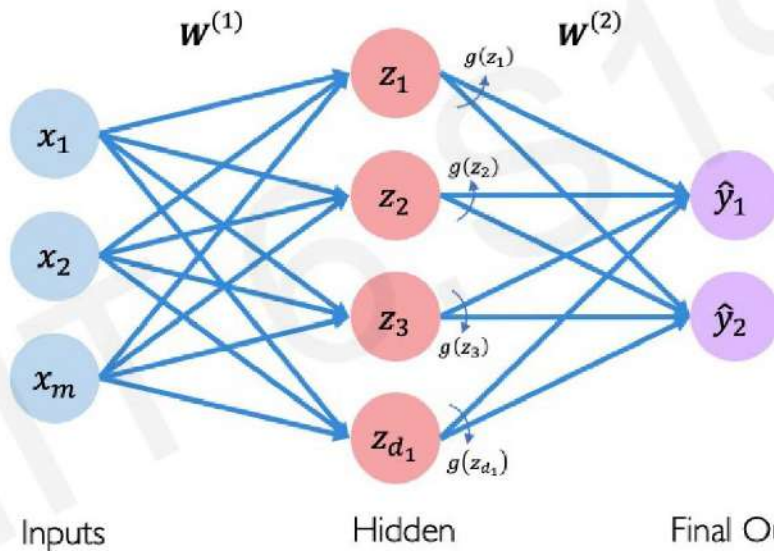


```
import tensorflow as tf
layer = tf.keras.layers.Dense(
    units=2)
```

$$z_i = w_{0,i} + \sum_{j=1}^m x_j w_{j,i}$$



Single Layer Neural Network

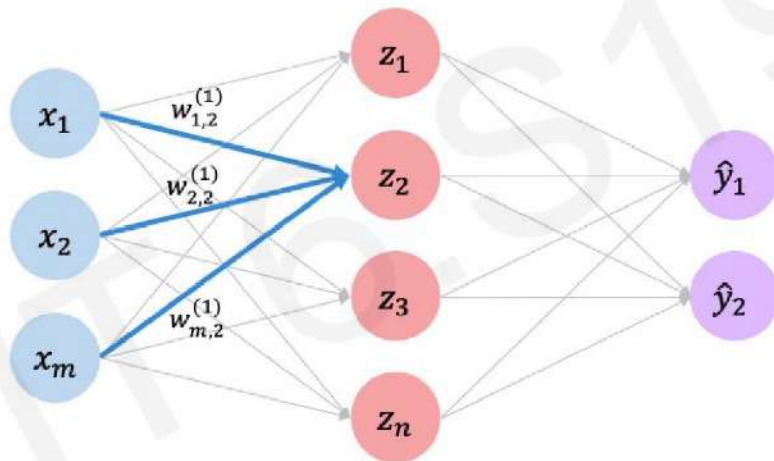


$$z_i = w_{0,i}^{(1)} + \sum_{j=1}^m x_j w_{j,i}^{(1)}$$

$$\hat{y}_i = g \left(w_{0,i}^{(2)} + \sum_{j=1}^{d_1} g(z_j) w_{j,i}^{(2)} \right)$$



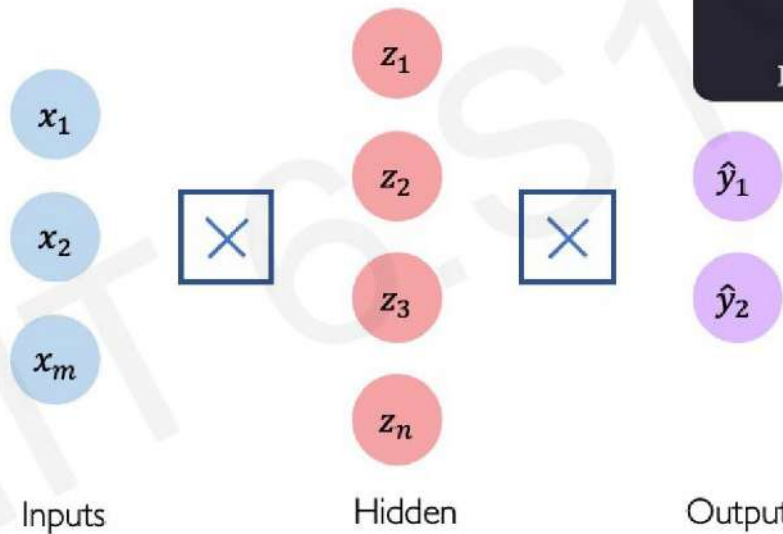
Single Layer Neural Network



$$\begin{aligned} z_2 &= w_{0,2}^{(1)} + \sum_{j=1}^m x_j w_{j,2}^{(1)} \\ &= w_{0,2}^{(1)} + x_1 w_{1,2}^{(1)} + x_2 w_{2,2}^{(1)} + x_m w_{m,2}^{(1)} \end{aligned}$$



Multi Output Perceptron

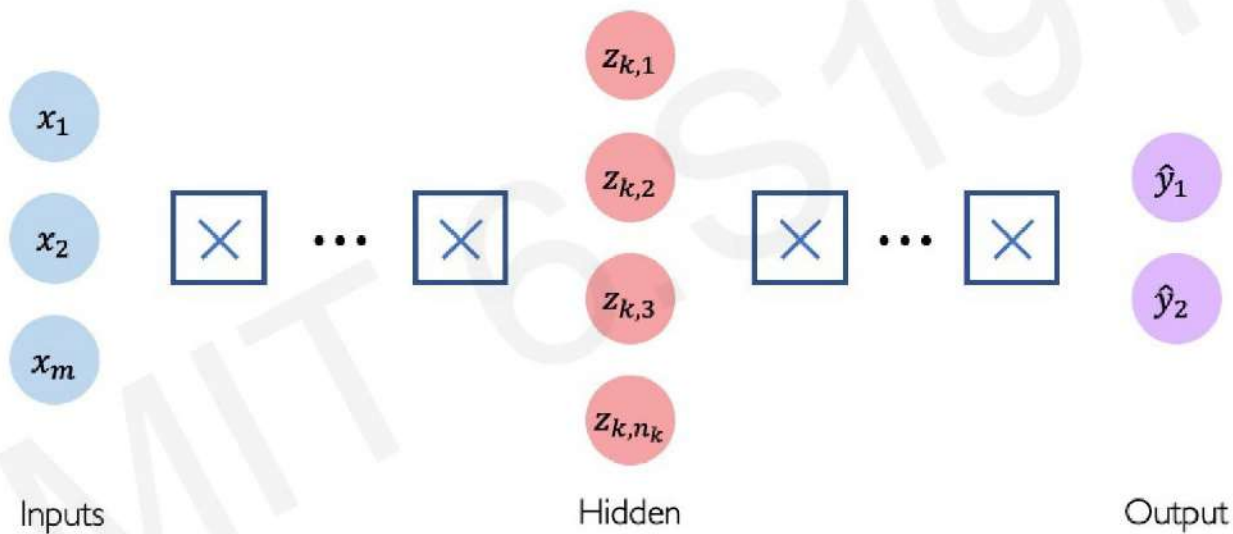


```
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(n),
    tf.keras.layers.Dense(2)
])
```



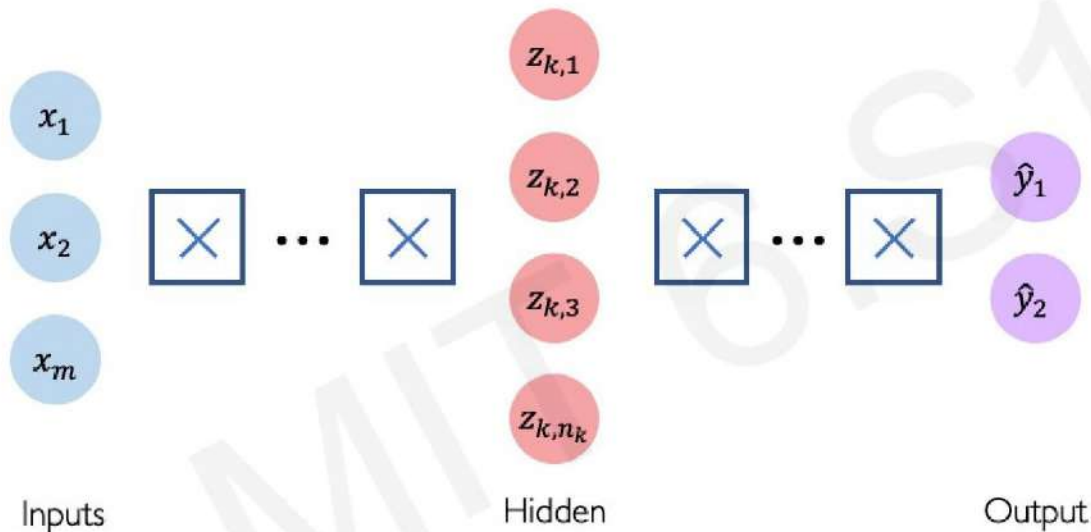
Deep Neural Network



$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$



Deep Neural Network



```
import tensorflow as tf

model = tf.keras.Sequential([
    tf.keras.layers.Dense(n1),
    tf.keras.layers.Dense(n2),
    :
    tf.keras.layers.Dense(2)
])
```

$$z_{k,i} = w_{0,i}^{(k)} + \sum_{j=1}^{n_{k-1}} g(z_{k-1,j}) w_{j,i}^{(k)}$$



人工智能基本理论

Applying Neural Networks



Example Problem

Will I pass this class?

Let's start with a simple two feature model

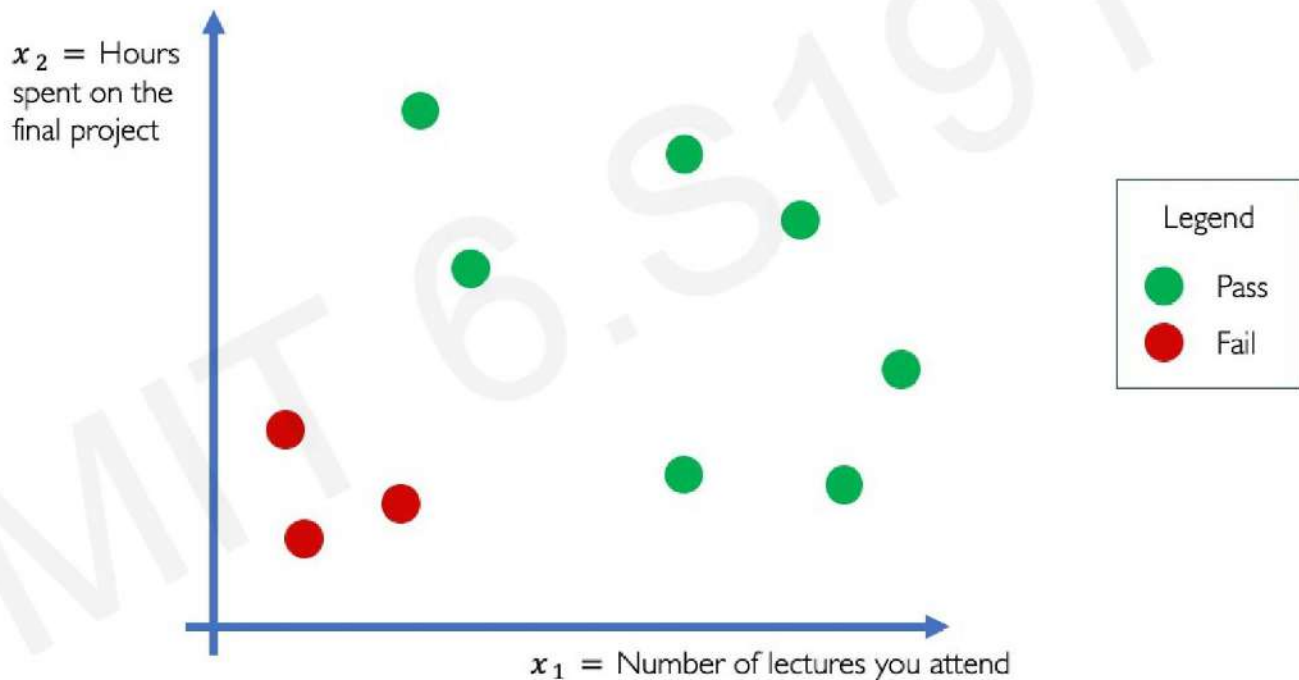
x_1 = Number of lectures you attend

x_2 = Hours spent on the final project



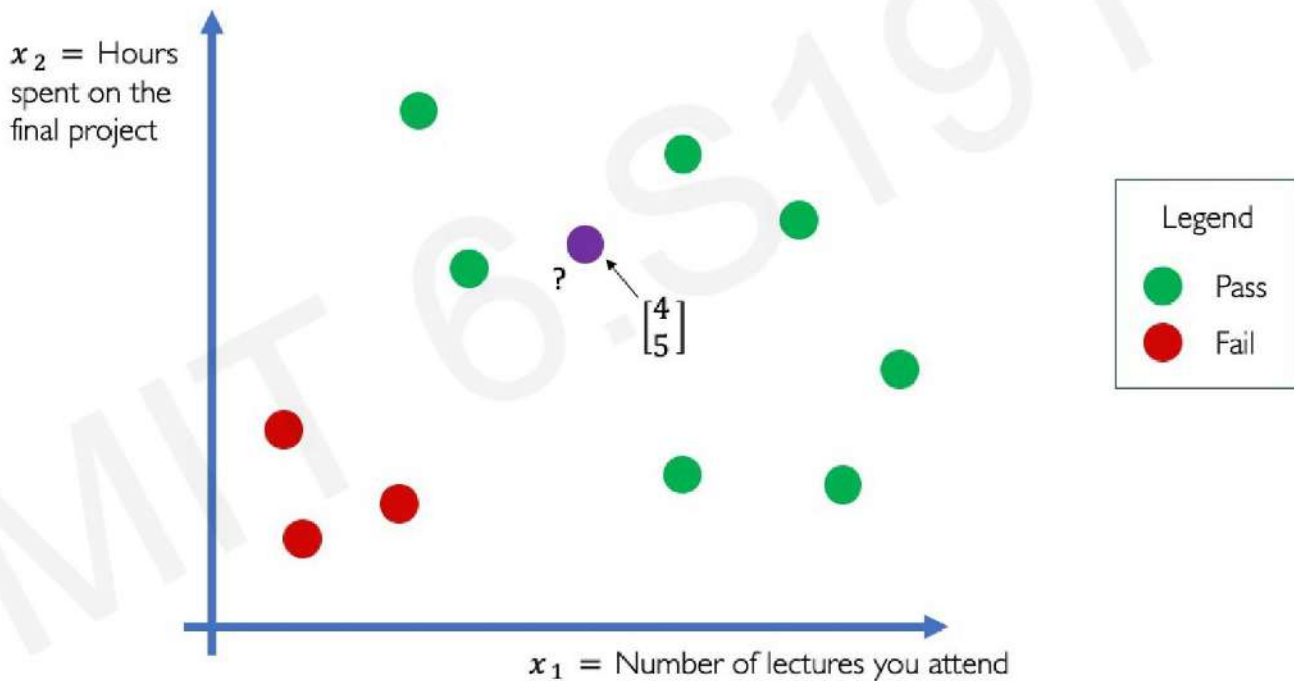
人工智能基本理论

Example Problem: Will I pass this class?





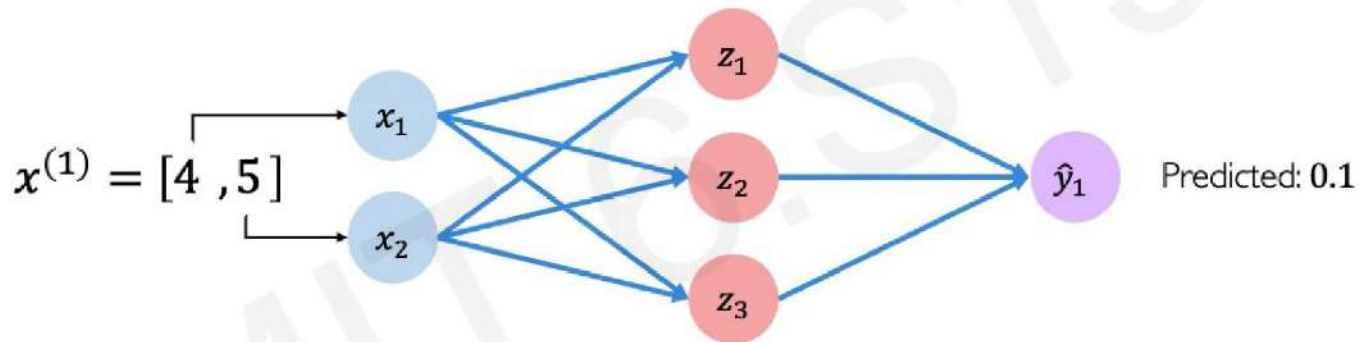
Example Problem: Will I pass this class?





人工智能基本理论

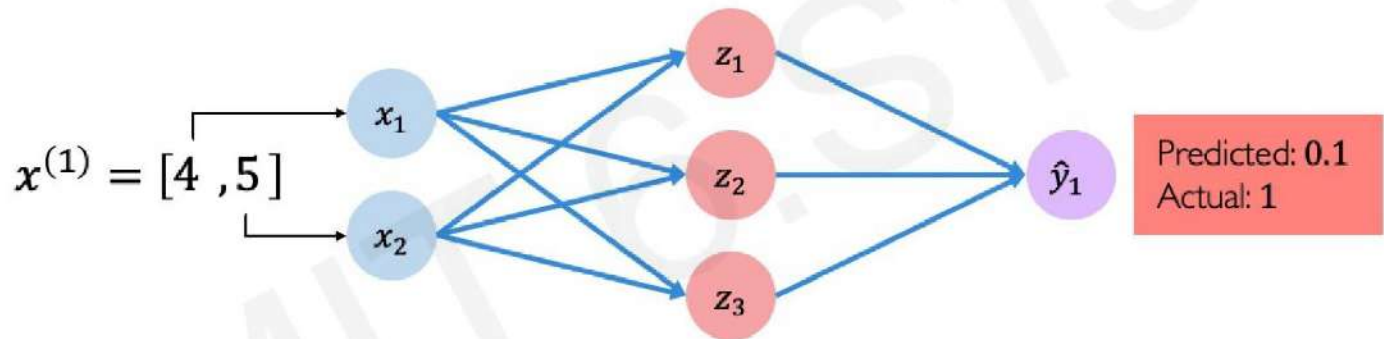
Example Problem: Will I pass this class?





人工智能基本理论

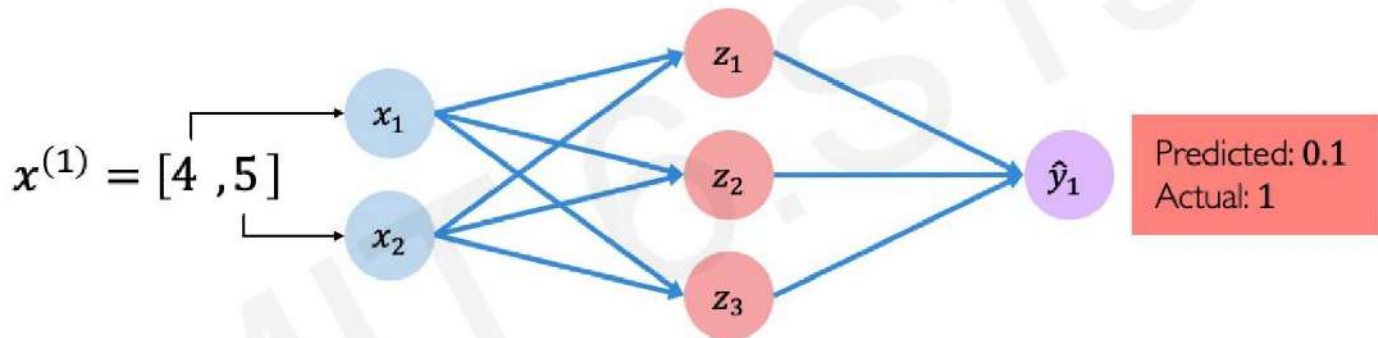
Example Problem: Will I pass this class?





Quantifying Loss

The **loss** of our network measures the cost incurred from incorrect predictions

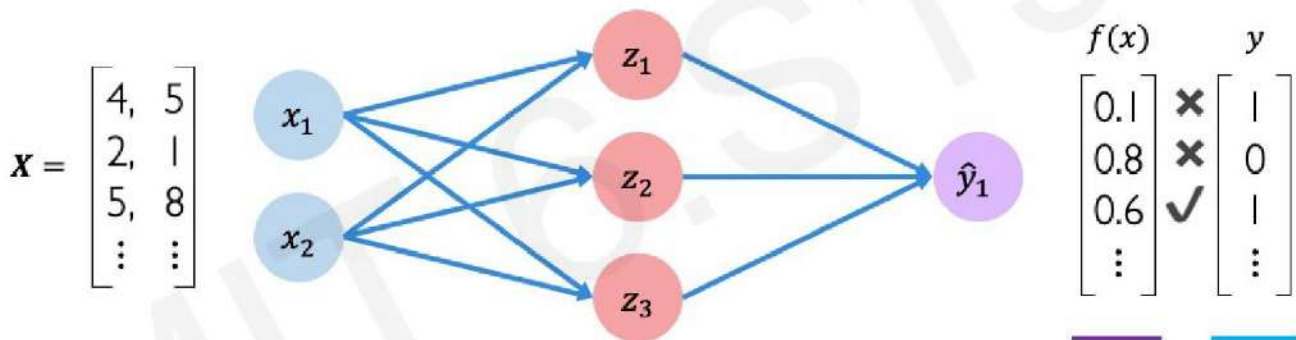


$$\mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$



Empirical Loss

The **empirical loss** measures the total loss over our entire dataset



Also known as:

- Objective function
- Cost function
- Empirical Risk

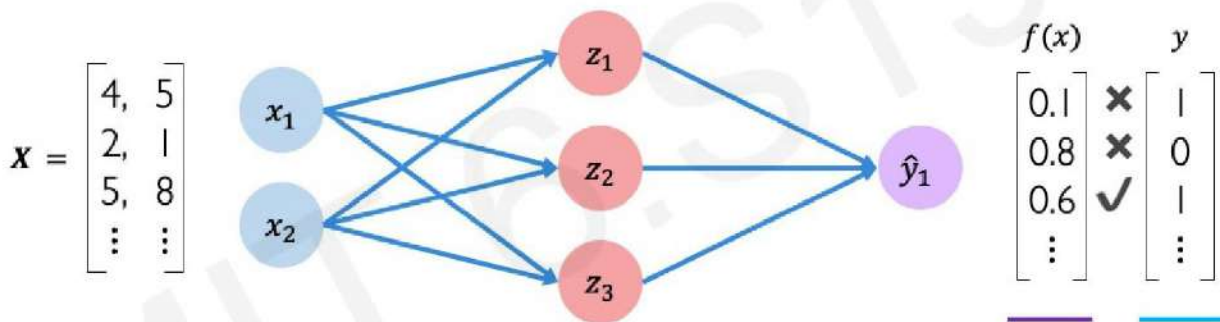
$$J(\mathbf{W}) = \frac{1}{n} \sum_{i=1}^n \mathcal{L}(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}}, \underbrace{y^{(i)}}_{\text{Actual}})$$

Predicted Actual



Binary Cross Entropy Loss

Cross entropy loss can be used with models that output a probability between 0 and 1



$$J(\mathbf{W}) = -\frac{1}{n} \sum_{i=1}^n \underbrace{y^{(i)}}_{\text{Actual}} \log \left(\underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}} \right) + (1 - \underbrace{y^{(i)}}_{\text{Actual}}) \log \left(1 - \underbrace{f(x^{(i)}; \mathbf{W})}_{\text{Predicted}} \right)$$

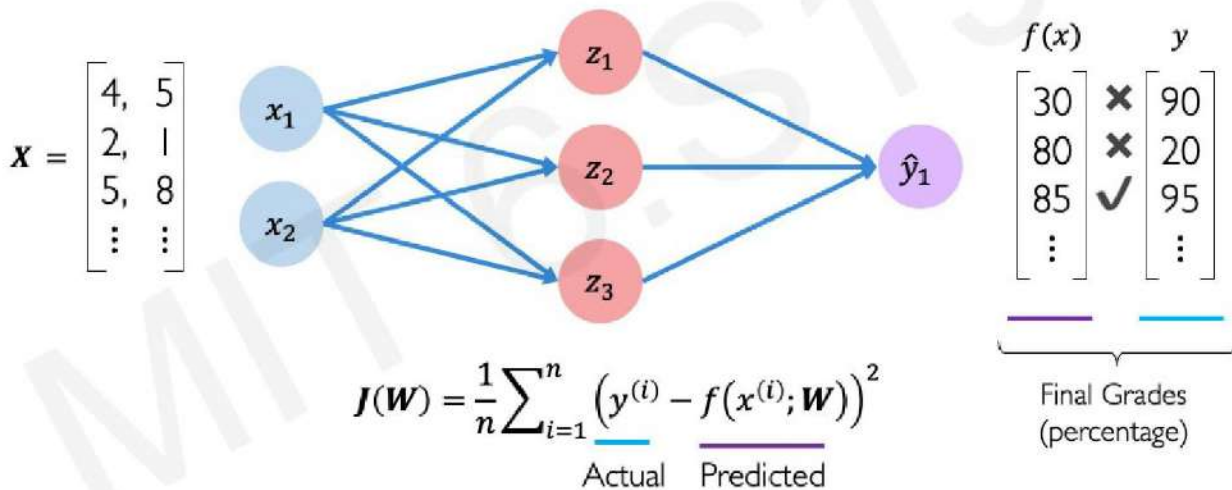


```
loss = tf.reduce_mean( tf.nn.softmax_cross_entropy_with_logits(y, predicted) )
```



Mean Squared Error Loss

Mean squared error loss can be used with regression models that output continuous real numbers



```
loss = tf.reduce_mean( tf.square(tf.subtract(y, predicted)) )  
loss = tf.keras.losses.MSE( y, predicted )
```



人工智能基本理论

Training Neural Networks



Loss Optimization

We want to find the network weights that **achieve the lowest loss**

$$W^* = \operatorname{argmin}_W \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; W), y^{(i)})$$

$$W^* = \operatorname{argmin}_W J(W)$$



Loss Optimization

We want to find the network weights that *achieve the lowest loss*

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} \frac{1}{n} \sum_{i=1}^n \mathcal{L}(f(x^{(i)}; \mathbf{W}), y^{(i)})$$

$$\mathbf{W}^* = \operatorname{argmin}_{\mathbf{W}} J(\mathbf{W})$$

Remember:

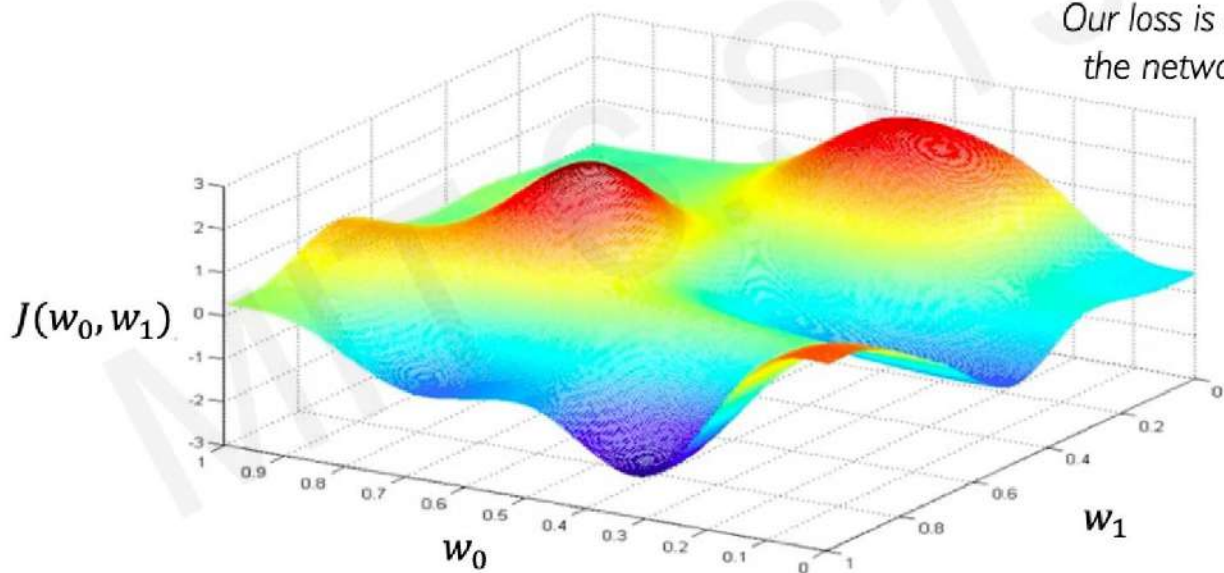
$$\mathbf{W} = \{\mathbf{W}^{(0)}, \mathbf{W}^{(1)}, \dots\}$$



Loss Optimization

$$W^* = \underset{W}{\operatorname{argmin}} J(W)$$

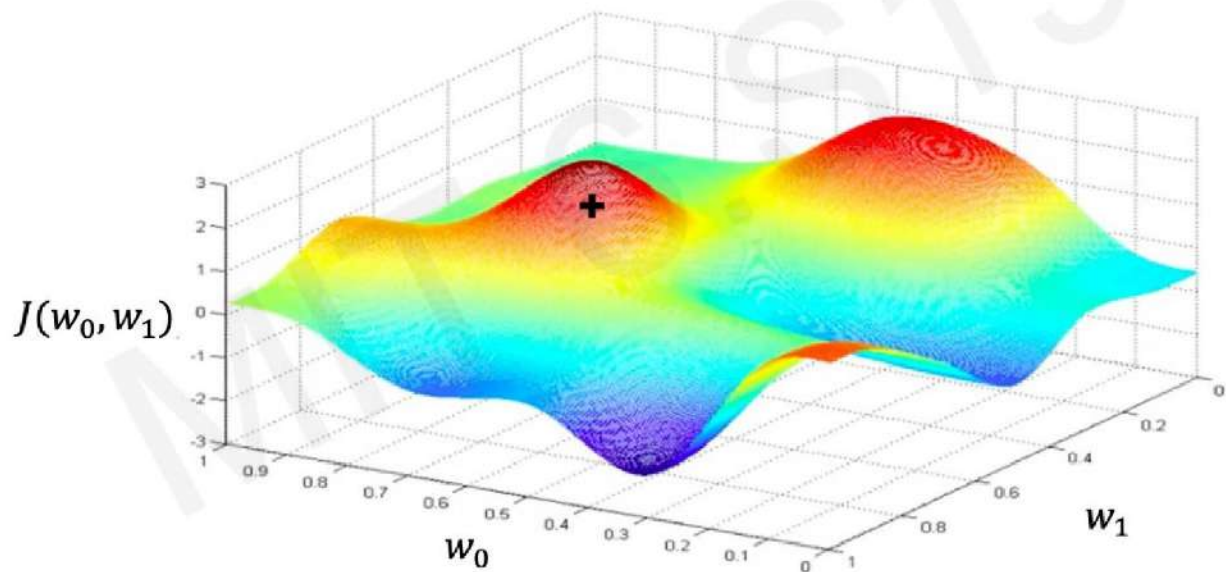
Remember:
*Our loss is a function of
the network weights!*





Loss Optimization

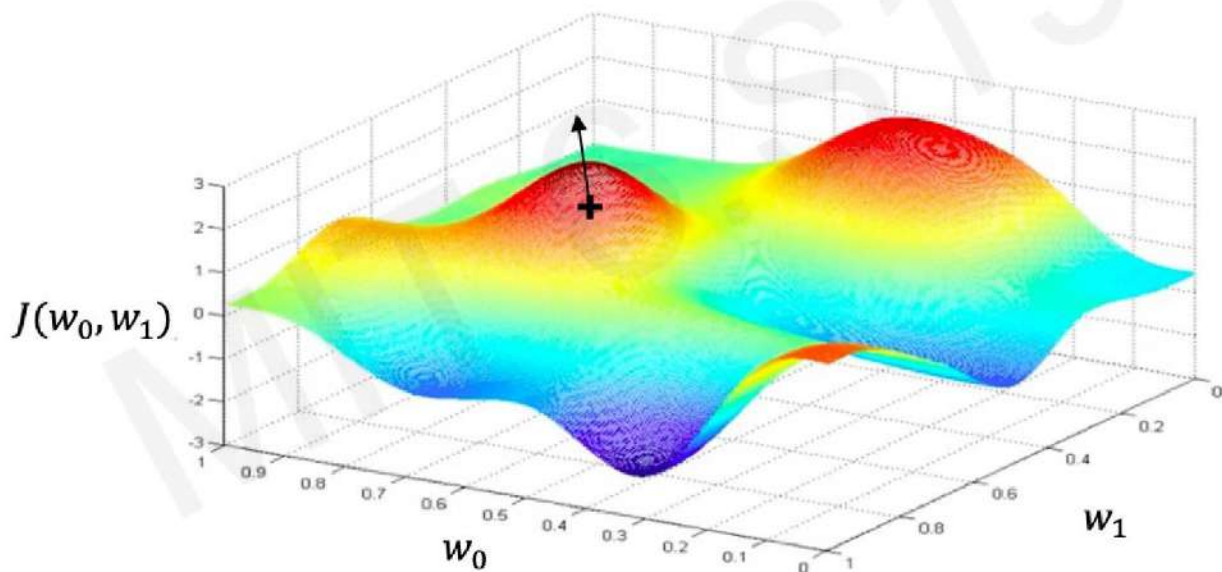
Randomly pick an initial (w_0, w_1)





Loss Optimization

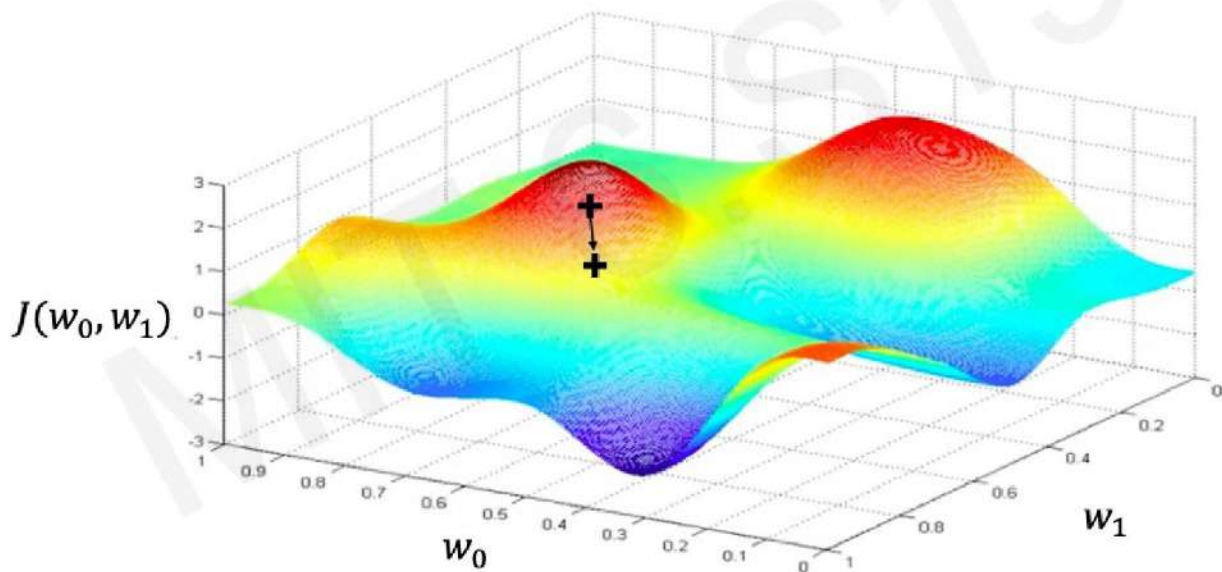
Compute gradient, $\frac{\partial J(w)}{\partial w}$





Loss Optimization

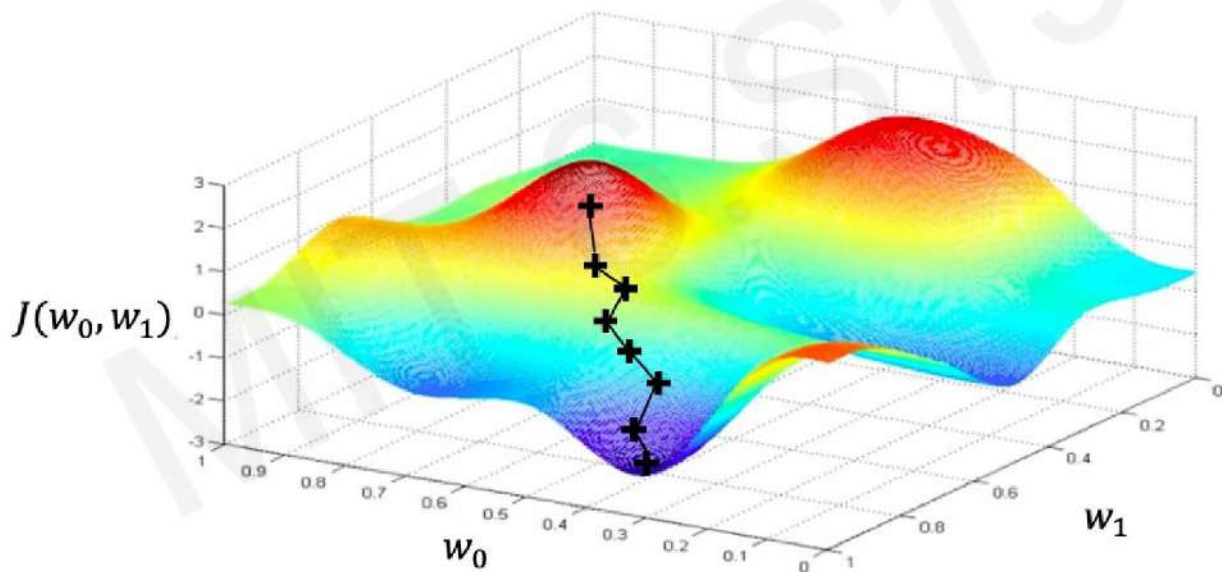
Take small step in opposite direction of gradient





Gradient Descent

Repeat until convergence





Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights



```
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True:    # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient
```



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
5. Return weights



```
import tensorflow as tf

weights = tf.Variable([tf.random.normal()])

while True: # loop forever
    with tf.GradientTape() as g:
        loss = compute_loss(weights)
        gradient = g.gradient(loss, weights)

    weights = weights - lr * gradient
```



Computing Gradients: Backpropagation



How does a small change in one weight (ex. w_2) affect the final loss $J(W)$?



Computing Gradients: Backpropagation

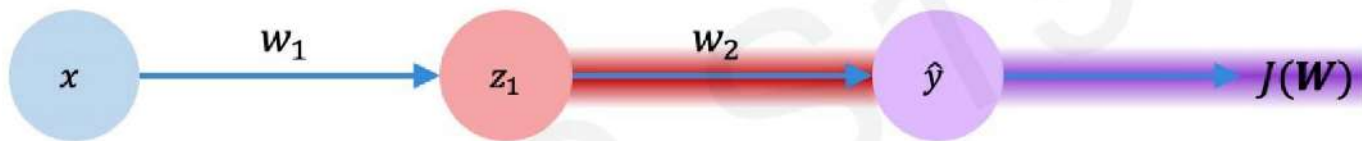


$$\frac{\partial J(W)}{\partial w_2} =$$

Let's use the chain rule!



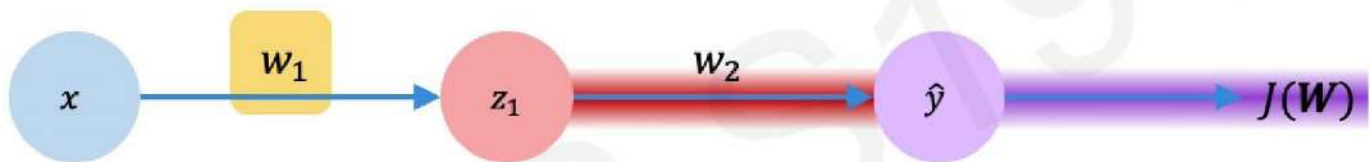
Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_2} = \underbrace{\frac{\partial J(W)}{\partial \hat{y}}}_{\text{purple}} * \underbrace{\frac{\partial \hat{y}}{\partial w_2}}_{\text{red}}$$



Computing Gradients: Backpropagation



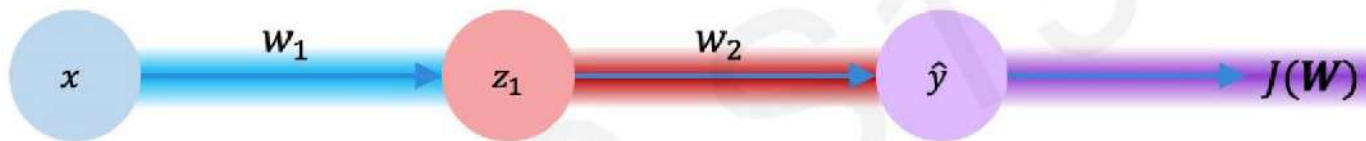
$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial w_1}$$

Apply chain rule!

Apply chain rule!



Computing Gradients: Backpropagation

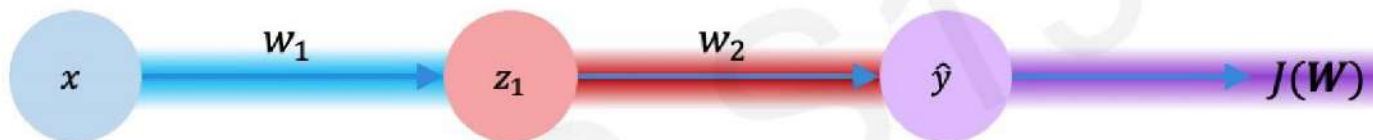


$$\frac{\partial J(\mathbf{W})}{\partial w_1} = \frac{\partial J(\mathbf{W})}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

(Purple bar under $\frac{\partial J(\mathbf{W})}{\partial \hat{y}}$, Red bar under $\frac{\partial \hat{y}}{\partial z_1}$, Blue bar under $\frac{\partial z_1}{\partial w_1}$)



Computing Gradients: Backpropagation



$$\frac{\partial J(W)}{\partial w_1} = \frac{\partial J(W)}{\partial \hat{y}} * \frac{\partial \hat{y}}{\partial z_1} * \frac{\partial z_1}{\partial w_1}$$

The equation shows the chain rule for computing the gradient of the loss function $J(W)$ with respect to the weight w_1 . The terms are color-coded: the gradient $\frac{\partial J(W)}{\partial \hat{y}}$ is purple, the gradient $\frac{\partial \hat{y}}{\partial z_1}$ is red, and the gradient $\frac{\partial z_1}{\partial w_1}$ is blue.

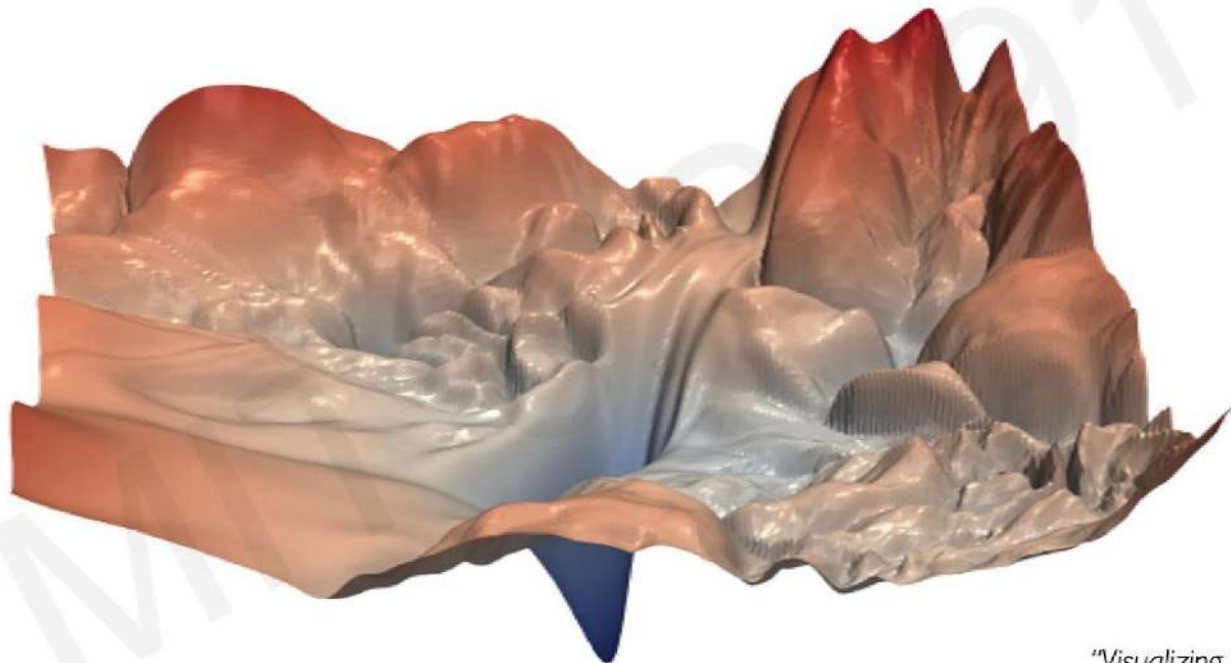
Repeat this for **every weight in the network** using gradients from later layers



Neural Networks in Practice: Optimization



Training Neural Networks is Difficult



"Visualizing the loss landscape of neural nets". Dec 2017.



Loss Functions Can Be Difficult to Optimize

Remember:

Optimization through gradient descent

$$\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$$



Loss Functions Can Be Difficult to Optimize

Remember:

Optimization through gradient descent

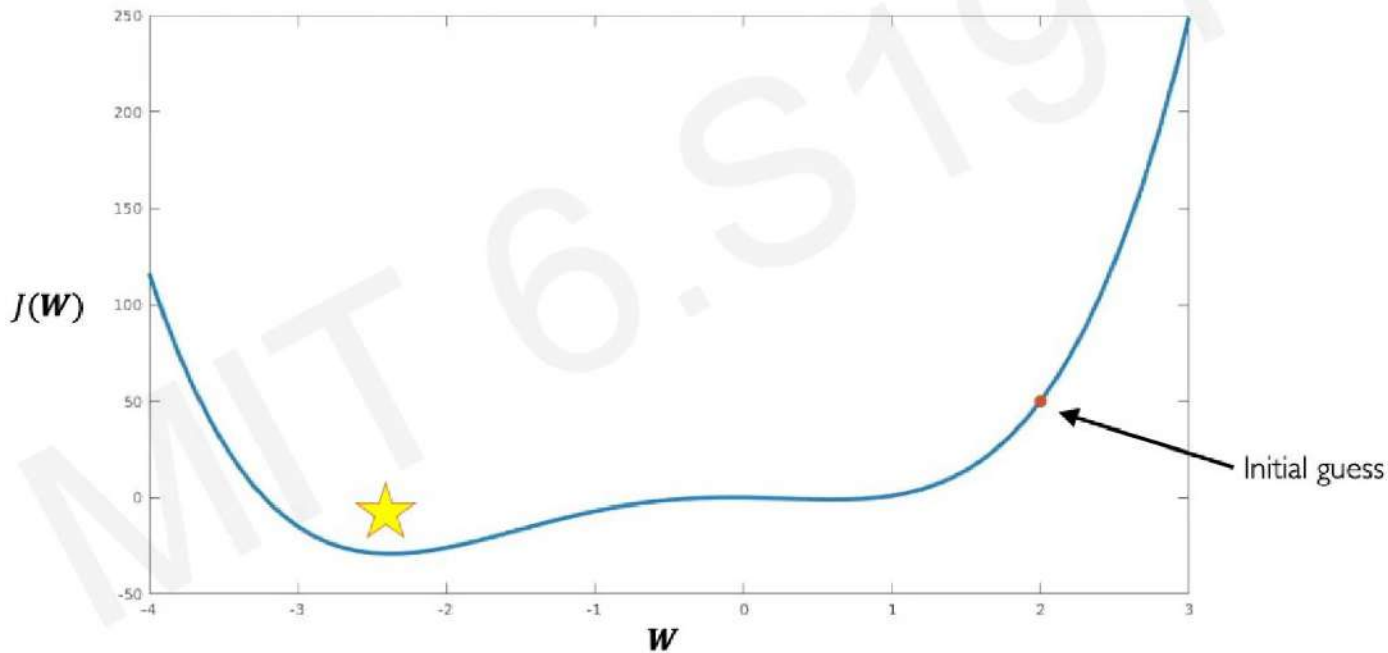
$$W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$$

How can we set the learning rate?



Setting the Learning Rate

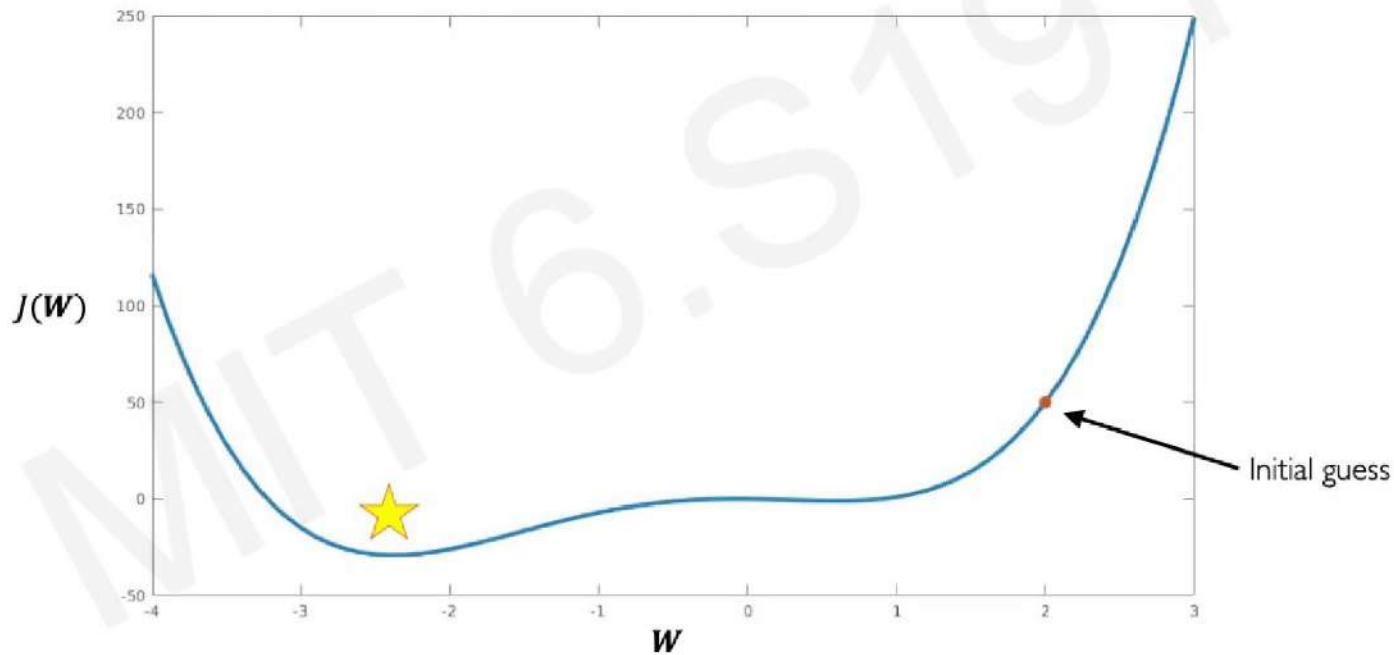
Small learning rate converges slowly and gets stuck in false local minima





Setting the Learning Rate

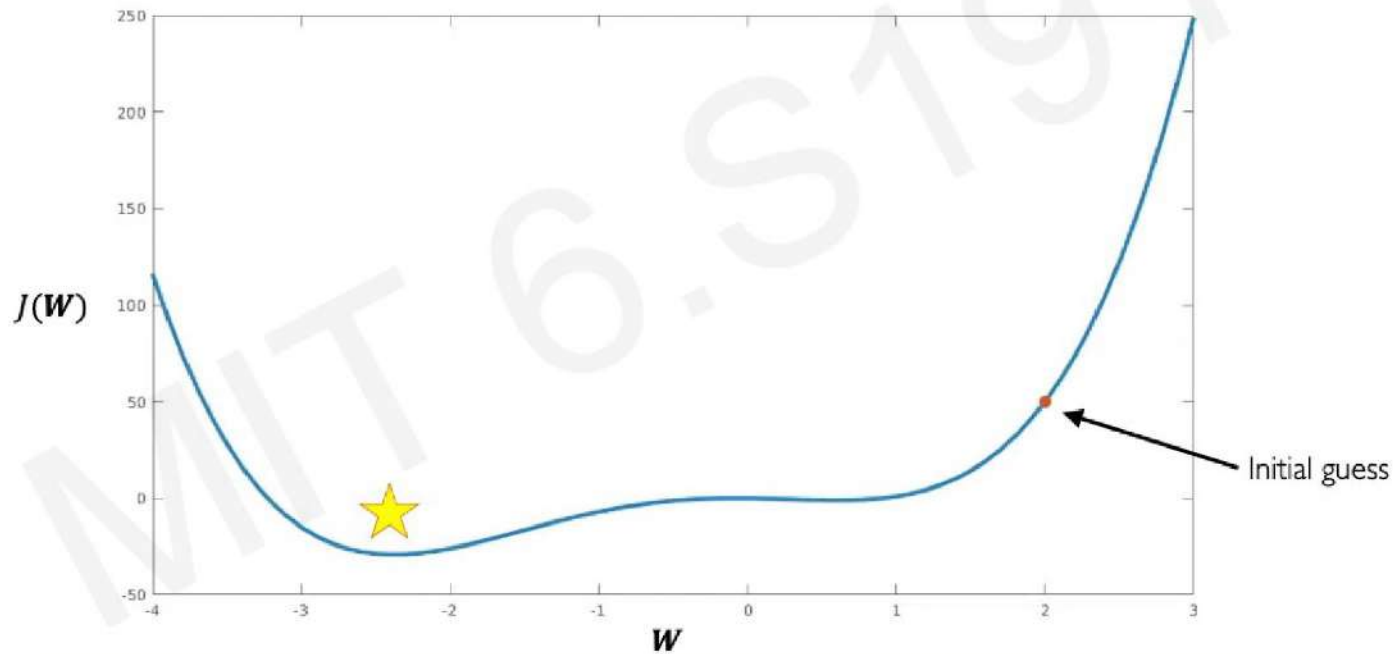
Large learning rates overshoot, become unstable and diverge





Setting the Learning Rate

Stable learning rates converge smoothly and avoid local minima





How to deal with this?

Idea 1:

Try lots of different learning rates and see what works “just right”

MIT 6.S191



How to deal with this?

Idea 1:

Try lots of different learning rates and see what works “just right”

Idea 2:

Do something smarter!

Design an adaptive learning rate that “adapts” to the landscape



Adaptive Learning Rates

- Learning rates are no longer fixed
- Can be made larger or smaller depending on:
 - how large gradient is
 - how fast learning is happening
 - size of particular weights
 - etc...



Gradient Descent Algorithms

Algorithm	TF Implementation	Reference
• SGD	 <code>tf.keras.optimizers.SGD</code>	Kiefer & Wolfowitz. "Stochastic Estimation of the Maximum of a Regression Function." 1952.
• Adam	 <code>tf.keras.optimizers.Adam</code>	Kingma et al. "Adam: A Method for Stochastic Optimization." 2014.
• Adadelata	 <code>tf.keras.optimizers.Adadelata</code>	Zeiler et al. "ADADELTA: An Adaptive Learning Rate Method." 2012.
• Adagrad	 <code>tf.keras.optimizers.Adagrad</code>	Duchi et al. "Adaptive Subgradient Methods for Online Learning and Stochastic Optimization." 2011.
• RMSProp	 <code>tf.keras.optimizers.RMSProp</code>	

Additional details: <http://ruder.io/optimizing-gradient-descent/>



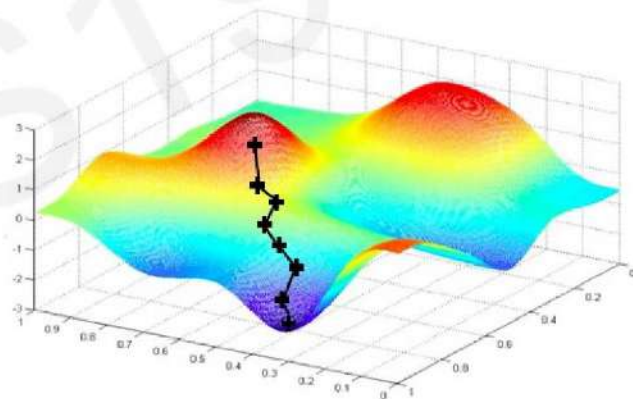
Neural Networks in Practice: Mini-batches



Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(W)}{\partial W}$
4. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights

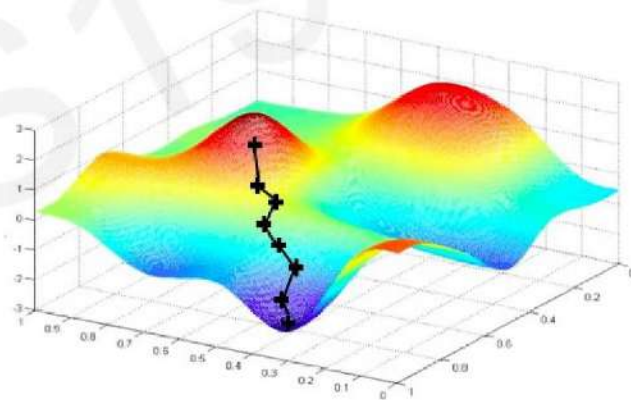




Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Compute gradient, $\frac{\partial J(W)}{\partial W}$
4. Update weights, $W \leftarrow W - \eta \frac{\partial J(W)}{\partial W}$
5. Return weights



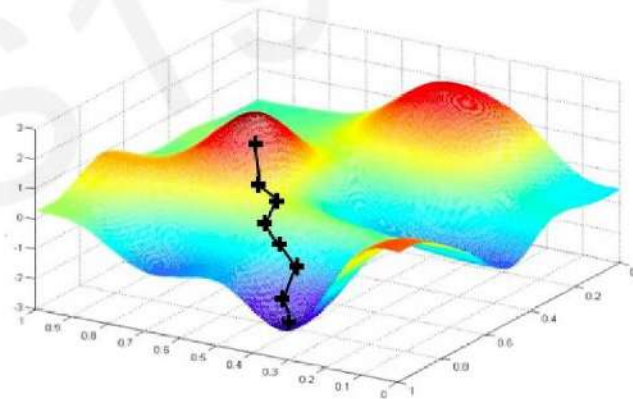
Can be very
computationally
intensive to compute!



Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
 3. Pick single data point i
 4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
 5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

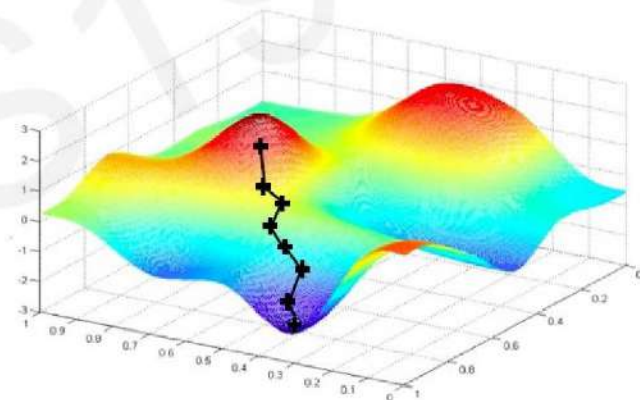




Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick single data point i
4. Compute gradient, $\frac{\partial J_i(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



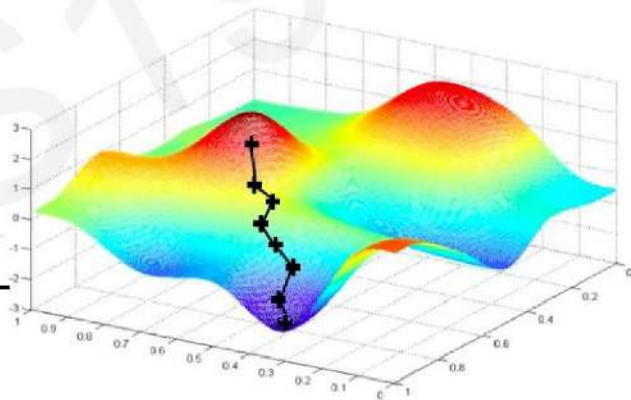
Easy to compute but
very noisy (stochastic)!



Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights

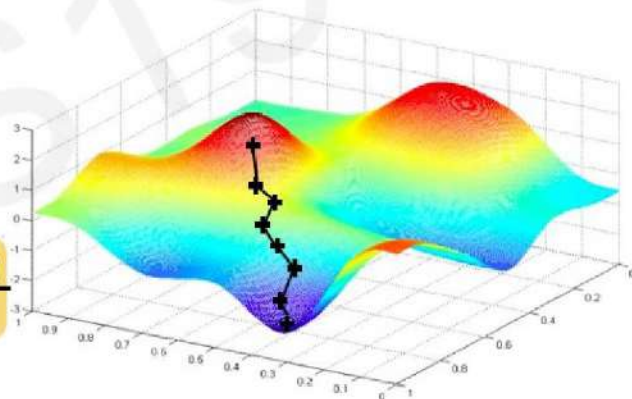




Stochastic Gradient Descent

Algorithm

1. Initialize weights randomly $\sim \mathcal{N}(0, \sigma^2)$
2. Loop until convergence:
3. Pick batch of B data points
4. Compute gradient, $\frac{\partial J(\mathbf{W})}{\partial \mathbf{W}} = \frac{1}{B} \sum_{k=1}^B \frac{\partial J_k(\mathbf{W})}{\partial \mathbf{W}}$
5. Update weights, $\mathbf{W} \leftarrow \mathbf{W} - \eta \frac{\partial J(\mathbf{W})}{\partial \mathbf{W}}$
6. Return weights



Fast to compute and a much better estimate of the true gradient!



Mini-batches while training

More accurate estimation of gradient

Smoother convergence

Allows for larger learning rates



Mini-batches while training

More accurate estimation of gradient

Smoother convergence

Allows for larger learning rates

Mini-batches lead to fast training!

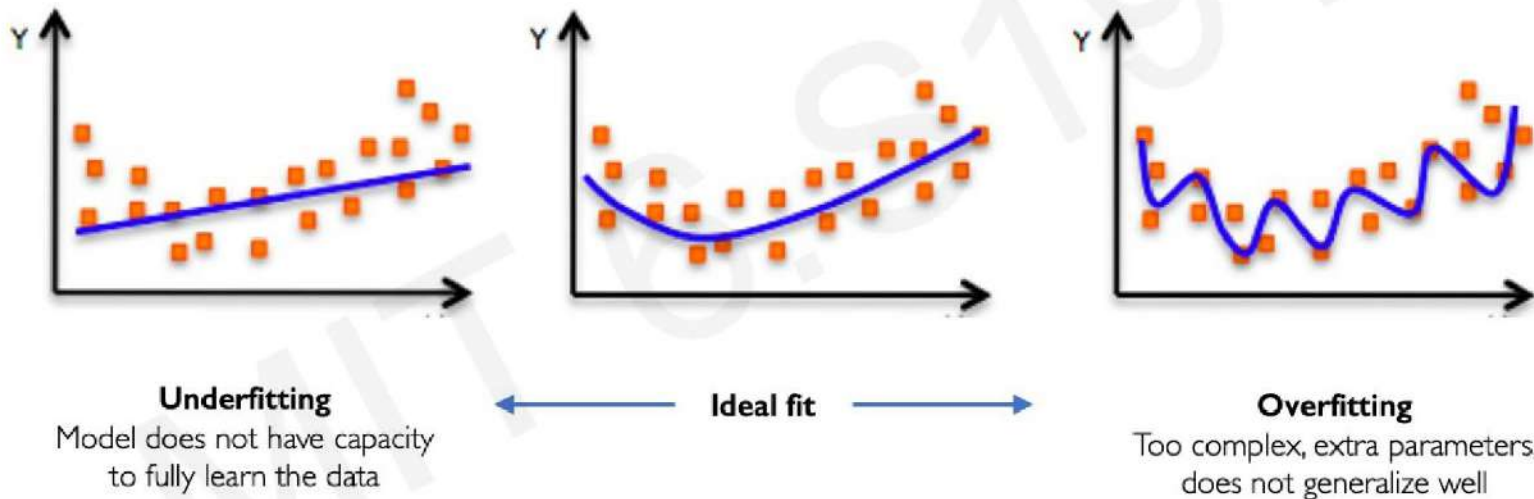
Can parallelize computation + achieve significant speed increases on GPU's



Neural Networks in Practice: Overfitting



The Problem of Overfitting





Regularization

What is it?

Technique that constrains our optimization problem to discourage complex models

MIT 6.S191



Regularization

What is it?

Technique that constrains our optimization problem to discourage complex models

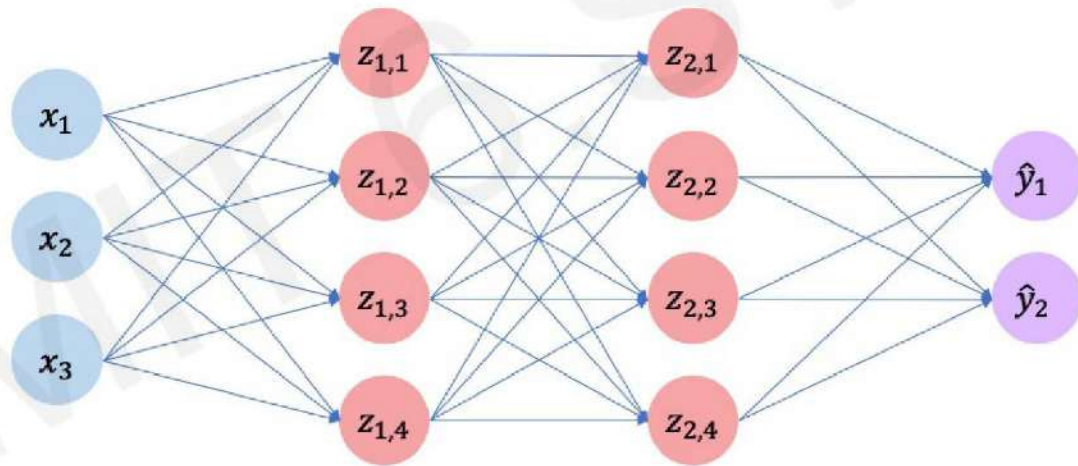
Why do we need it?

Improve generalization of our model on unseen data



Regularization I: Dropout

- During training, randomly set some activations to 0

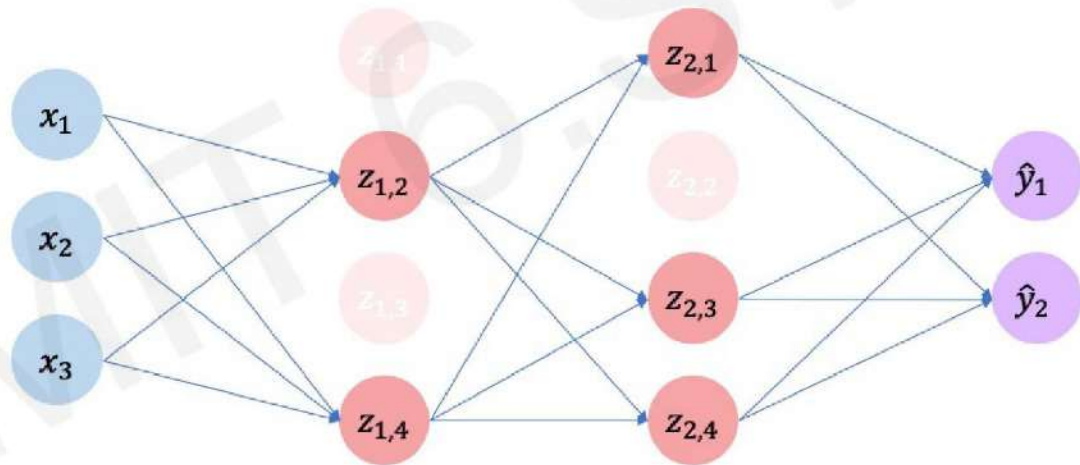




Regularization I: Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

```
tf.keras.layers.Dropout(p=0.5)
```

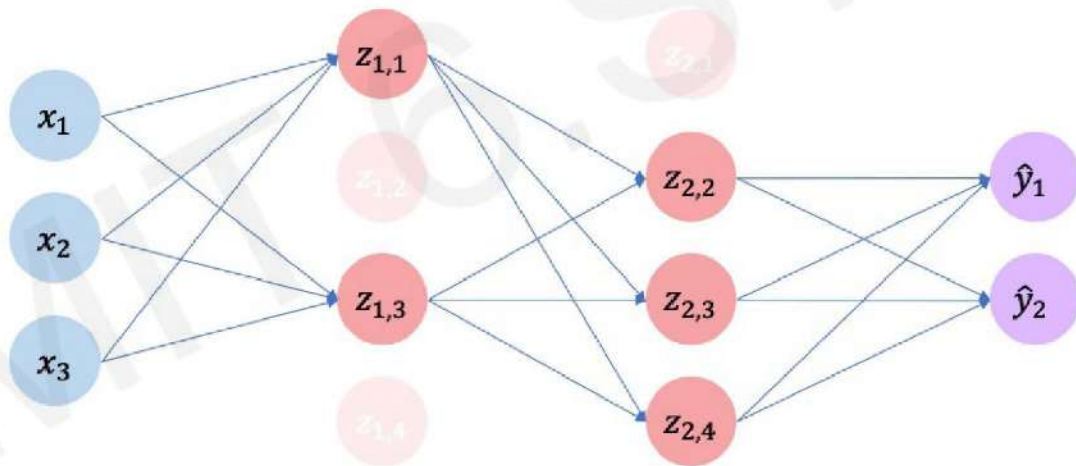




Regularization I: Dropout

- During training, randomly set some activations to 0
 - Typically 'drop' 50% of activations in layer
 - Forces network to not rely on any 1 node

```
tf.keras.layers.Dropout(p=0.5)
```





Regularization 2: Early Stopping

- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

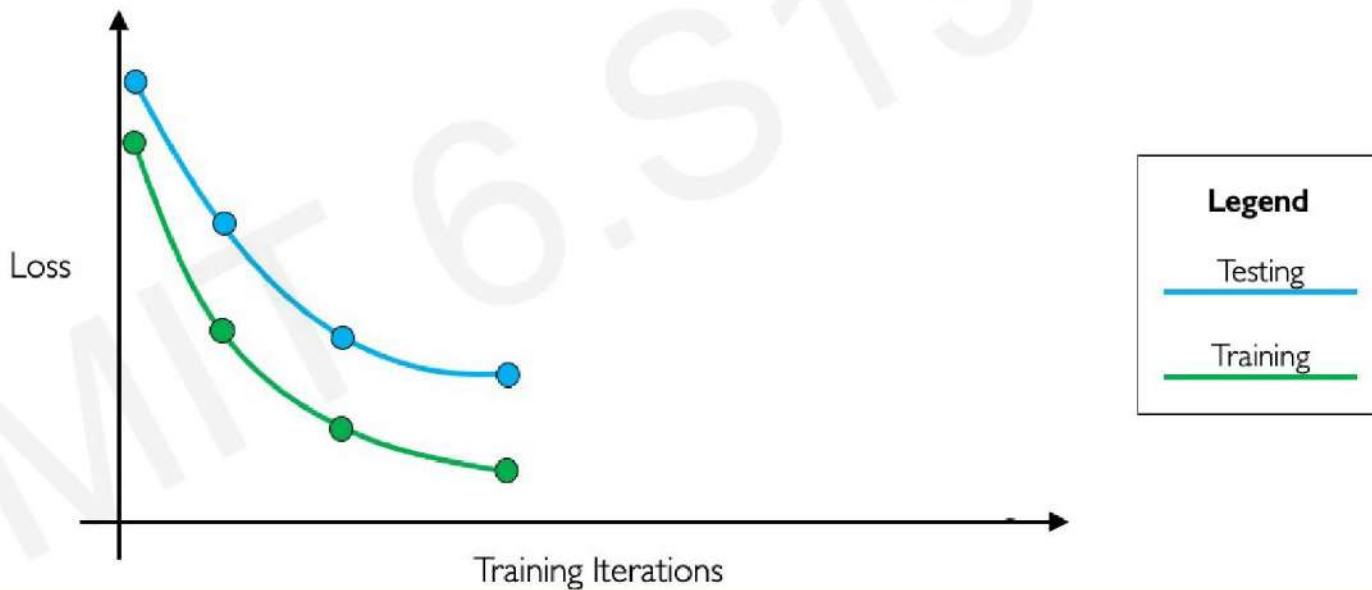
- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

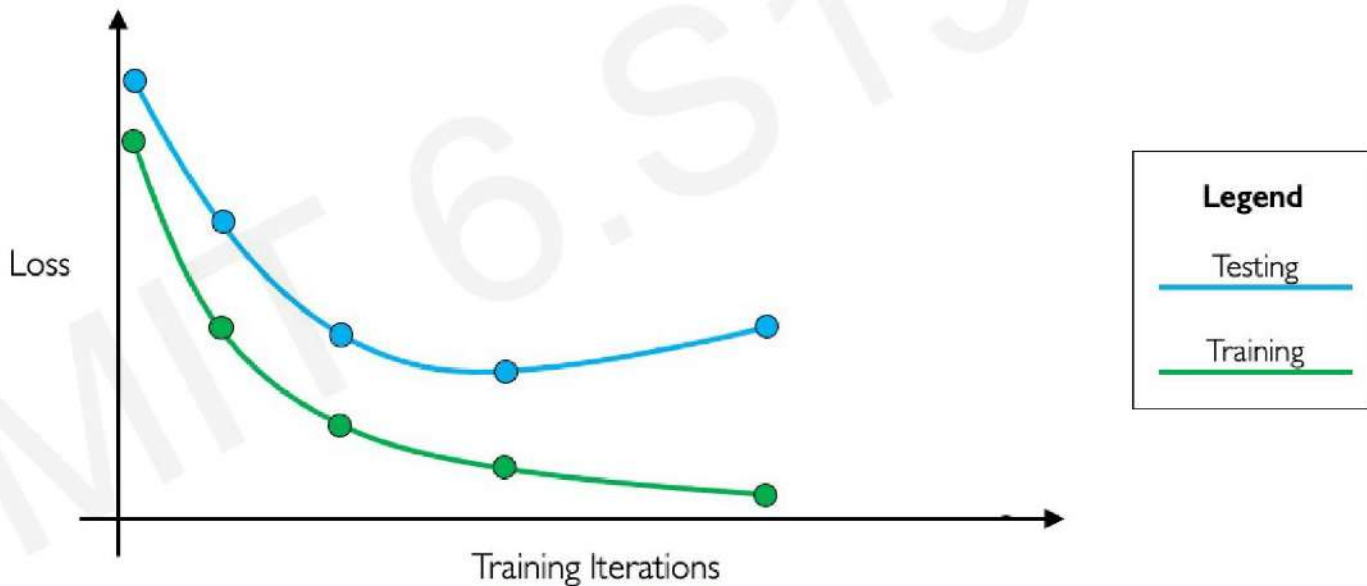
- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

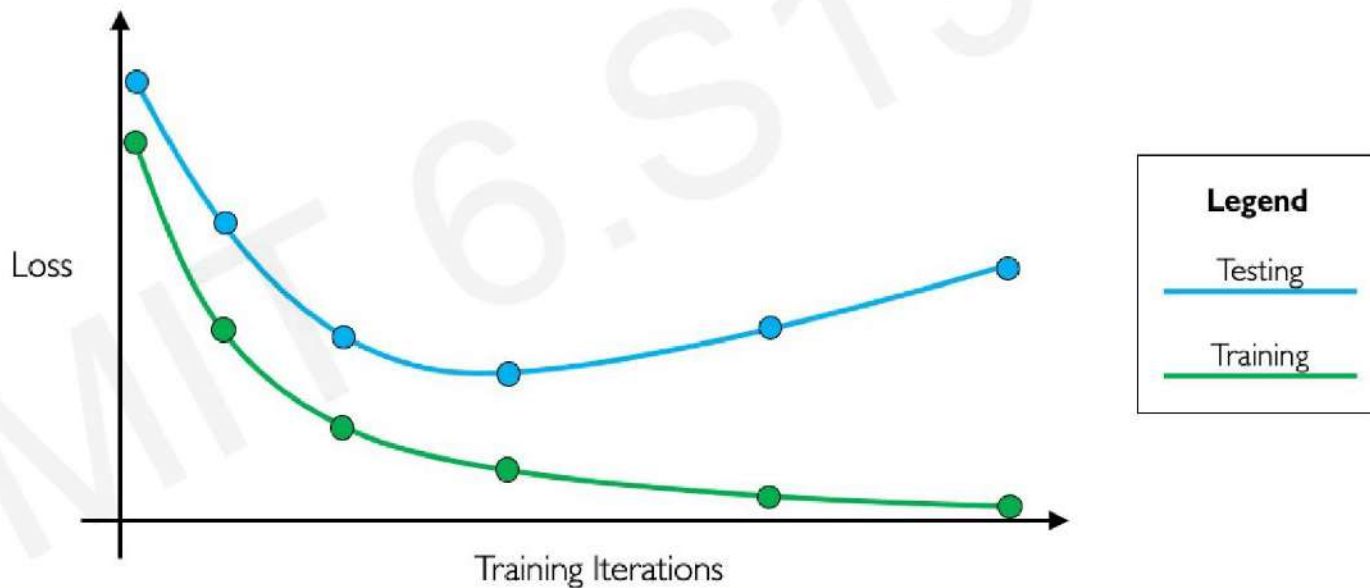
- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

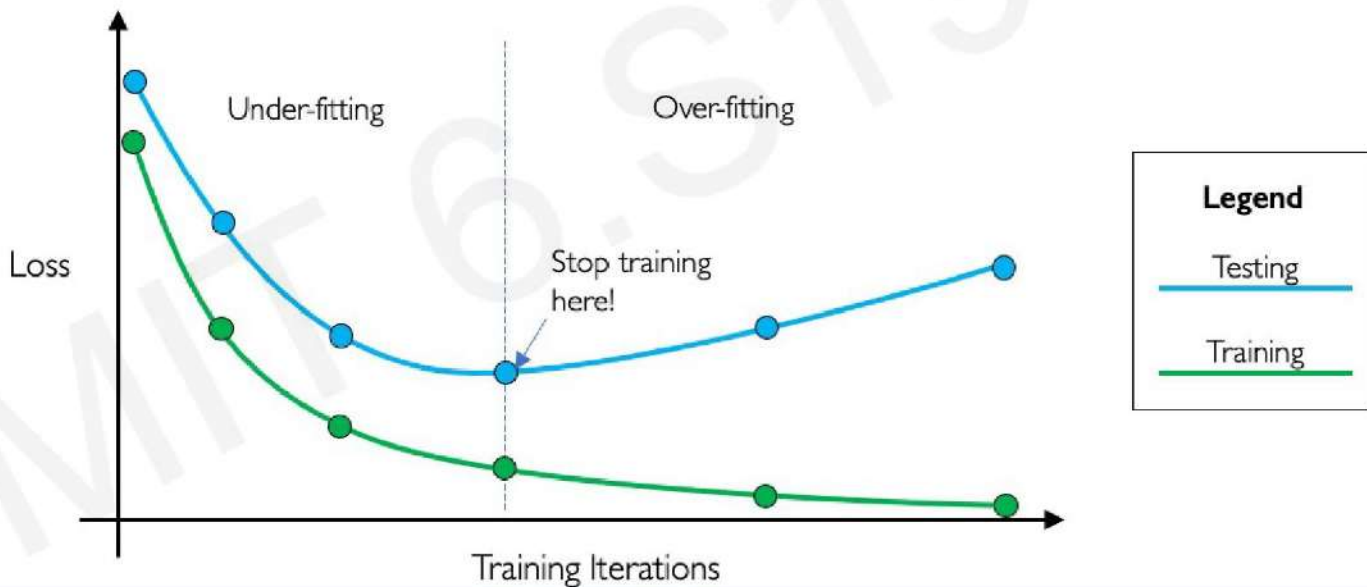
- Stop training before we have a chance to overfit





Regularization 2: Early Stopping

- Stop training before we have a chance to overfit



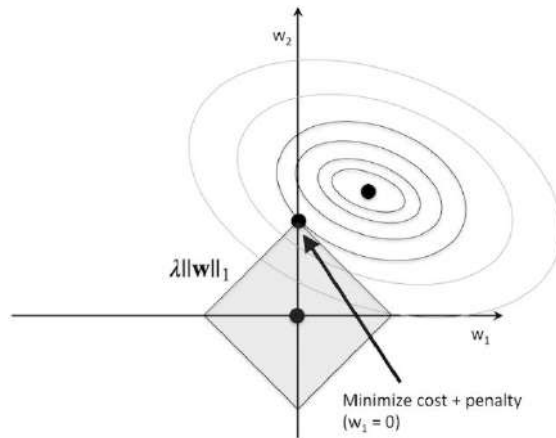


L1 Regularization

$$\text{loss}(\text{truth}, \text{predictions}) + \lambda |W|$$

Or, if we number each individual weight w in the model with a number $i = 0, 1, 2, \dots, n$:

$$\text{loss}(\text{truth}, \text{predictions}) + \lambda \sum_{i=0}^n |w_i|$$

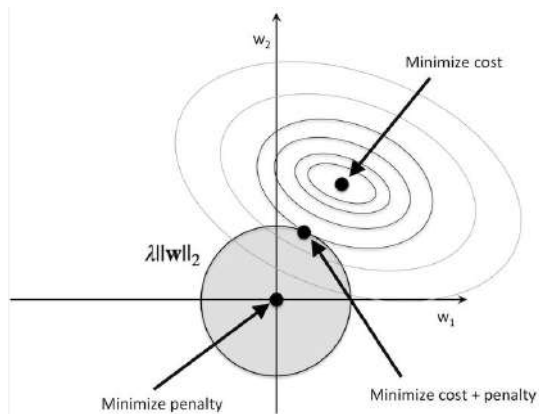


L2 Regularization

$$\text{loss}(\text{truth}, \text{predictions}) + \lambda W^2$$

Or, if we number each individual weight w in the model with a number $i = 0, 1, 2, \dots, n$:

$$\text{loss}(\text{truth}, \text{predictions}) + \lambda \sum_{i=0}^n w_i^2$$

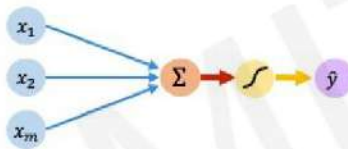




Core Foundation Review

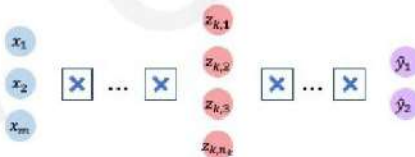
The Perceptron

- Structural building blocks
- Nonlinear activation functions



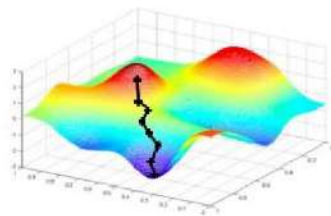
Neural Networks

- Stacking Perceptrons to form neural networks
- Optimization through backpropagation



Training in Practice

- Adaptive learning
- Batching
- Regularization



提纲

一、前馈网络

二、循环网络

三、卷积网络



上海大学
SHANGHAI UNIVERSITY



Sequences in the Wild





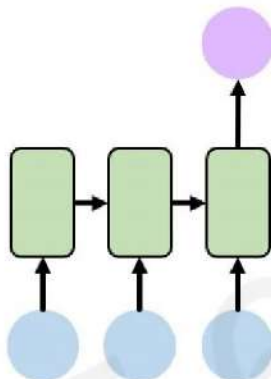
Sequence Modeling Applications



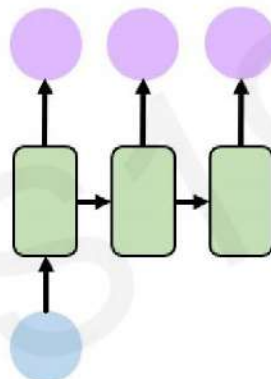
One to One
Binary Classification



"Will I pass this class?"
Student → Pass?



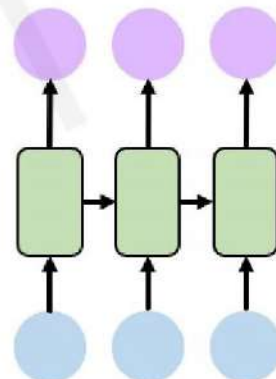
Many to One
Sentiment Classification



One to Many
Image Captioning



"A baseball player throws a ball."



Many to Many
Machine Translation





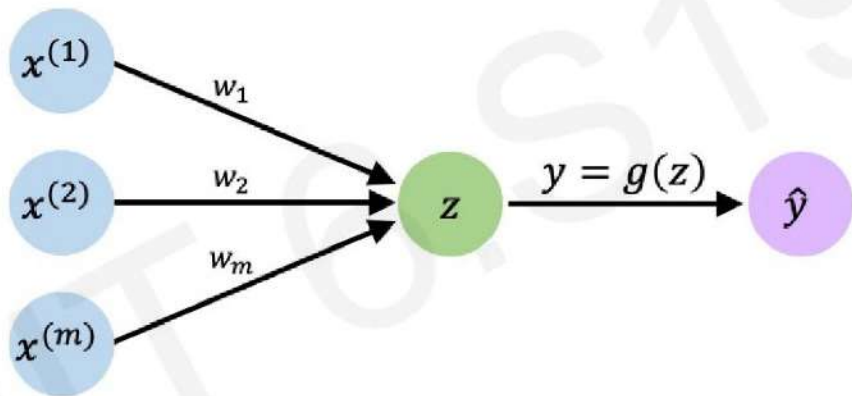
人工智能基本理论

Neurons with Recurrence

MIT 6.S191

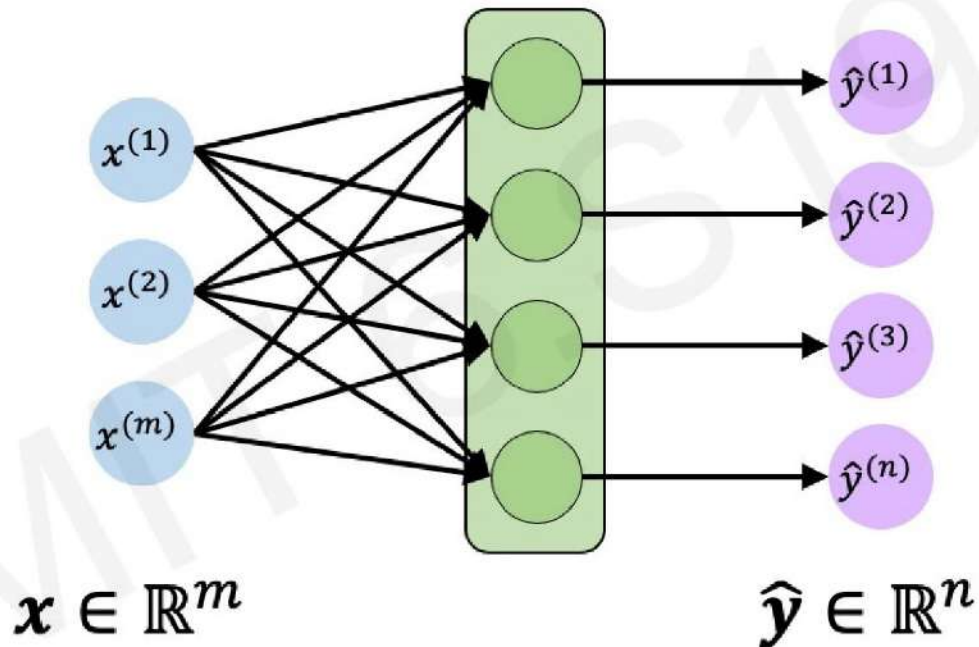


The Perceptron Revisited



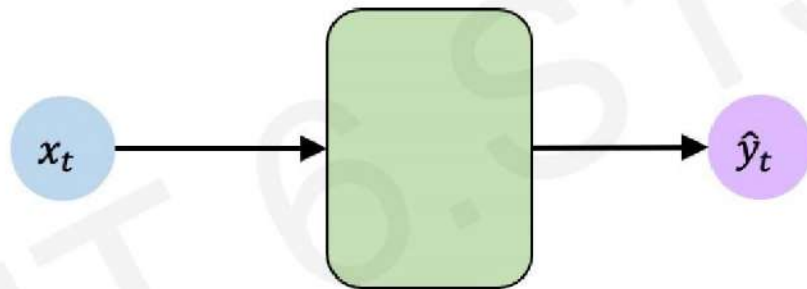


Feed-Forward Networks Revisited





Feed-Forward Networks Revisited

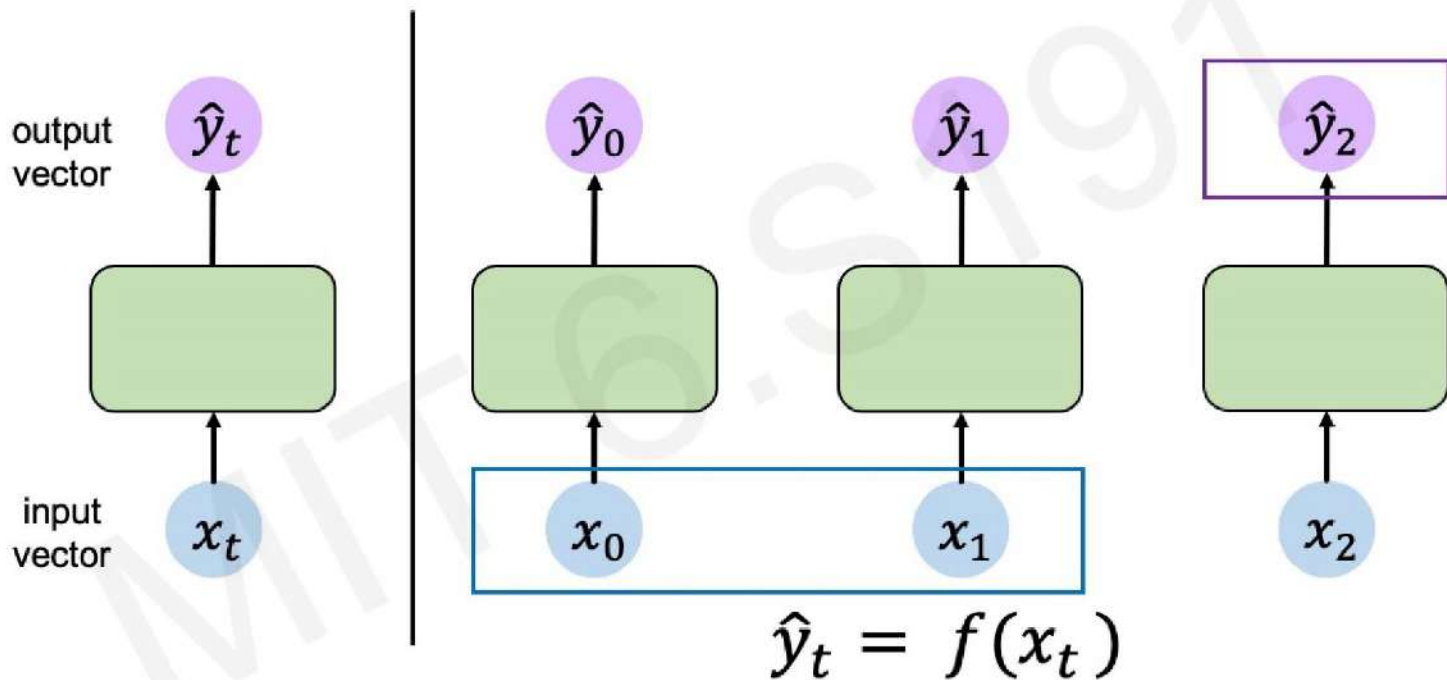


$$x_t \in \mathbb{R}^m$$

$$\hat{y}_t \in \mathbb{R}^n$$

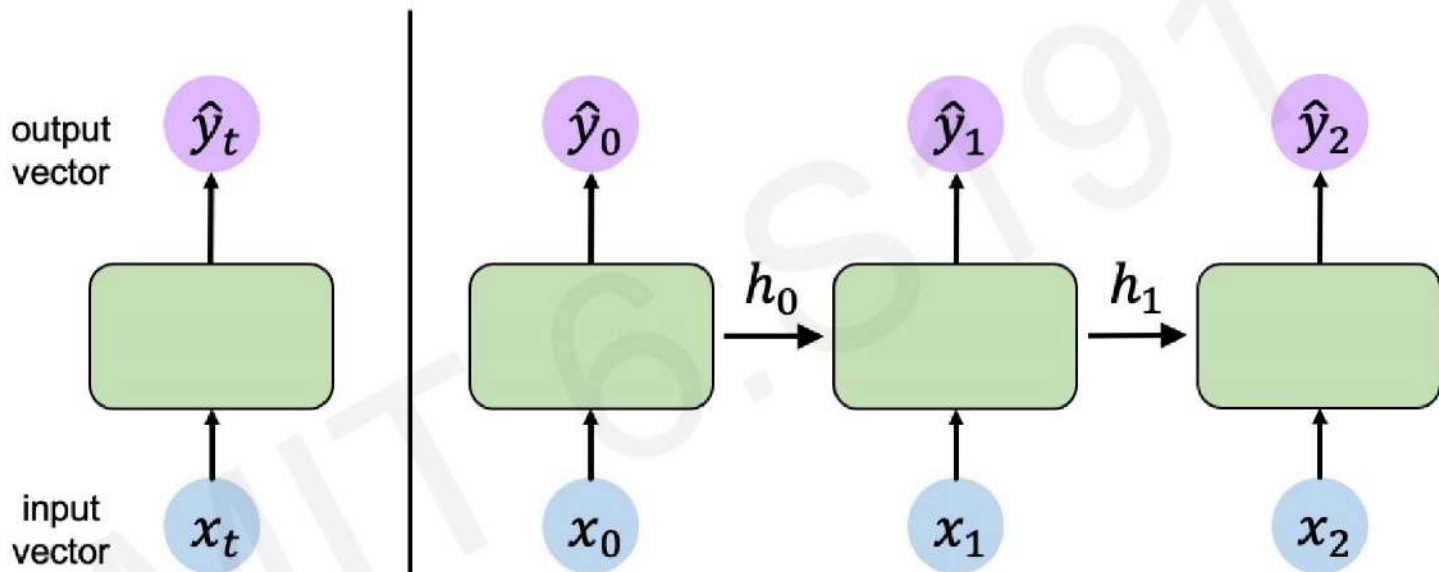


Handling Individual Time Steps





Neurons with Recurrence

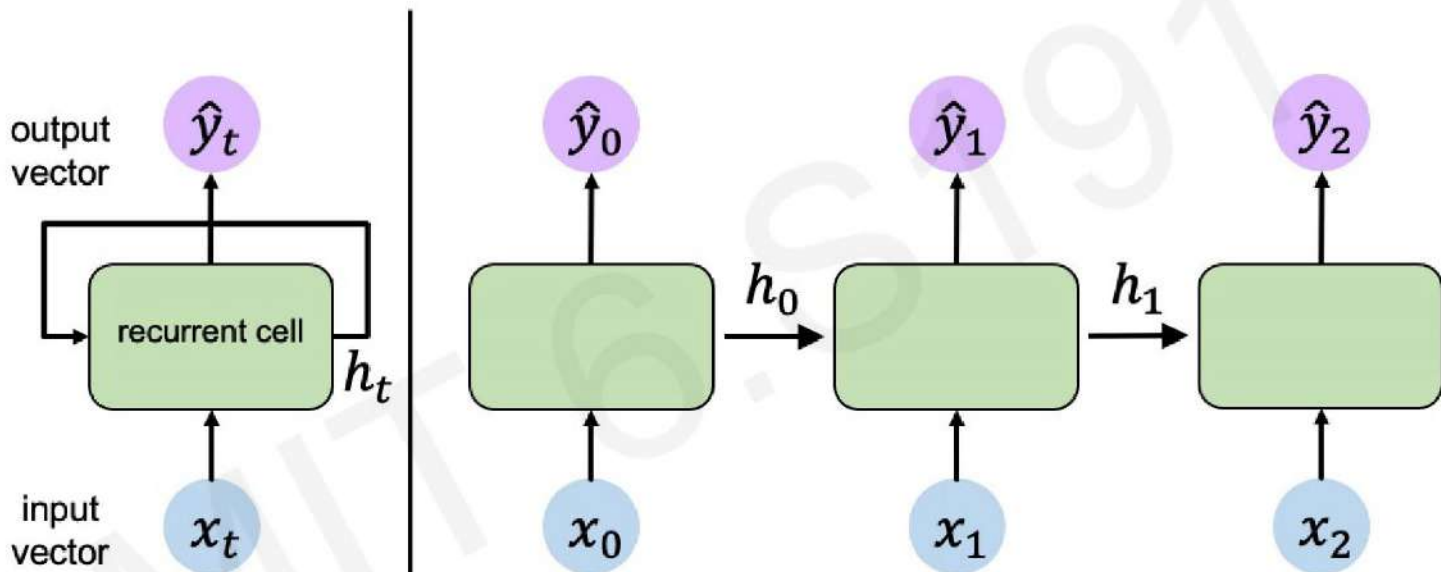


$$\hat{y}_t = f(x_t, h_{t-1})$$

output input past memory



Neurons with Recurrence



$$\hat{y}_t = f(x_t, h_{t-1})$$

output input past memory

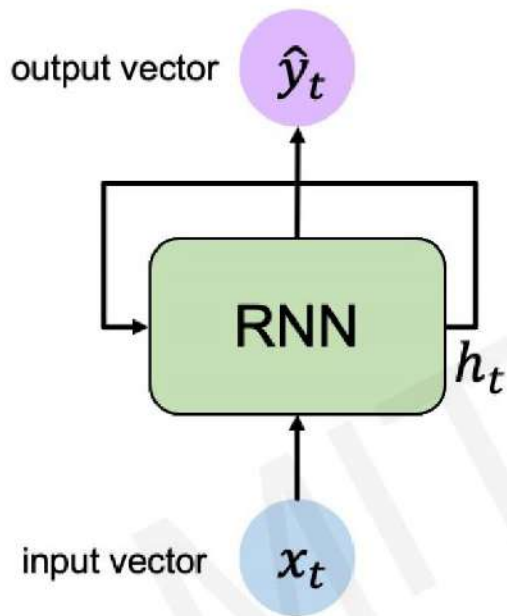


Recurrent Neural Networks (RNNs)

MIT 6.S191



Recurrent Neural Networks (RNNs)



Apply a **recurrence relation** at every time step to process a sequence:

$$h_t = f_W(x_t, h_{t-1})$$

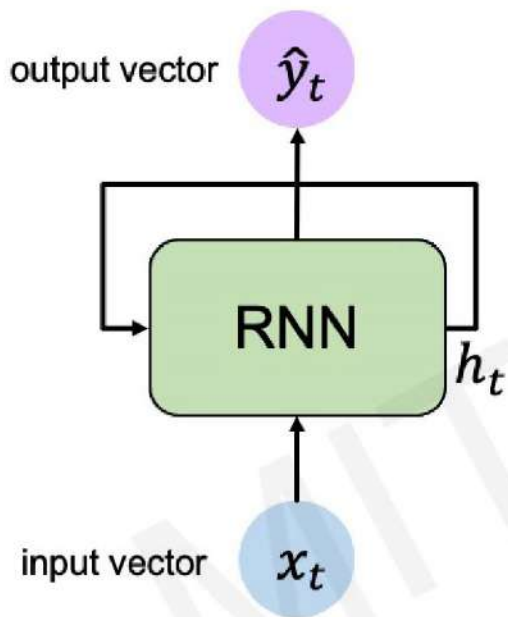
cell state function with weights W input old state

Note: the same function and set of parameters are used at every time step

RNNs have a **state**, h_t , that is updated **at each time step** as a sequence is processed



RNN State Update and Output



Update Hidden State

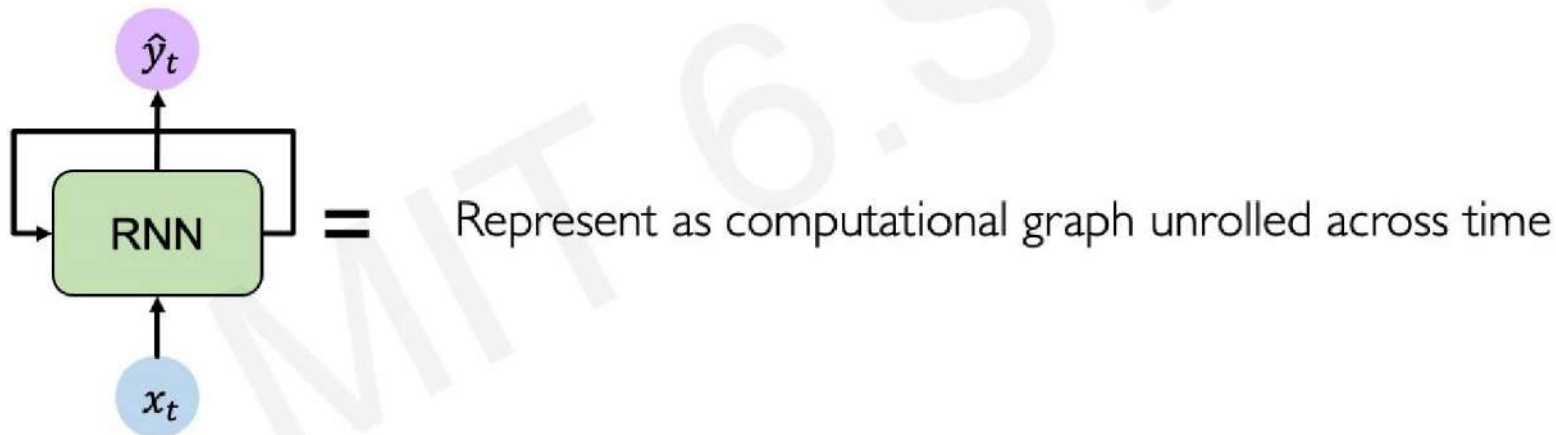
$$h_t = \tanh(W_{hh}^T h_{t-1} + W_{xh}^T x_t)$$

Input Vector

x_t

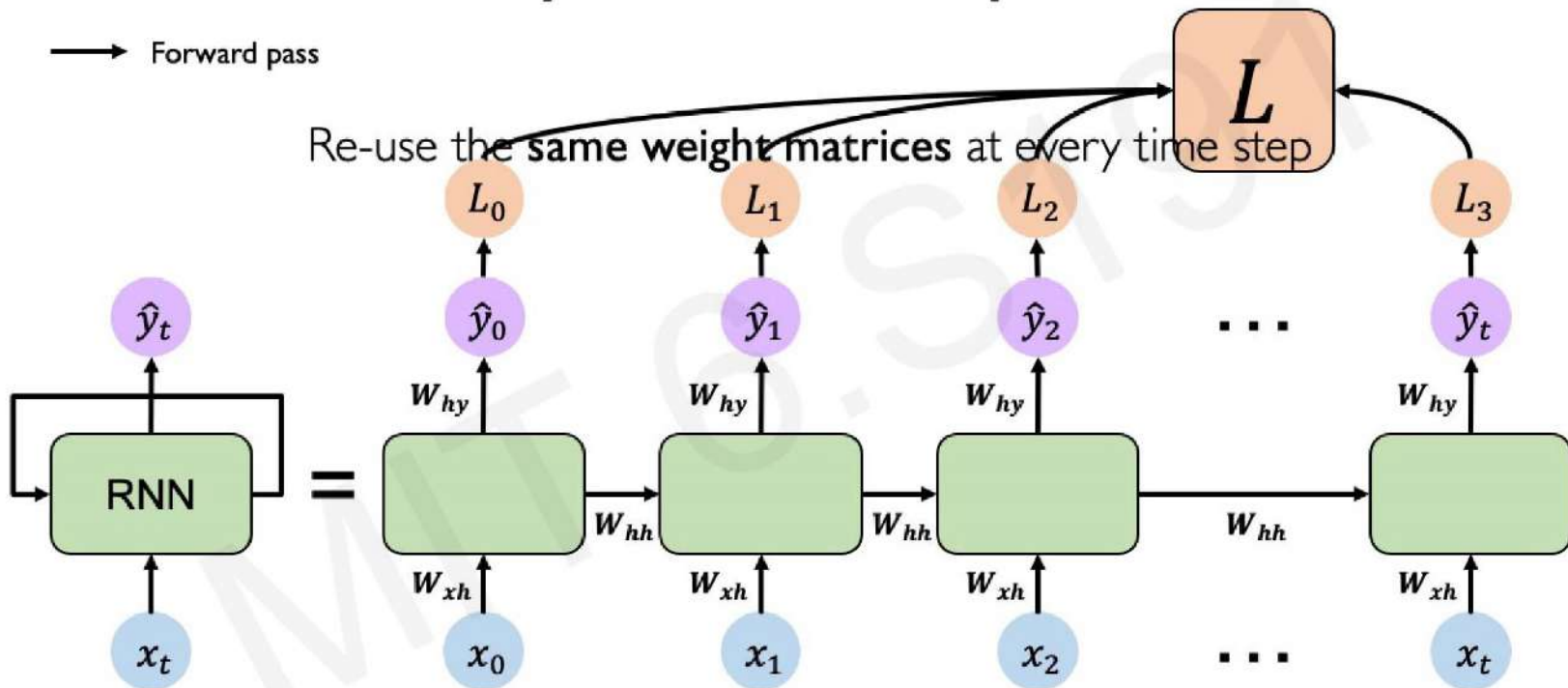


RNNs: Computational Graph Across Time





RNNs: Computational Graph Across Time

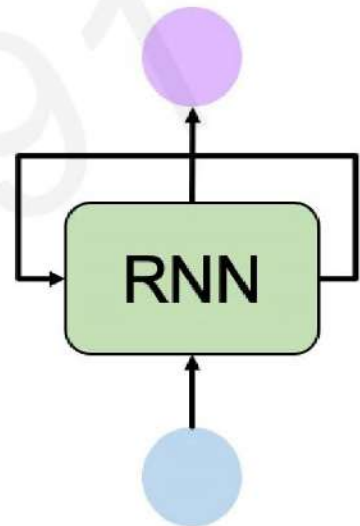




Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences
2. Track **long-term** dependencies
3. Maintain information about **order**
4. **Share parameters** across the sequence



Recurrent Neural Networks (RNNs) meet these sequence modeling design criteria



A Sequence Modeling Problem:
Predict the Next Word



A Sequence Modeling Problem: Predict the Next Word

“This morning I took my cat for a walk.”

MIT 6.S191



A Sequence Modeling Problem: Predict the Next Word

“This morning I took my cat for a walk.”

given these words

MIT 6.S191



A Sequence Modeling Problem: Predict the Next Word

“This morning I took my cat for a walk.”

given these words

predict the
next word

MIT 6.S191



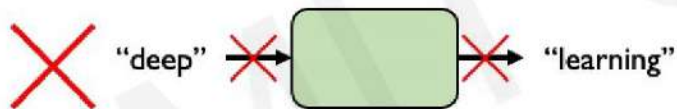
A Sequence Modeling Problem: Predict the Next Word

“This morning I took my cat for a walk.”

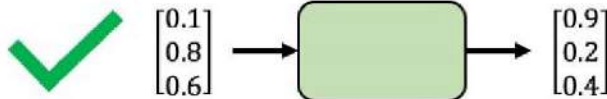
given these words

predict the
next word

Representing Language to a Neural Network



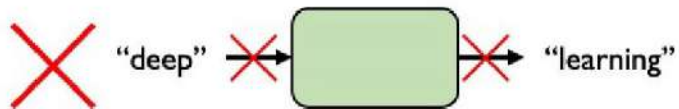
Neural networks cannot interpret words



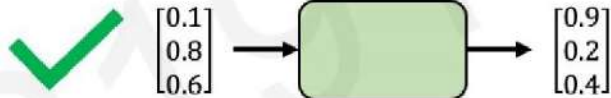
Neural networks require numerical inputs



Encoding Language for a Neural Network

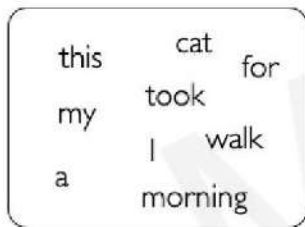


Neural networks cannot interpret words

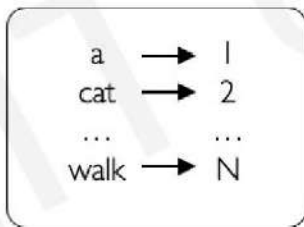


Neural networks require numerical inputs

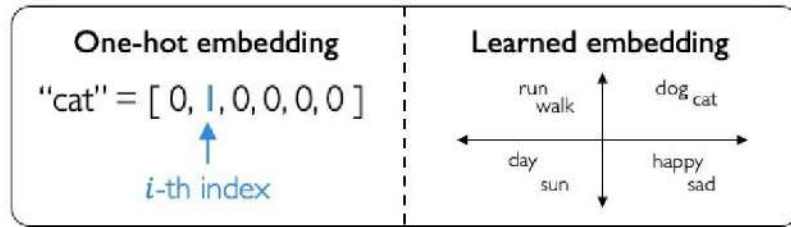
Embedding: transform indexes into a vector of fixed size.



1. Vocabulary:
Corpus of words



2. Indexing:
Word to index



3. Embedding:
Index to fixed-sized vector



Handle Variable Sequence Lengths

The food was great

vs.

We visited a restaurant for lunch

vs.

We were hungry but cleaned the house before eating



Model Long-Term Dependencies

“**France** is where I grew up, but I now live in Boston. I speak fluent ____.”



We need information from **the distant past** to accurately predict the correct word.



Capture Differences in Sequence Order



The food was good, not bad at all.

vs.

The food was bad, not good at all.

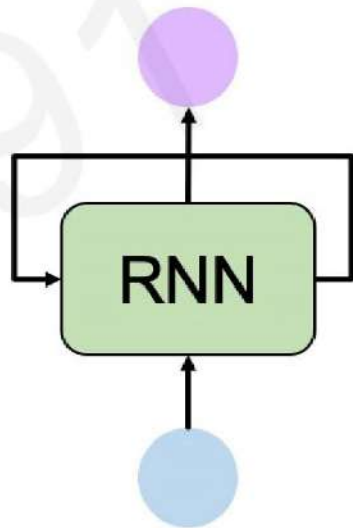




Sequence Modeling: Design Criteria

To model sequences, we need to:

1. Handle **variable-length** sequences
2. Track **long-term** dependencies
3. Maintain information about **order**
4. **Share parameters** across the sequence



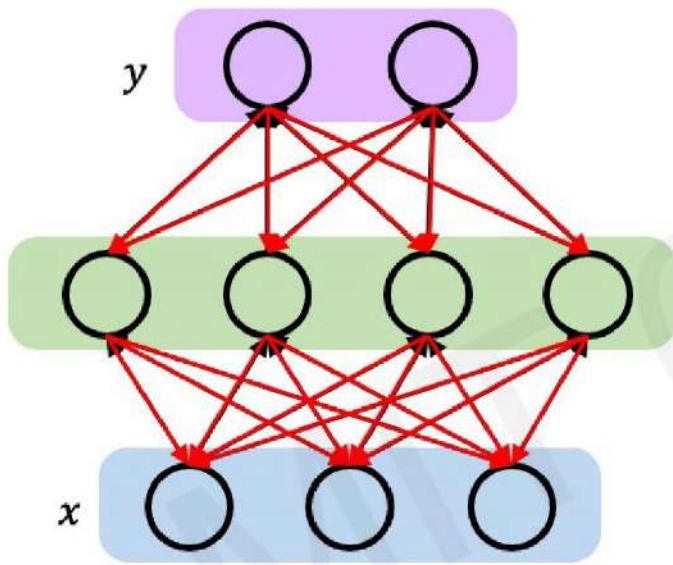
Recurrent Neural Networks (RNNs) meet these sequence modeling design criteria



Backpropagation Through Time (BPTT)



Recall: Backpropagation in Feed Forward Models

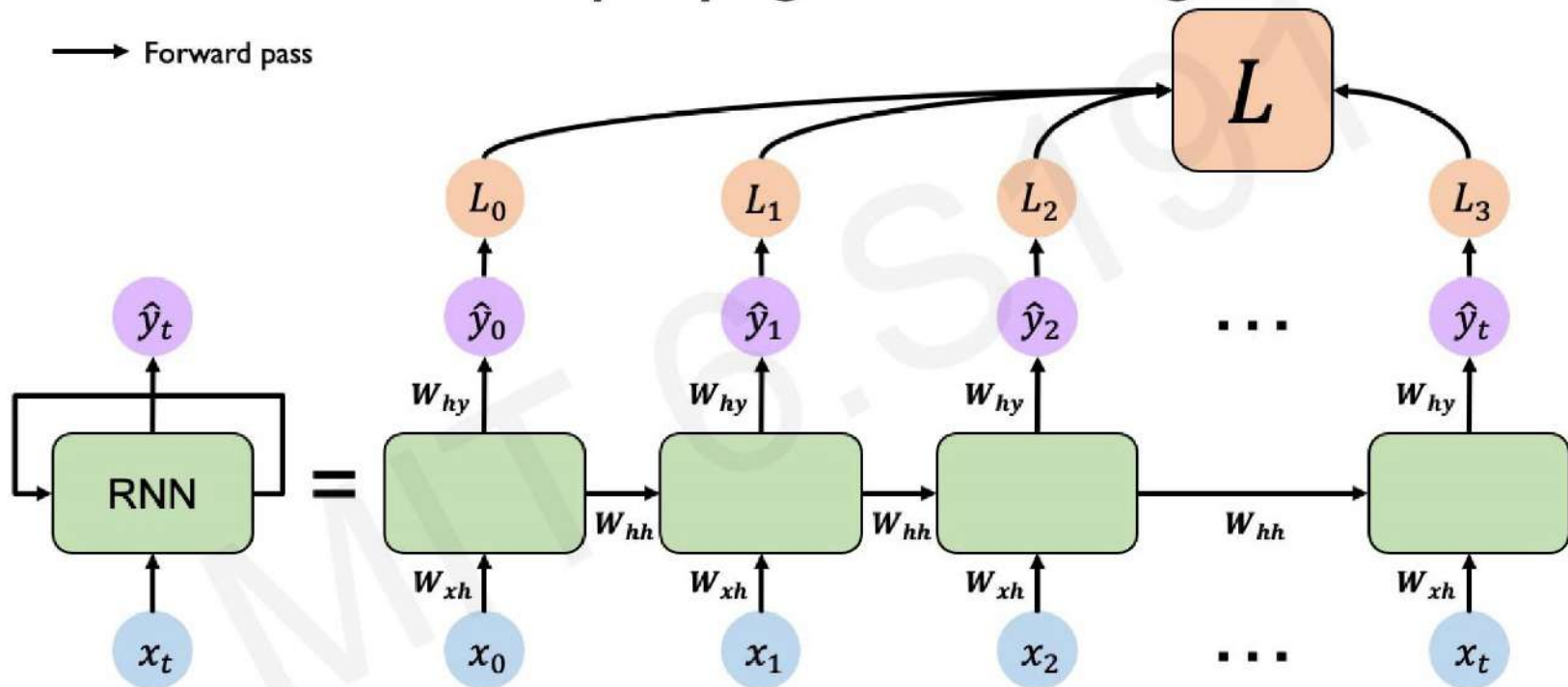


Backpropagation algorithm:

1. Take the derivative (gradient) of the loss with respect to each parameter
2. Shift parameters in order to minimize loss

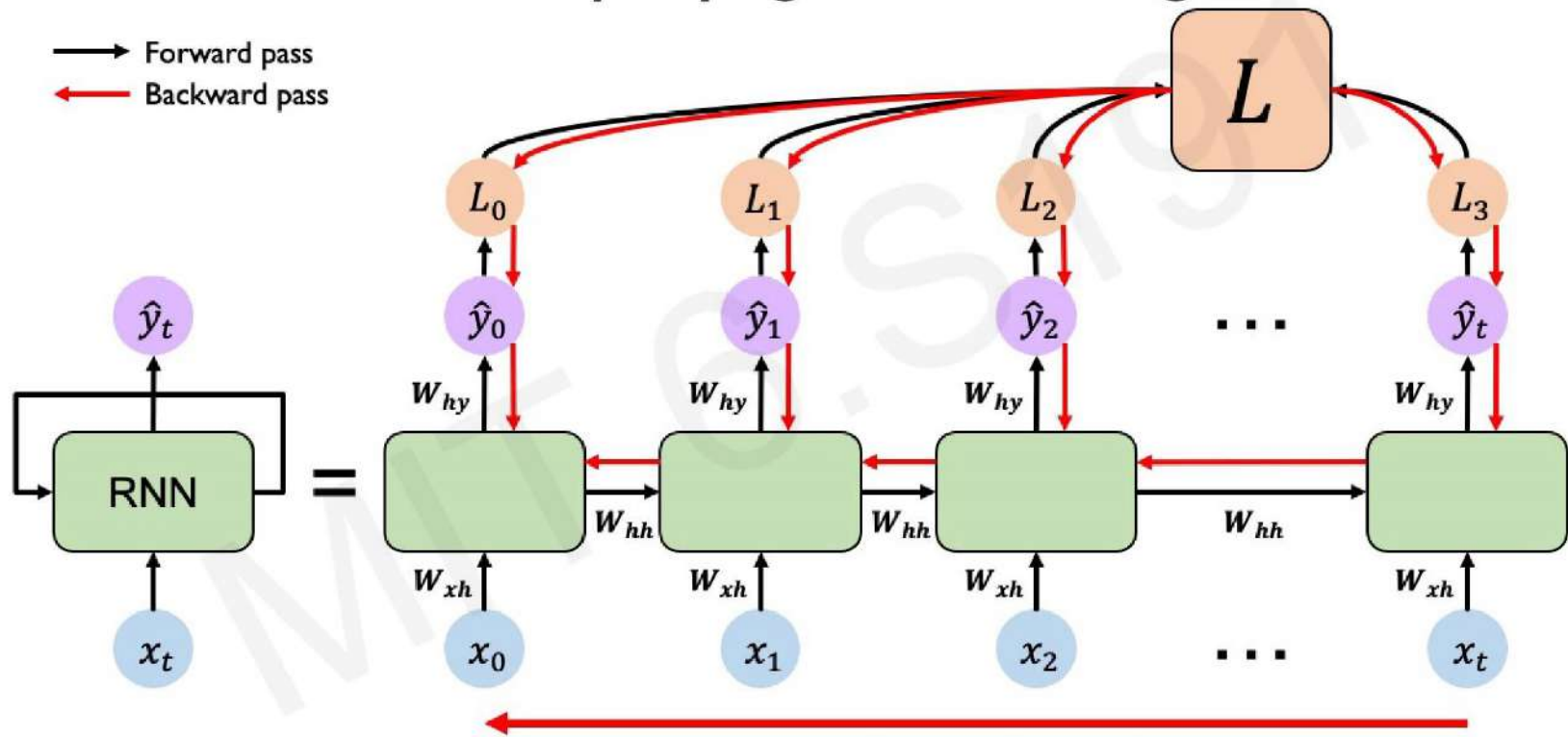


RNNs: Backpropagation Through Time



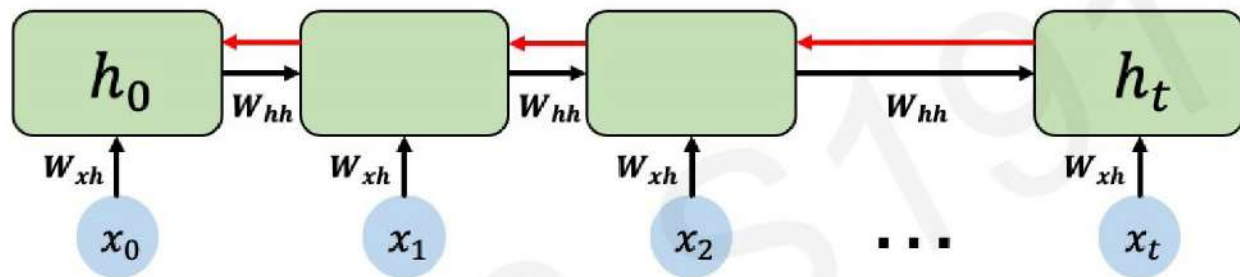


RNNs: Backpropagation Through Time



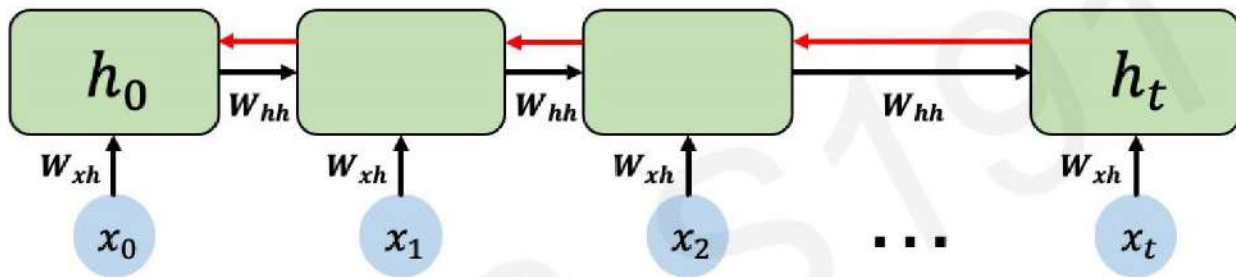


Standard RNN Gradient Flow





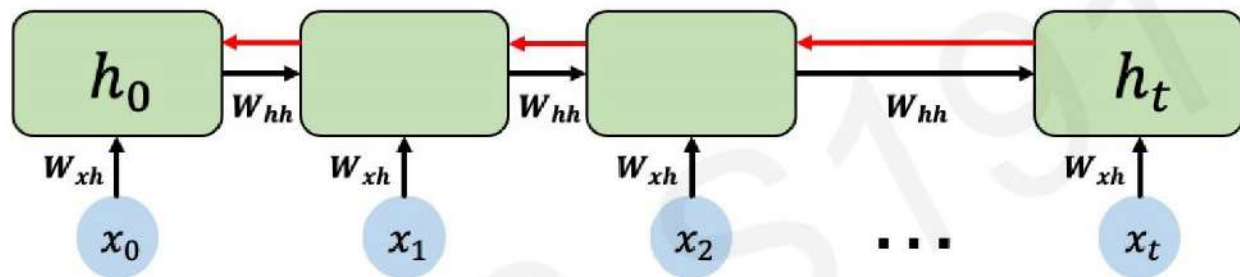
Standard RNN Gradient Flow



Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!



Standard RNN Gradient Flow: Exploding Gradients



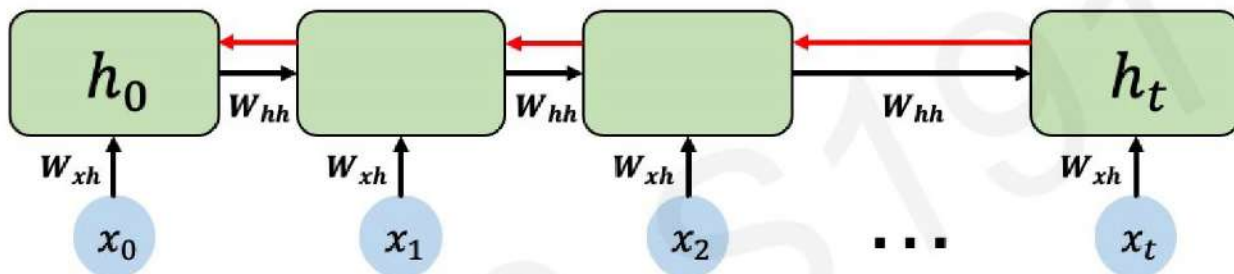
Computing the gradient wrt h_0 involves many factors of W_{hh} + repeated gradient computation!

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients



Standard RNN Gradient Flow: Vanishing Gradients



Computing the gradient wrt h_0 involves **many factors of W_{hh}** + repeated gradient computation!

Many values > 1 :
exploding gradients

Gradient clipping to
scale big gradients

Many values < 1 :
vanishing gradients

1. Activation function
2. Weight initialization
3. Network architecture



The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



The Problem of Long-Term Dependencies

"The clouds are in the ____"

Why are vanishing gradients a problem?

Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies



The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

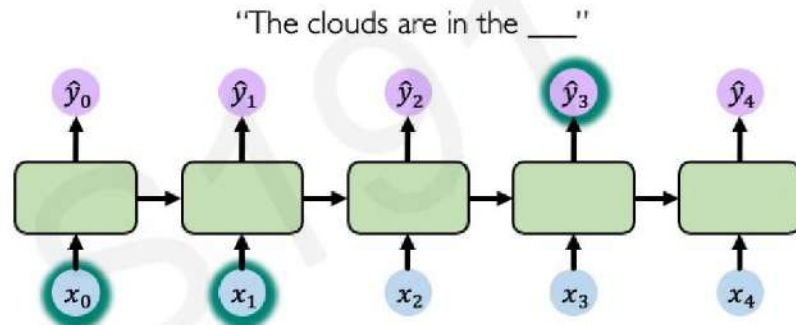
Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies





The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

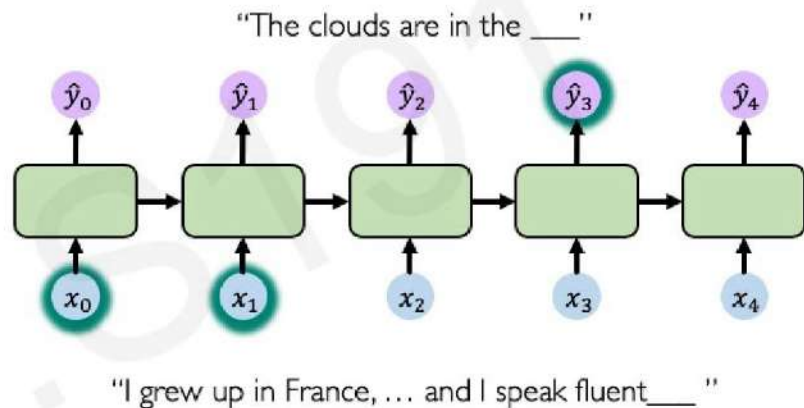
Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients



Bias parameters to capture short-term
dependencies





The Problem of Long-Term Dependencies

Why are vanishing gradients a problem?

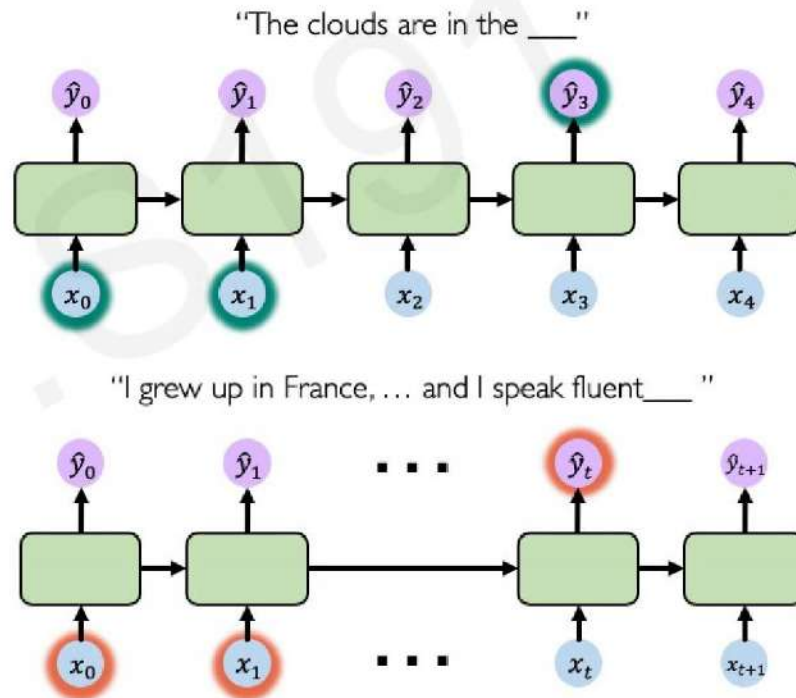
Multiply many **small numbers** together



Errors due to further back time steps
have smaller and smaller gradients

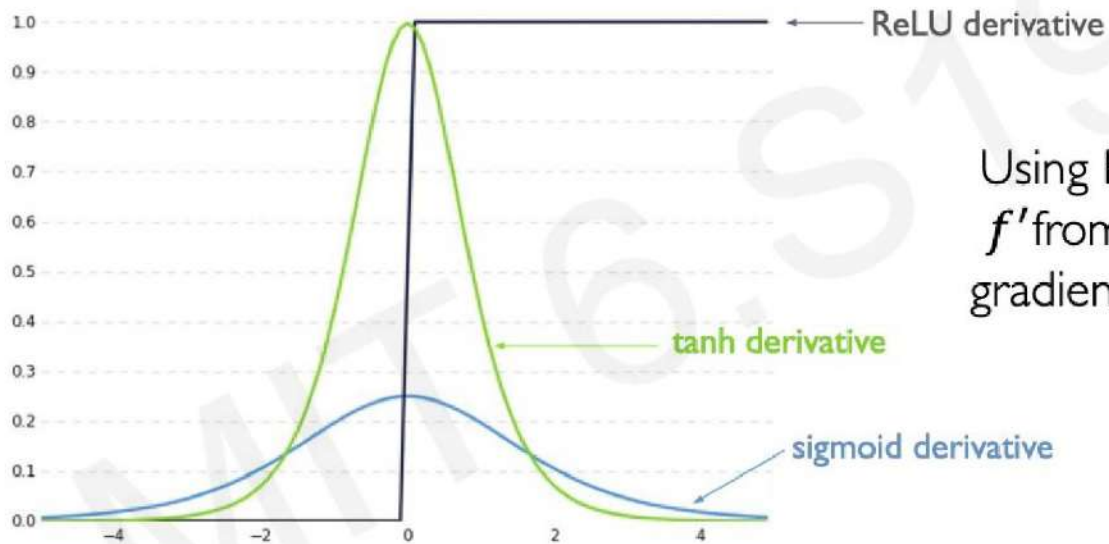


Bias parameters to capture short-term
dependencies





Trick #1: Activation Functions



Using ReLU prevents f' from shrinking the gradients when $x > 0$



Trick #2: Parameter Initialization

Initialize **weights** to identity matrix

Initialize **biases** to zero

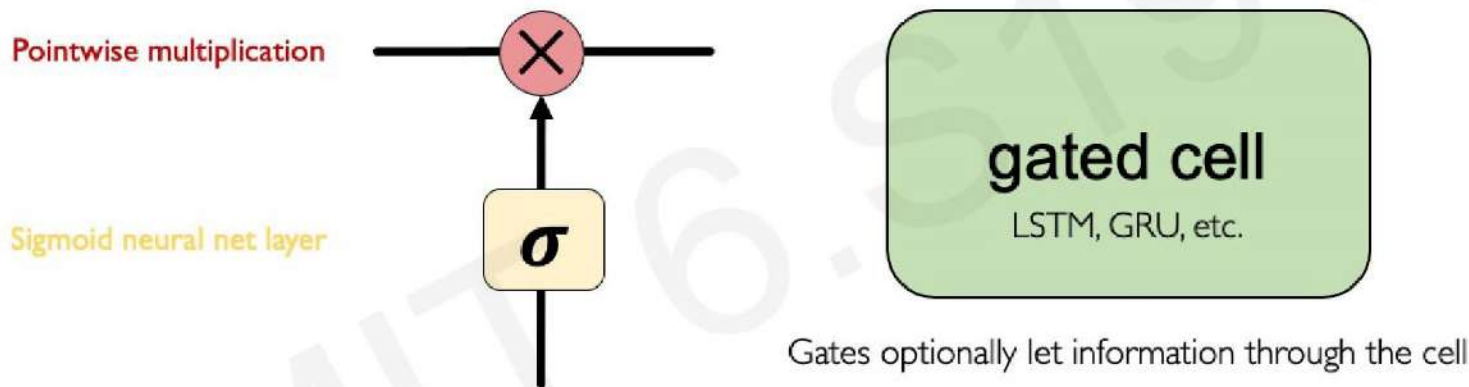
$$I_n = \begin{pmatrix} 1 & 0 & 0 & \cdots & 0 \\ 0 & 1 & 0 & \cdots & 0 \\ 0 & 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \cdots & 1 \end{pmatrix}$$

This helps prevent the weights from shrinking to zero.



Trick #3: Gated Cells

Idea: use **gates** to selectively **add** or **remove** information within **each recurrent unit** with



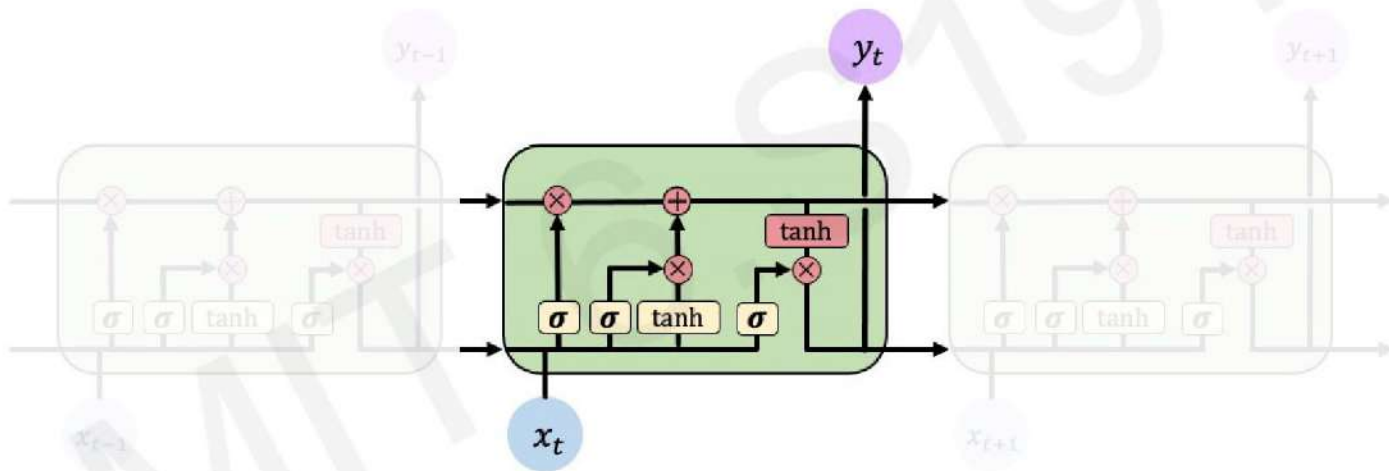
Long Short Term Memory (LSTMs) networks rely on a gated cell to track information throughout many time steps.



Long Short Term Memory (LSTMs)

Gated LSTM cells control information flow:

- 1) Forget
- 2) Store
- 3) Update
- 4) Output



LSTM cells are able to track information throughout many timesteps



```
tf.keras.layers.LSTM(num_units)
```



LSTMs: Key Concepts

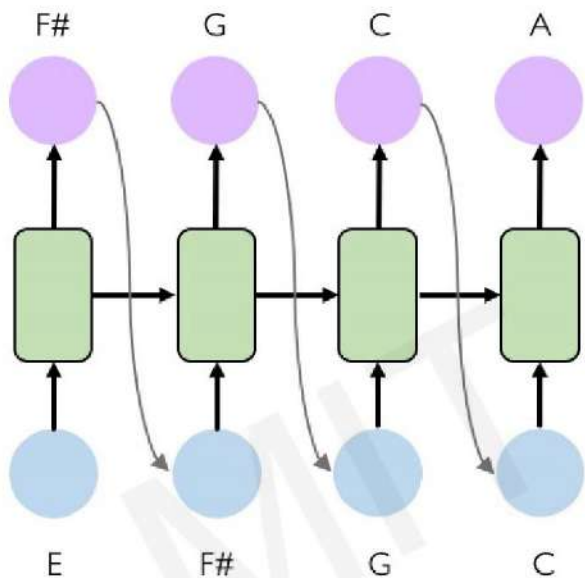
1. Maintain a **cell state**
2. Use **gates** to control the **flow of information**
 - **Forget** gate gets rid of irrelevant information
 - **Store** relevant information from current input
 - Selectively **update** cell state
 - **Output** gate returns a filtered version of the cell state
3. Backpropagation through time with partially **uninterrupted gradient flow**



RNN Applications & Limitations



Example Task: Music Generation



Input: sheet music

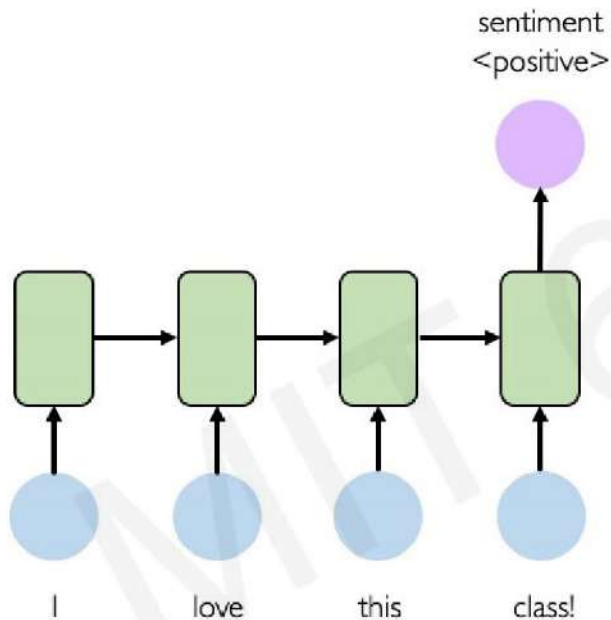
Output: next character in sheet music

Listening to
3rd movement






Example Task: Sentiment Classification



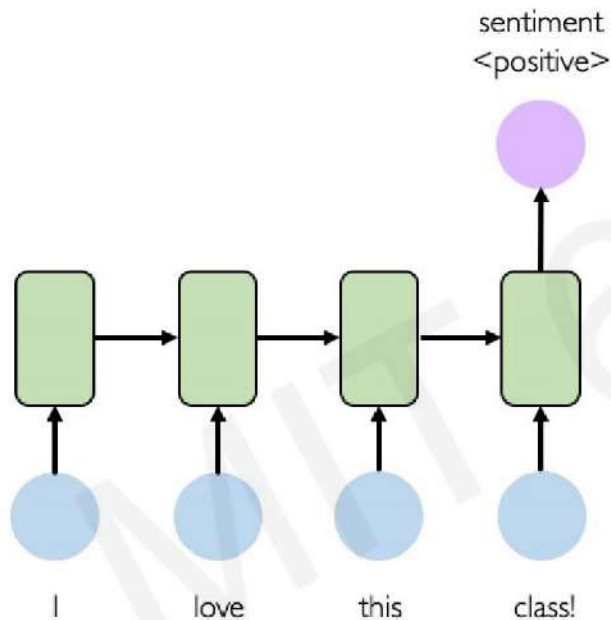
Input: sequence of words

Output: probability of having positive sentiment

```
 loss = tf.nn.softmax_cross_entropy_with_logits(y, predicted)
```



Example Task: Sentiment Classification



Tweet sentiment classification



Ivar Hagendoorn
@IvarHagendoorn

Follow



The @MIT Introduction to #DeepLearning is definitely one of the best courses of its kind currently available online
introtodeeplearning.com

12:45 PM - 12 Feb 2018



Angels-Cave
@AngelsCave

Follow



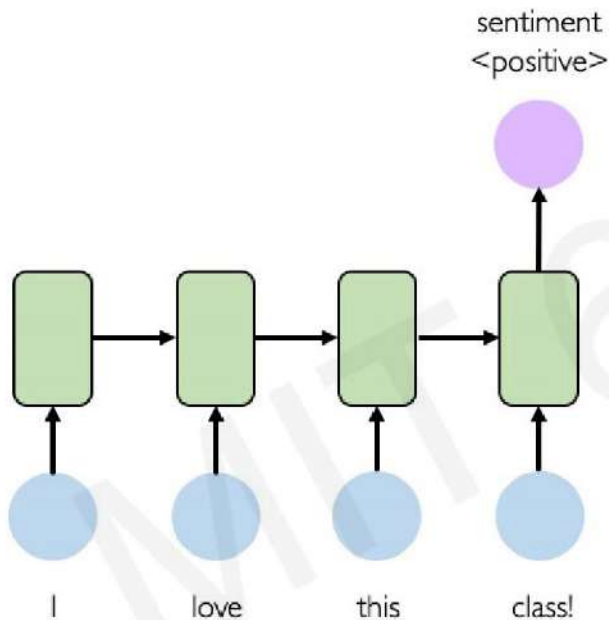
Replying to @Kazuki2048

I wouldn't mind a bit of snow right now. We haven't had any in my bit of the Midlands this winter! :(




2:19 AM - 25 Jan 2019



Limitations of Recurrent Models



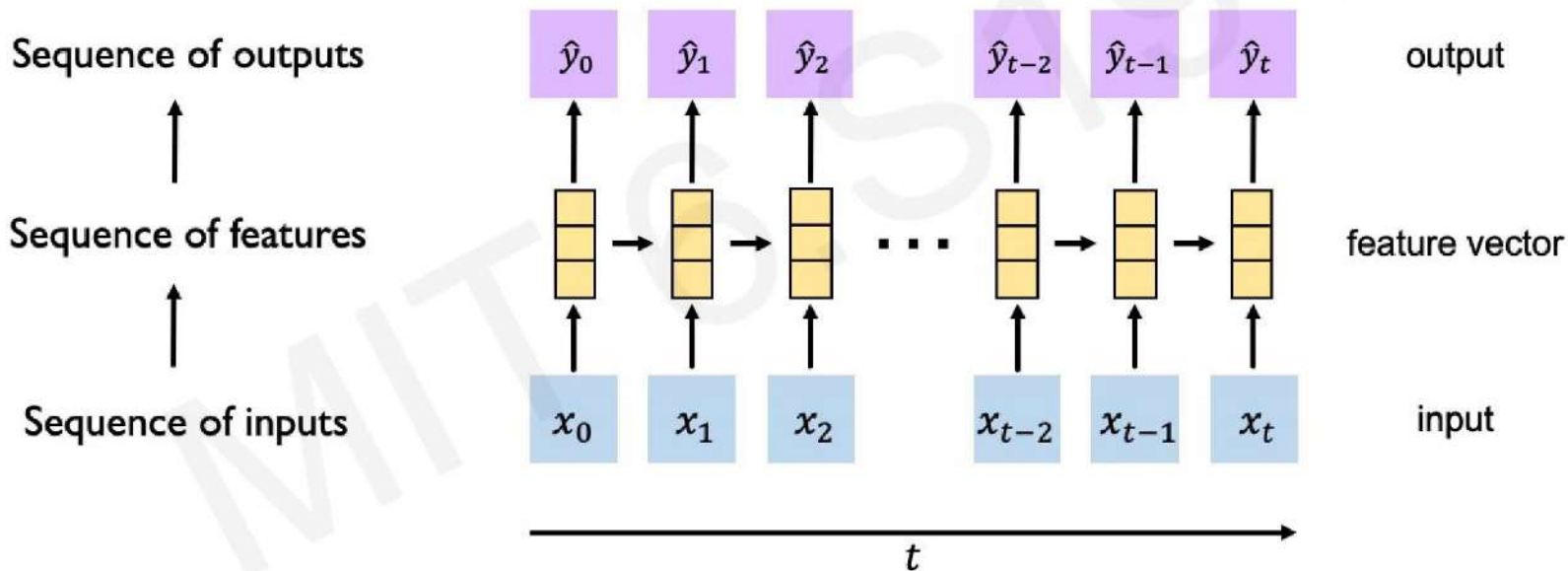
Limitations of RNNs

-  Encoding bottleneck
-  Slow, no parallelization
-  Not long memory



Goal of Sequence Modeling

RNNs: recurrence to model sequence dependencies






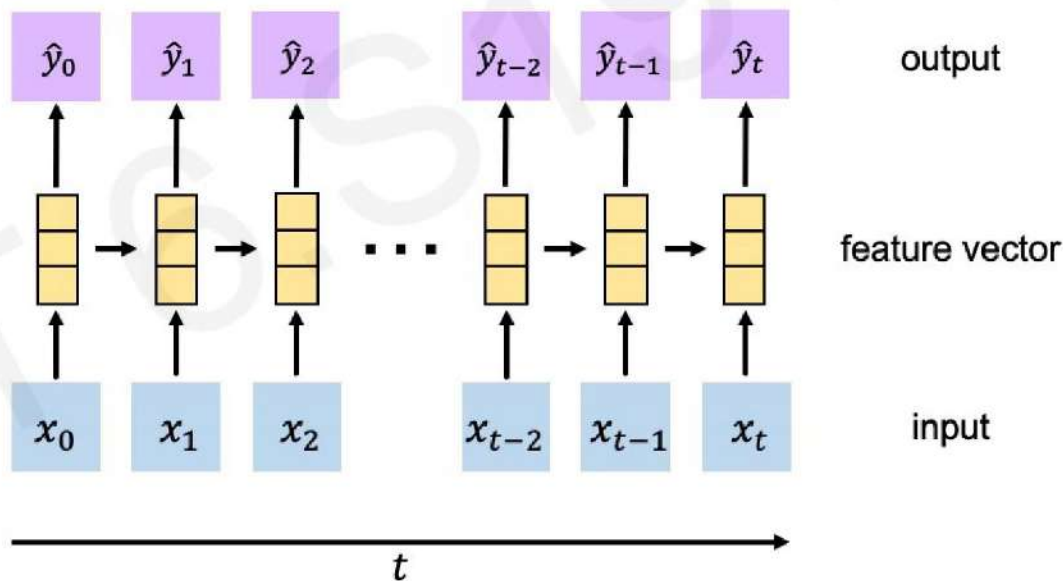


Goal of Sequence Modeling

RNNs: recurrence to model sequence dependencies

Limitations of RNNs

-  Encoding bottleneck
-  Slow, no parallelization
-  Not long memory





Goal of Sequence Modeling

Can we eliminate the need for recurrence entirely?

Desired Capabilities



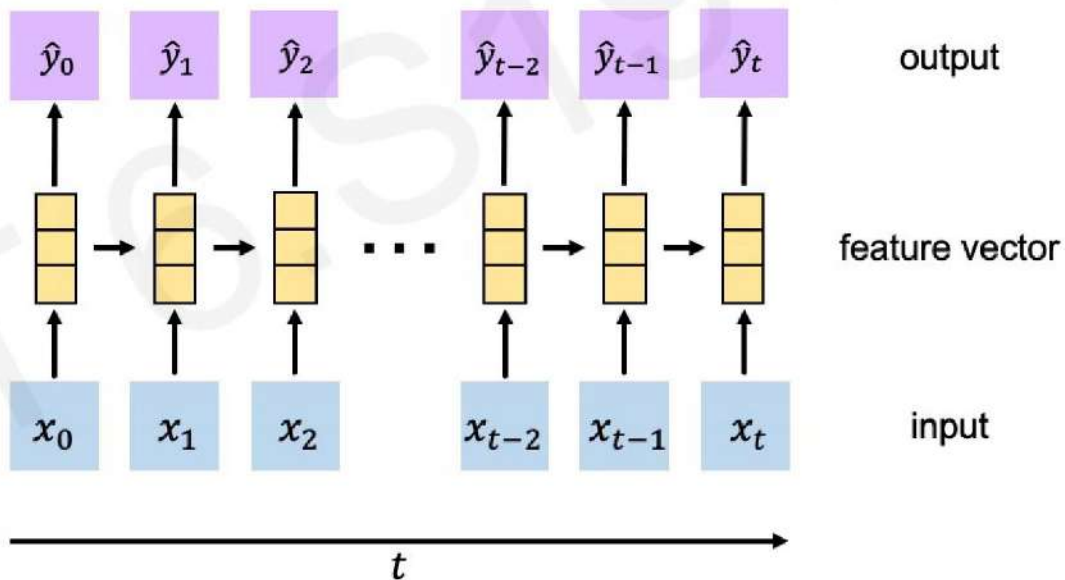
Continuous stream



Parallelization



Long memory





Goal of Sequence Modeling

Can we eliminate the need for recurrence entirely?

Desired Capabilities



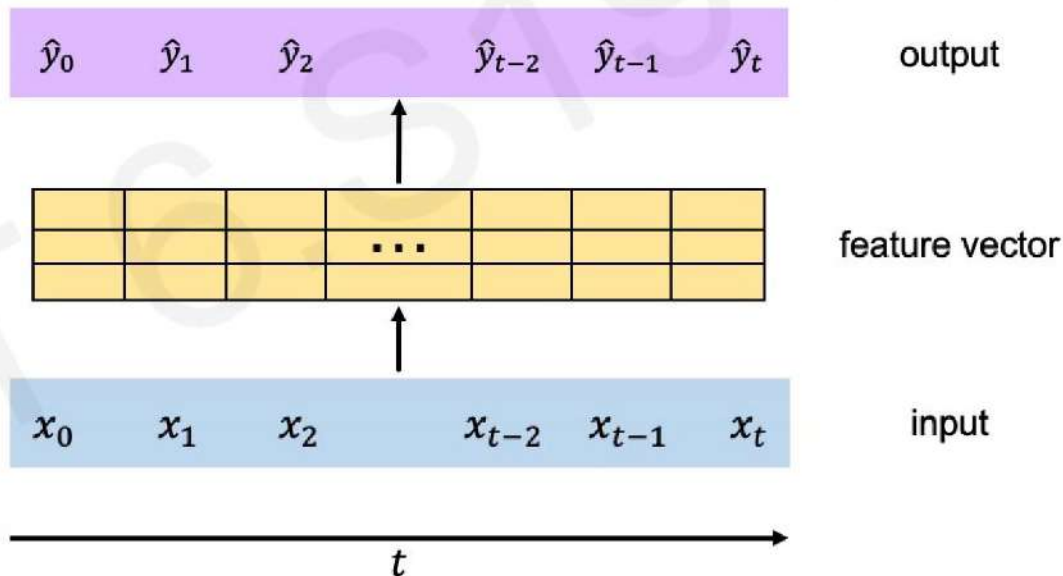
Continuous stream



Parallelization



Long memory





Goal of Sequence Modeling

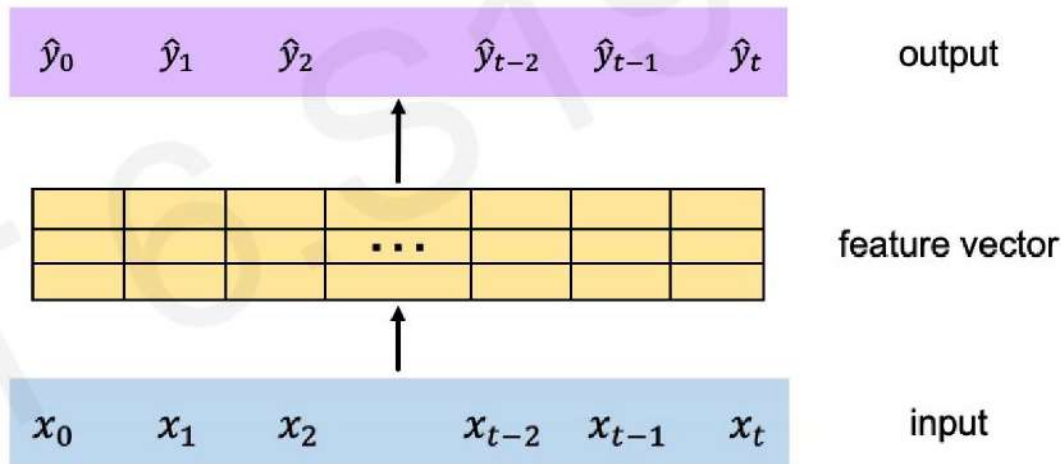
Idea 1: Feed everything into dense network

- ✓ No recurrence
- ✗ Not scalable
- ✗ No order
- ✗ No long memory



Idea: Identify and attend to what's important

Can we eliminate the need for recurrence entirely?





人工智能基本理论

Attention Is All You Need

MIT 6.S191



Intuition Behind Self-Attention

Attending to the most important parts of an input.



1. Identify which parts to attend to
2. Extract the features with high attention

Similar to a search problem!

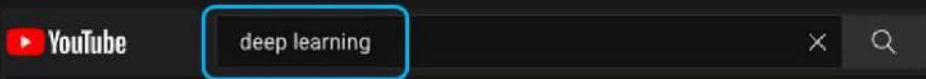


A Simple Example: Search





Understanding Attention with Search



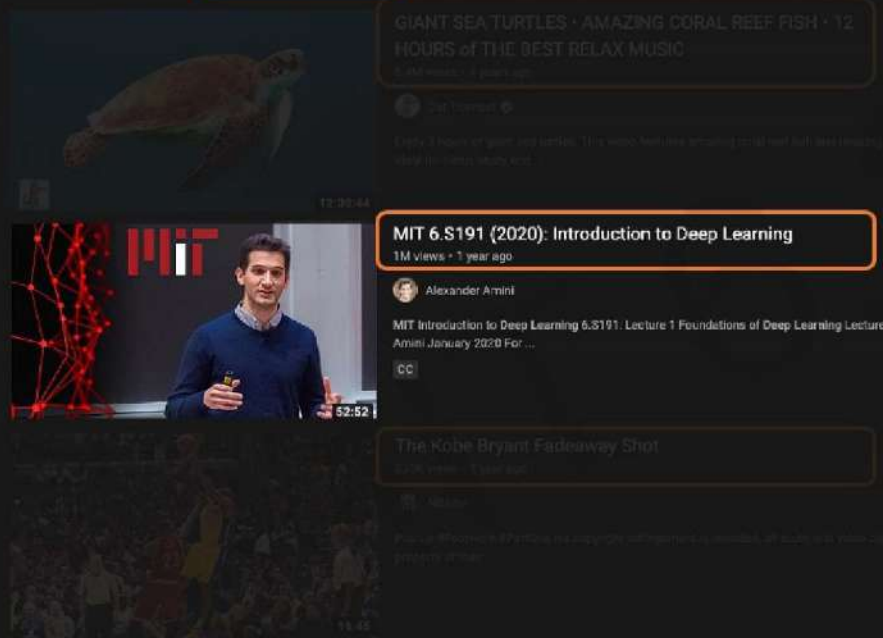
Query (Q)

How similar is the key to the query?

Key (K₁)

Key (K₂)

Key (K₃)



1. **Compute attention mask:** how similar is each key to the desired query?



Understanding Attention with Search

YouTube

deep learning



Query (Q)

Key (K_1)

Key (K_2)

Value (V)

Key (K_3)

2. Extract values based on attention:
Return the values highest attention



Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention



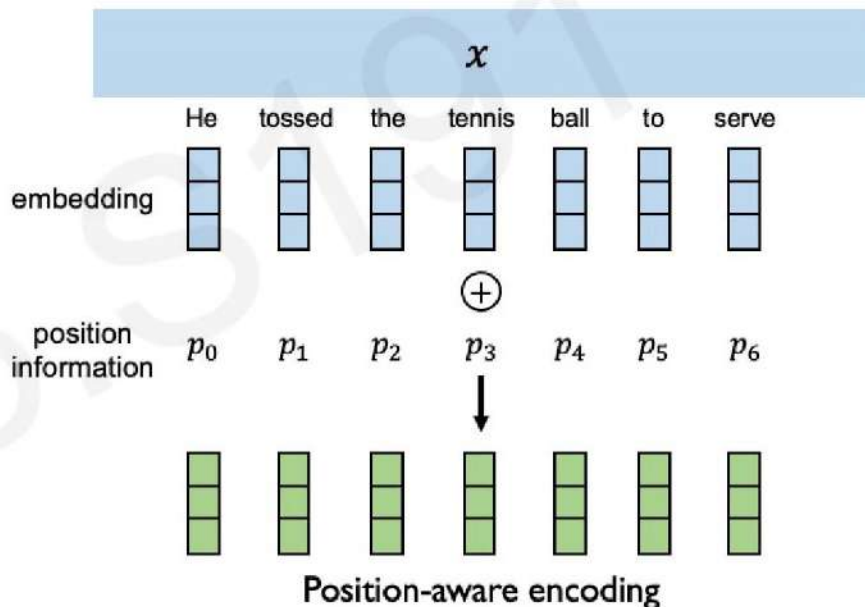
Data is fed in all at once! Need to encode position information to understand order.



Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute attention weighting
4. Extract features with high attention



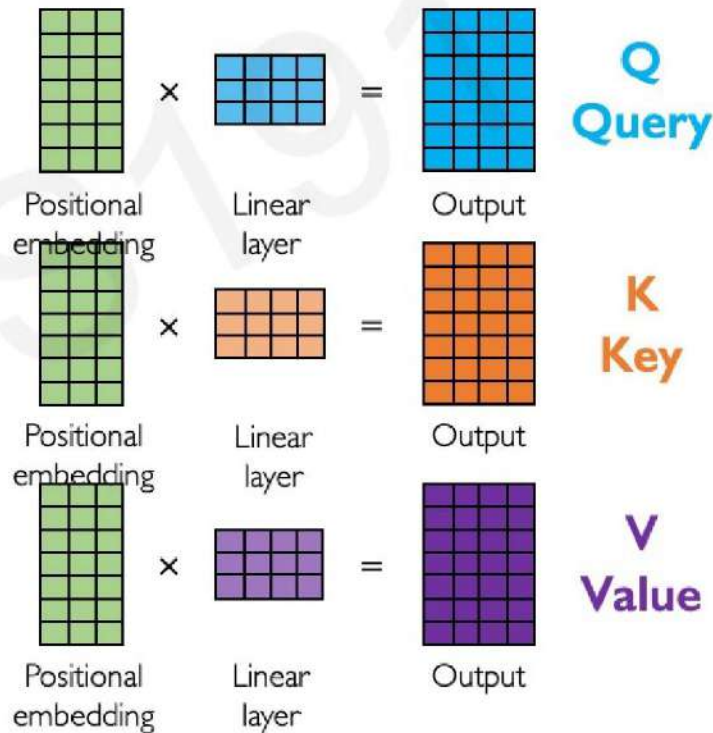
Data is fed in all at once! Need to encode position information to understand order.



Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute attention weighting
4. Extract features with high attention





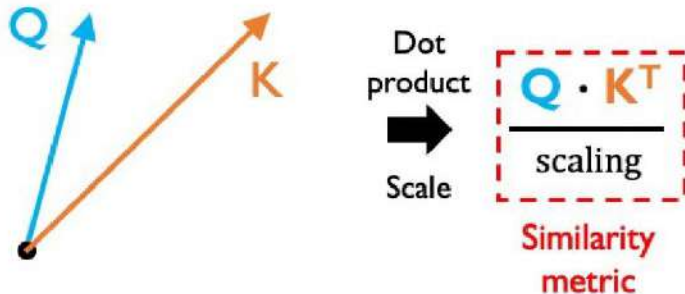
Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

Attention score: compute pairwise similarity between each **query** and **key**

How to compute similarity between two sets of features?



Also known as the “cosine similarity”



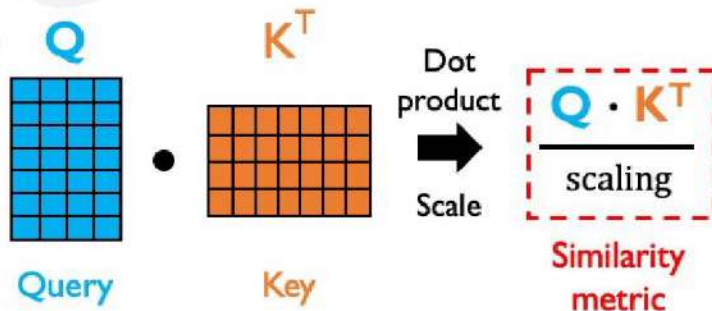
Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

Attention score: compute pairwise similarity between each **query** and **key**

How to compute similarity between two sets of features?



Also known as the “cosine similarity”

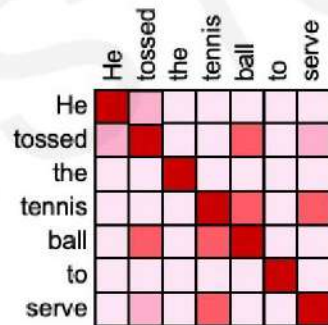


Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract features with high attention

Attention weighting: where to attend to!
How similar is the key to the query?



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right)$$

Attention weighting

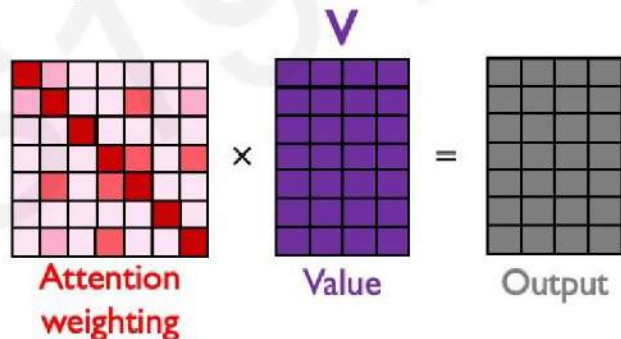


Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract **features with high attention**

Last step: self-attend to extract features



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V = A(Q, K, V)$$

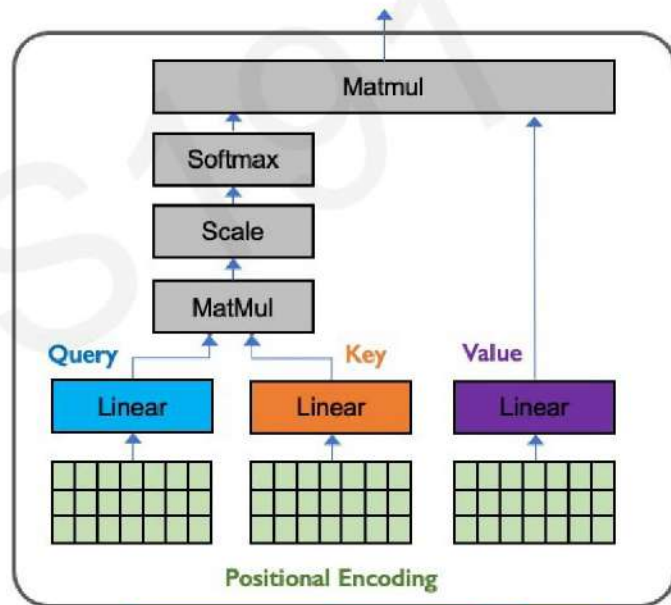


Learning Self-Attention with Neural Networks

Goal: identify and attend to most important features in input.

1. Encode **position** information
2. Extract **query, key, value** for search
3. Compute **attention weighting**
4. Extract **features with high attention**

These operations form a self-attention head that can plug into a larger network. Each head attends to a different part of input.



$$\text{softmax} \left(\frac{Q \cdot K^T}{\text{scaling}} \right) \cdot V$$



Applying Multiple Self-Attention Heads



Attention weighting

×



Value

=



Output



Output of attention head 1



Output of attention head 2



Output of attention head 3



Self-Attention Applied

Language Processing

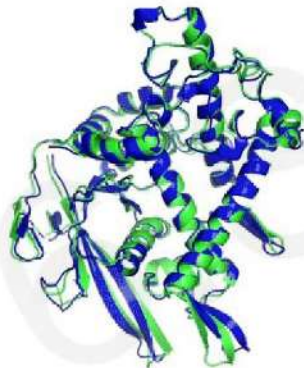


An armchair in the shape of an avocado

BERT, GPT-3

Devlin et al., *NAACL* 2019
Brown et al., *NeurIPS* 2020

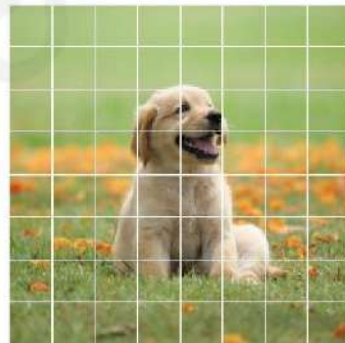
Biological Sequences



AlphaFold2

Jumper et al., *Nature* 2021

Computer Vision



Vision Transformers

Dosovitskiy et al., *ICLR* 2020



Deep Learning for Sequence Modeling: Summary

1. RNNs are well suited for **sequence modeling** tasks
2. Model sequences via a **recurrence relation**
3. Training RNNs with **backpropagation through time**
4. Models for **music generation**, classification, machine translation, and more
5. Self-attention to model **sequences without recurrence**



提纲

一、前馈网络

二、循环网络

三、卷积网络



上海大学
SHANGHAI UNIVERSITY



人工智能基本理论

从一个故事讲起.....



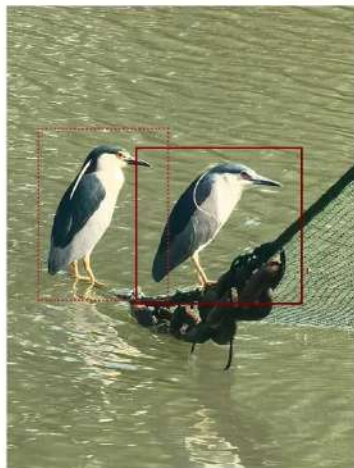


人工智能基本理论



15:38 网络信号 24

懂鸟·全球鸟类识别



拍摄地点: 未指定



yè lù
夜鹭

Black-crowned Night Heron
Nycticorax nycticorax
识别可信度: 92.8%



拍照识别



相册识别



群图片



听音辨鸟



浏览搜索



| 人工智能基本理论

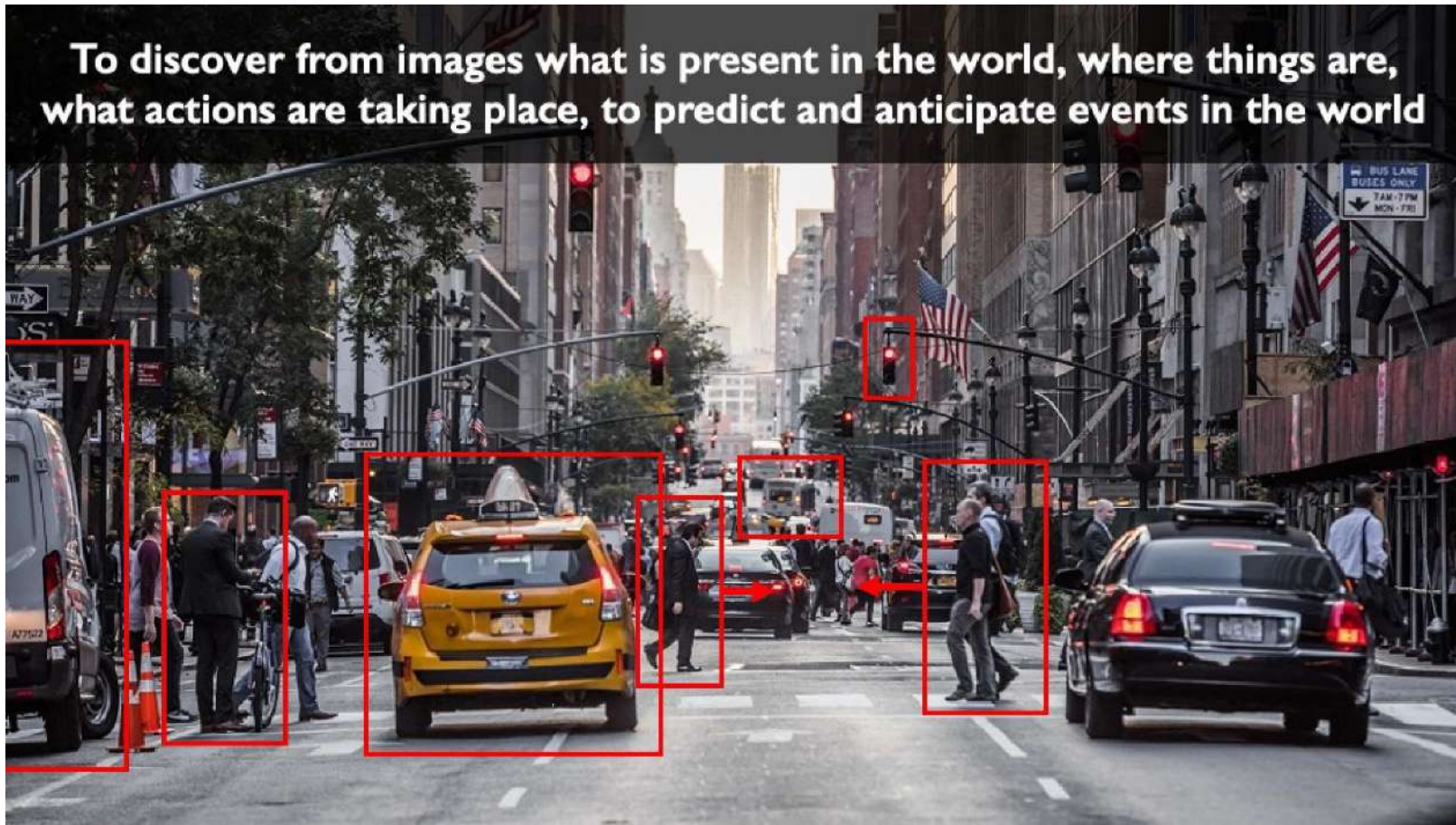
**“To know what is
where by looking.”**





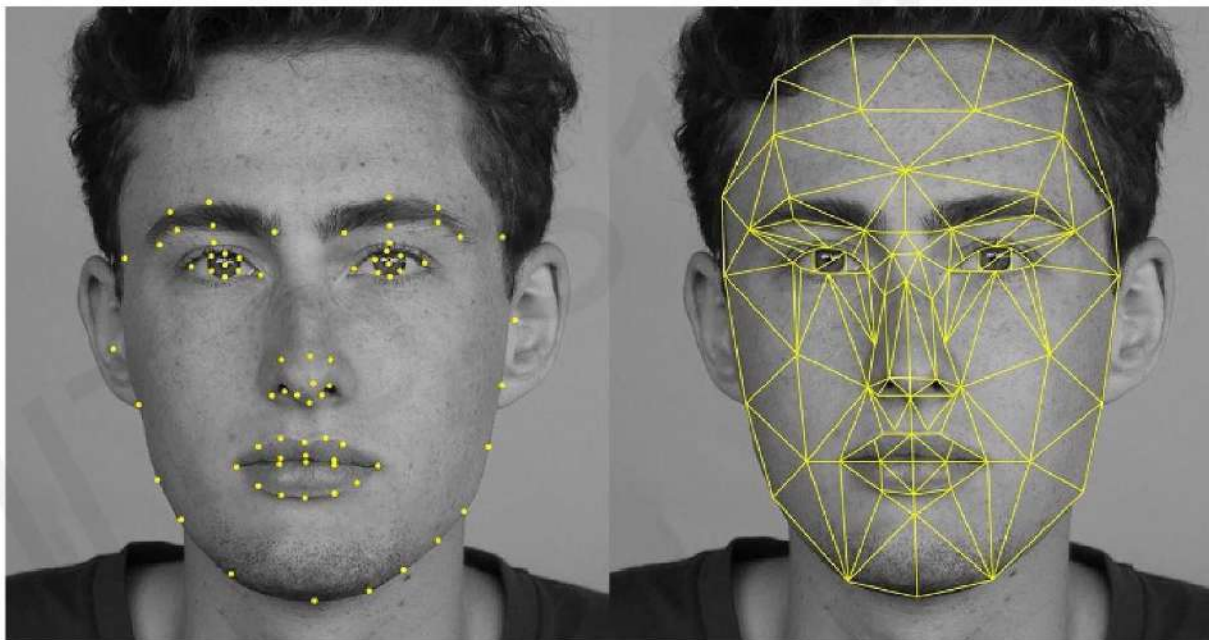
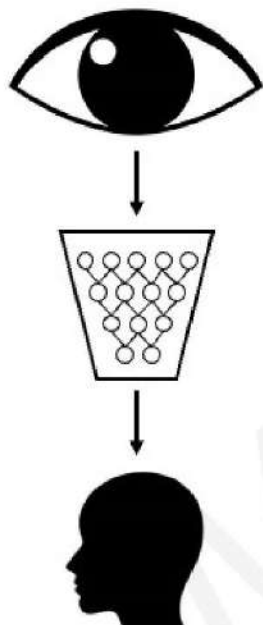
人工智能基本理论

To discover from images what is present in the world, where things are, what actions are taking place, to predict and anticipate events in the world



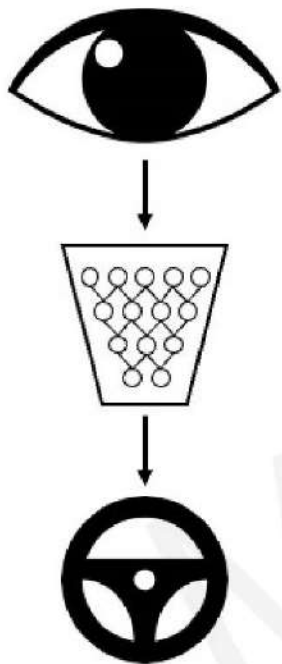


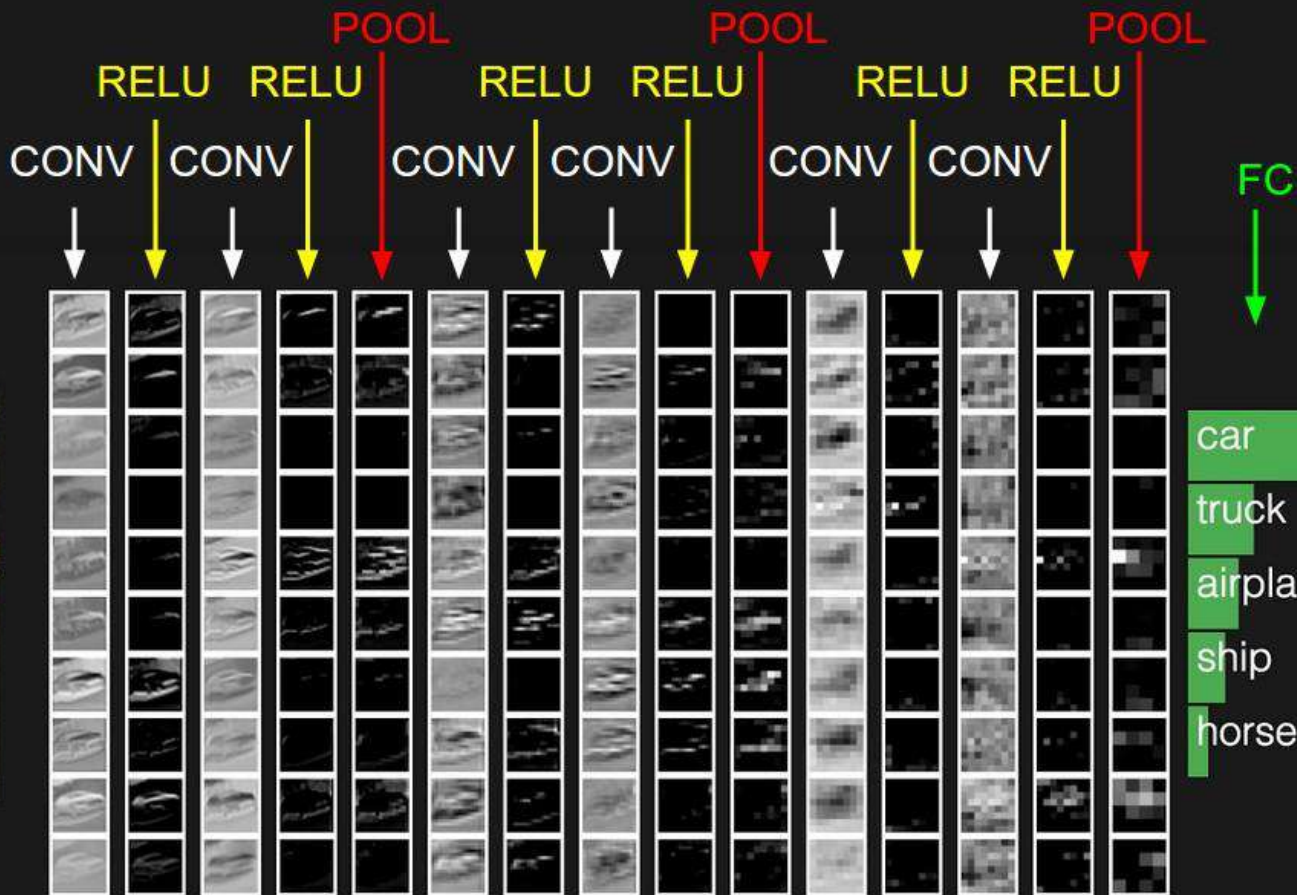
Impact: Facial Detection & Recognition





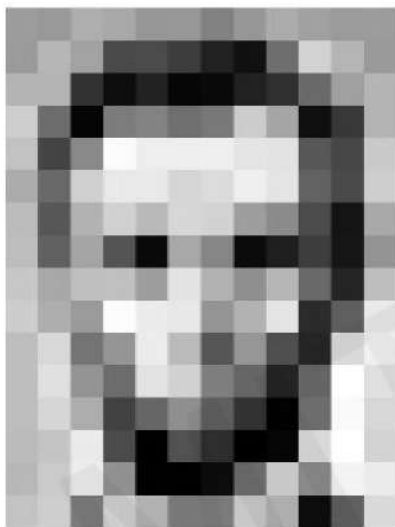
Impact: Self-Driving Cars







Images are Numbers



157	153	174	168	150	152	129	151	172	167	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	35	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	239	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

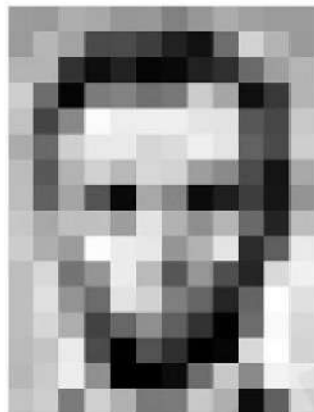
What the computer sees

157	153	174	168	150	152	129	151	172	167	155	156
155	182	163	74	75	62	33	17	110	210	180	154
180	180	50	14	34	6	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	137	251	237	239	239	228	227	87	71	201
172	105	207	233	233	214	220	239	228	98	74	206
188	88	179	209	185	215	211	158	139	75	20	169
189	97	165	84	10	168	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	36	190
205	174	155	252	236	231	149	178	228	43	95	234
190	216	116	149	236	187	85	150	79	38	218	241
190	224	147	108	227	210	127	102	35	101	255	224
190	214	173	66	103	143	95	50	2	109	249	215
187	196	239	75	1	81	47	0	6	217	255	211
183	202	237	145	0	0	12	108	209	138	243	236
195	206	123	207	177	121	123	200	175	13	96	218

An image is just a matrix of numbers $[0,255]$!
i.e., $1080 \times 1080 \times 3$ for an RGB image



Tasks in Computer Vision



Input Image



147	163	174	188	156	162	129	161	172	163	166	156
195	182	163	74	76	62	35	17	113	216	180	154
180	180	50	14	24	0	10	33	48	106	159	181
206	109	5	124	131	111	120	204	166	15	56	180
194	68	127	201	237	230	239	228	227	87	71	201
172	106	207	233	233	214	230	239	239	84	74	206
188	88	178	208	186	215	211	168	138	76	20	169
189	97	165	84	16	148	134	11	31	62	22	148
199	168	191	193	158	227	178	143	182	106	26	190
205	134	155	202	236	231	143	178	228	45	95	234
190	210	118	149	236	187	85	100	79	28	216	241
190	234	147	108	227	210	127	102	36	161	226	224
190	214	173	66	103	140	96	63	3	108	240	215
187	196	205	76	1	81	47	0	6	217	246	211
188	202	237	145	0	0	12	108	208	138	243	236
195	206	123	207	177	123	123	200	175	13	96	218

Pixel Representation

classification

Lincoln

0.8

Washington

0.1

Jefferson

0.05

Obama

0.05

- **Regression:** output variable takes continuous value
- **Classification:** output variable takes class label. Can produce probability of belonging to a particular class



High Level Feature Detection

Let's identify key features in each image category



Nose,
Eyes,
Mouth



Wheels,
License Plate,
Headlights



Door,
Windows,
Steps



Manual Feature Extraction

Domain knowledge

Define features

Detect features to classify

Viewpoint variation



Scale variation



Deformation



Occlusion



Illumination conditions



Background clutter



Intra-class variation

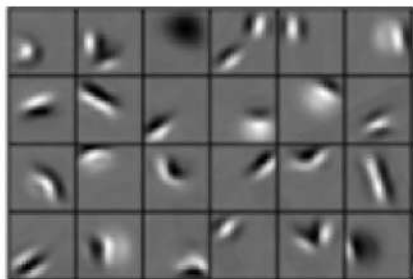




Learning Feature Representations

Can we learn a **hierarchy of features** directly from the data instead of hand engineering?

Low level features



Edges, dark spots

Mid level features



Eyes, ears, nose

High level features



Facial structure



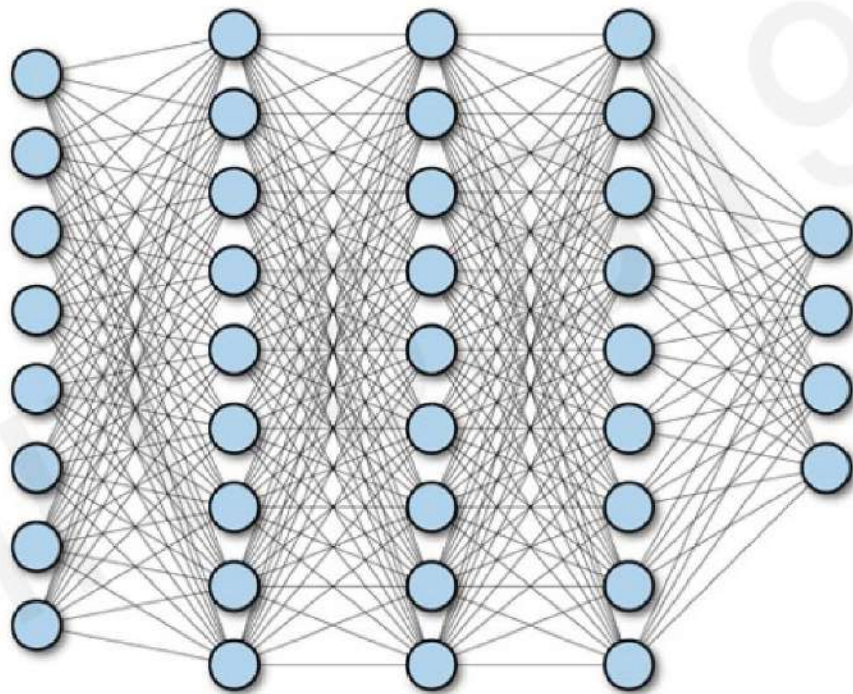
人工智能基本理论

Learning Visual Features

MIT 6.S191



Fully Connected Neural Network

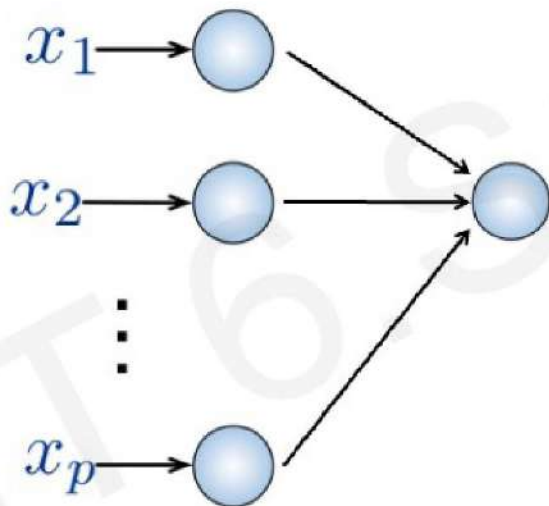




Fully Connected Neural Network

Input:

- 2D image
- Vector of pixel values



Fully Connected:

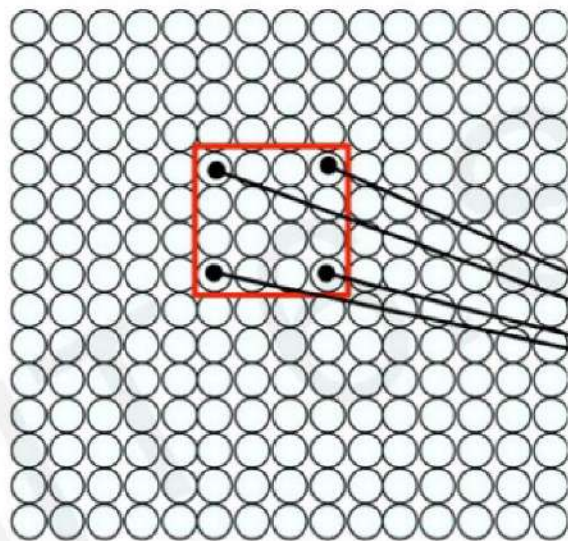
- Connect neuron in hidden layer to all neurons in input layer
- No spatial information!
- And many, many parameters!

How can we use **spatial structure** in the input to inform the architecture of the network?



Using Spatial Structure

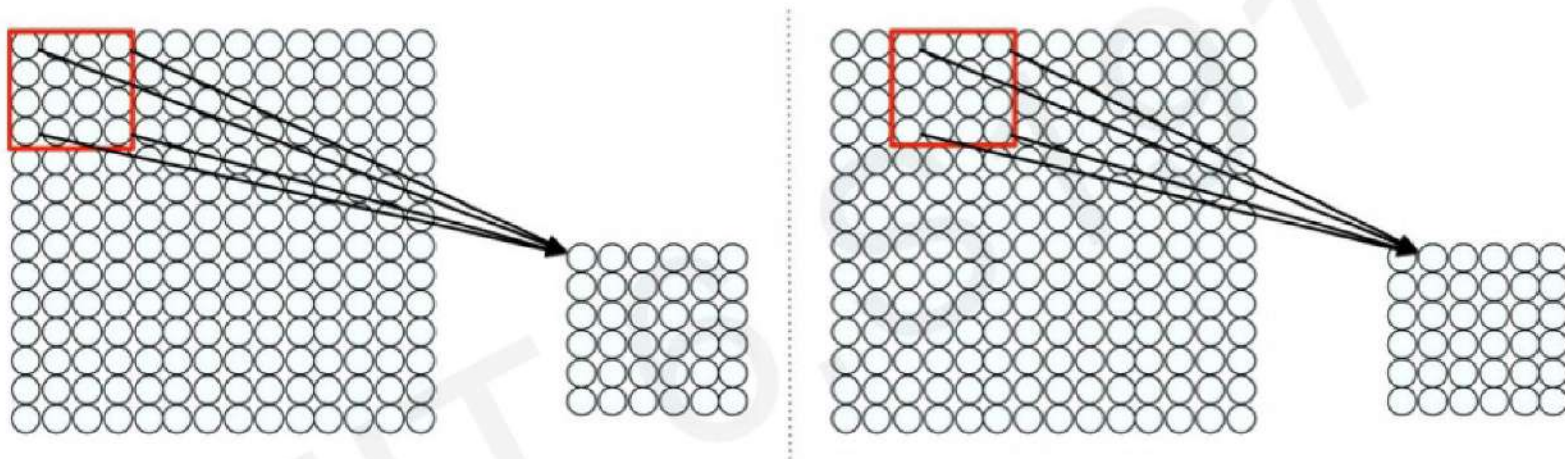
Input: 2D image.
Array of pixel values



Idea: connect patches of input
to neurons in hidden layer.
Neuron connected to region of
input. Only “sees” these values.



Using Spatial Structure



Connect patch in input layer to a single neuron in subsequent layer.

Use a sliding window to define connections.

How can we **weight** the patch to detect particular features?

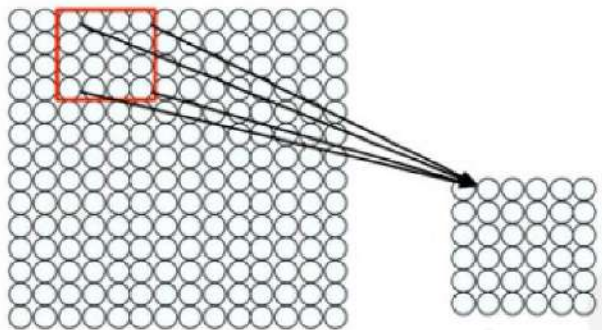


Applying Filters to Extract Features

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) Spatially **share** parameters of each filter
(features that matter in one part of the input should matter elsewhere)



Feature Extraction with Convolution



- Filter of size 4x4 : 16 different weights
- Apply this same filter to 4x4 patches in input
- Shift by 2 pixels for next patch

This “patchy” operation is **convolution**

- 1) Apply a set of weights – a filter – to extract **local features**
- 2) Use **multiple filters** to extract different features
- 3) **Spatially share** parameters of each filter



The Convolution Operation

We slide the 3x3 filter over the input image, element-wise multiply, and add the outputs:

1	1	1	0	0
0	1	1	1	0
0	0	1 _{x1}	1 _{x0}	1 _{x1}
0	0	1 _{x0}	1 _{x1}	0 _{x0}
0	1	1 _{x1}	0 _{x0}	0 _{x1}



1	0	1
0	1	0
1	0	1

filter



4	3	4
2	4	3
2	3	4

feature map



Producing Feature Maps



Original



Sharpen



Edge Detect



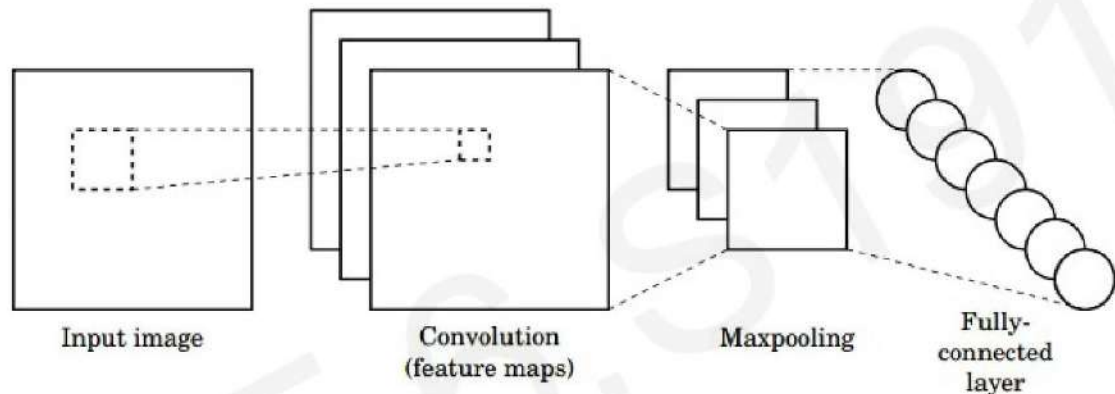
“Strong” Edge Detect



Convolutional Neural Networks (CNNs)



CNNs for Classification



1. **Convolution**: Apply filters to generate feature maps.
2. **Non-linearity**: Often ReLU.
3. **Pooling**: Downsampling operation on each feature map.



```
tf.keras.layers.Conv2D
```



```
tf.keras.activations.*
```

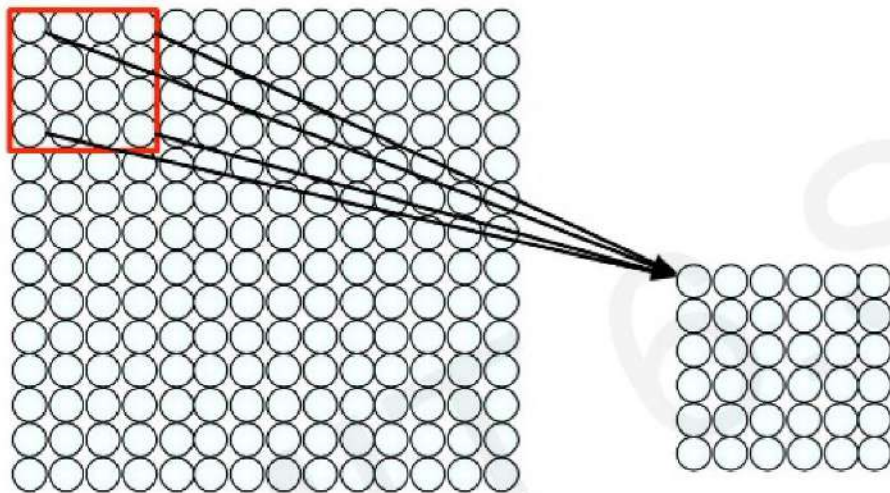



```
tf.keras.layers.MaxPool2D
```

Train model with image data.
Learn weights of filters in convolutional layers.



Convolutional Layers: Local Connectivity



 `tf.keras.layers.Conv2D`

For a neuron in hidden layer:

- Take inputs from patch
- Compute weighted sum
- Apply bias

4x4 filter: matrix
of weights w_{ij}

$$\sum_{i=1}^4 \sum_{j=1}^4 w_{ij} x_{i+p,j+q} + b$$

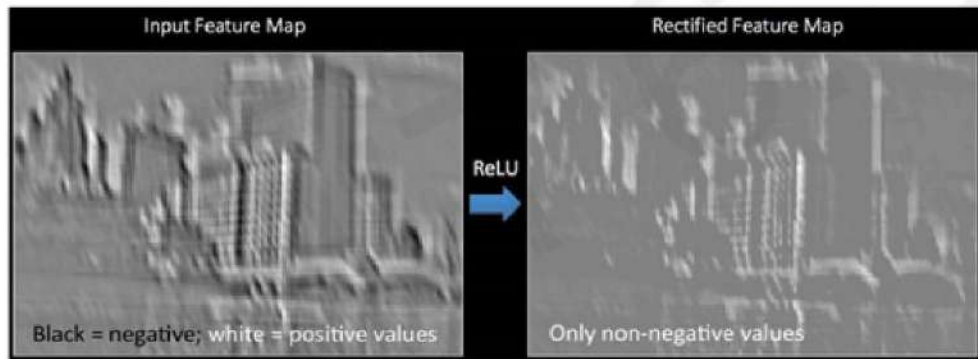
for neuron (p,q) in hidden layer

- 1) applying a window of weights
- 2) computing linear combinations
- 3) activating with non-linear function

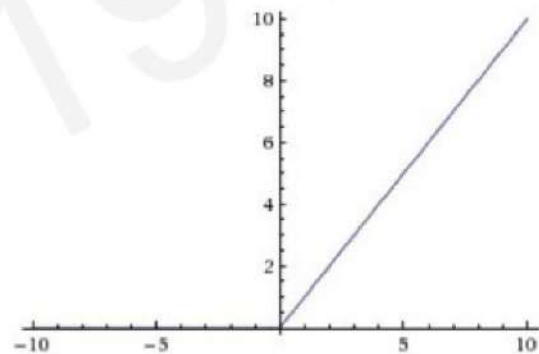


Introducing Non-Linearity

- Apply after every convolution operation (i.e., after convolutional layers)
- ReLU: pixel-by-pixel operation that replaces all negative values by zero. **Non-linear operation!**



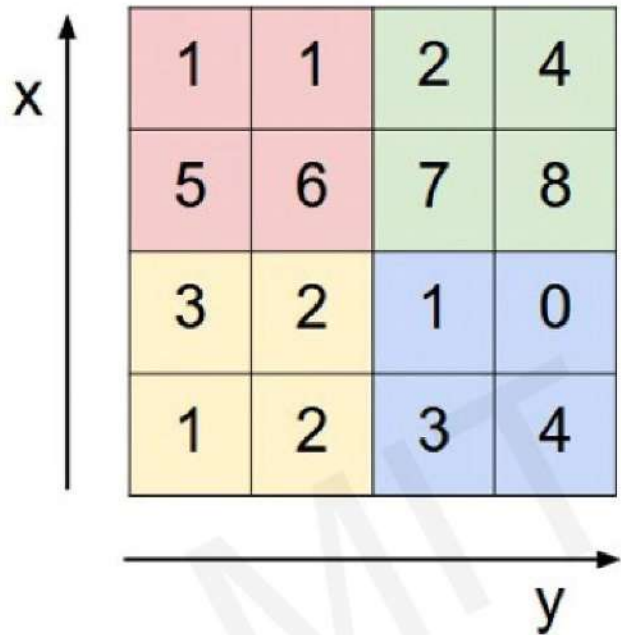
Rectified Linear Unit (ReLU)



`tf.keras.layers.ReLU`



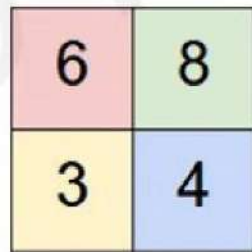
Pooling



max pool with 2x2 filters
and stride 2



```
tf.keras.layers.MaxPool2D(  
    pool_size=(2,2),  
    strides=2  
)
```

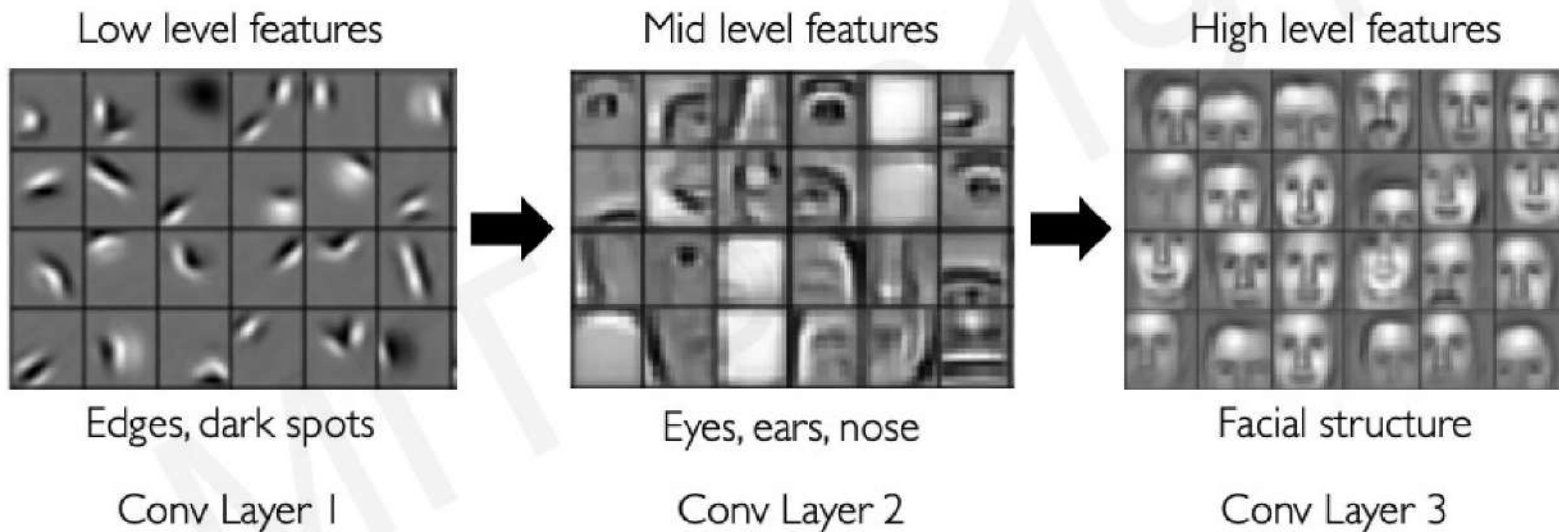


- 1) Reduced dimensionality
- 2) Spatial invariance

How else can we downsample and preserve spatial invariance?

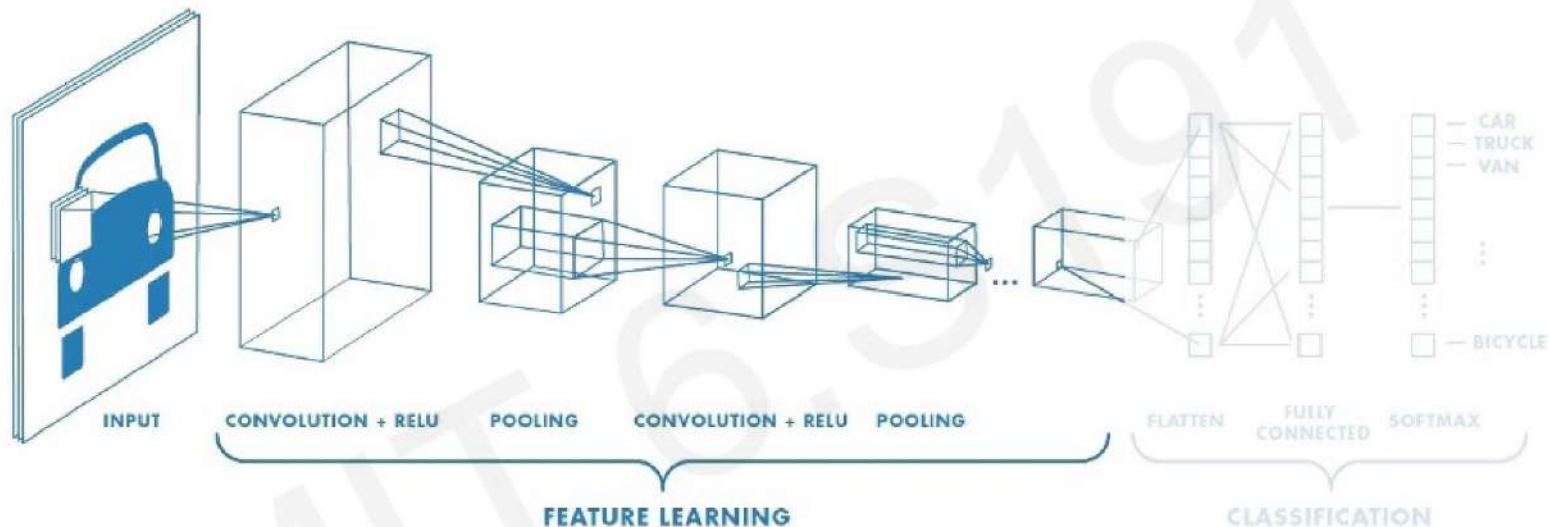


Representation Learning in Deep CNNs





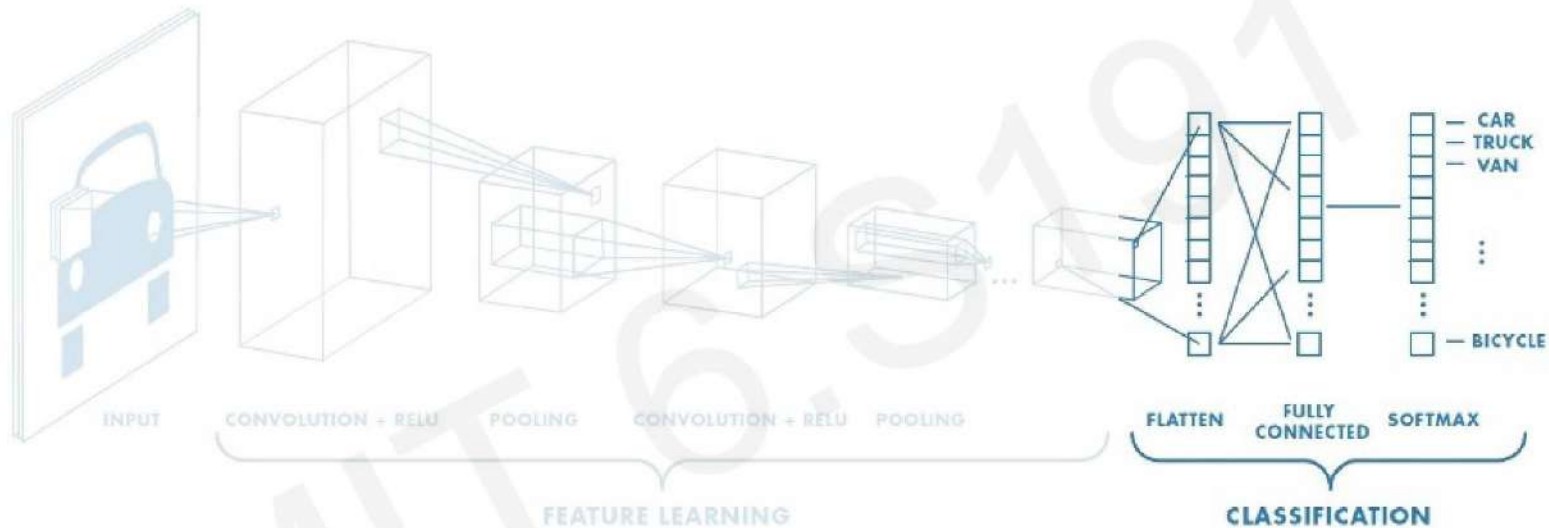
CNNs for Classification: Feature Learning



1. Learn features in input image through **convolution**
2. Introduce **non-linearity** through activation function (real-world data is non-linear!)
3. Reduce dimensionality and preserve spatial invariance with **pooling**



CNNs for Classification: Class Probabilities



- CONV and POOL layers output high-level features of input
- Fully connected layer uses these features for classifying input image
- Express output as **probability** of image belonging to a particular class

$$\text{softmax}(y_i) = \frac{e^{y_i}}{\sum_j e^{y_j}}$$

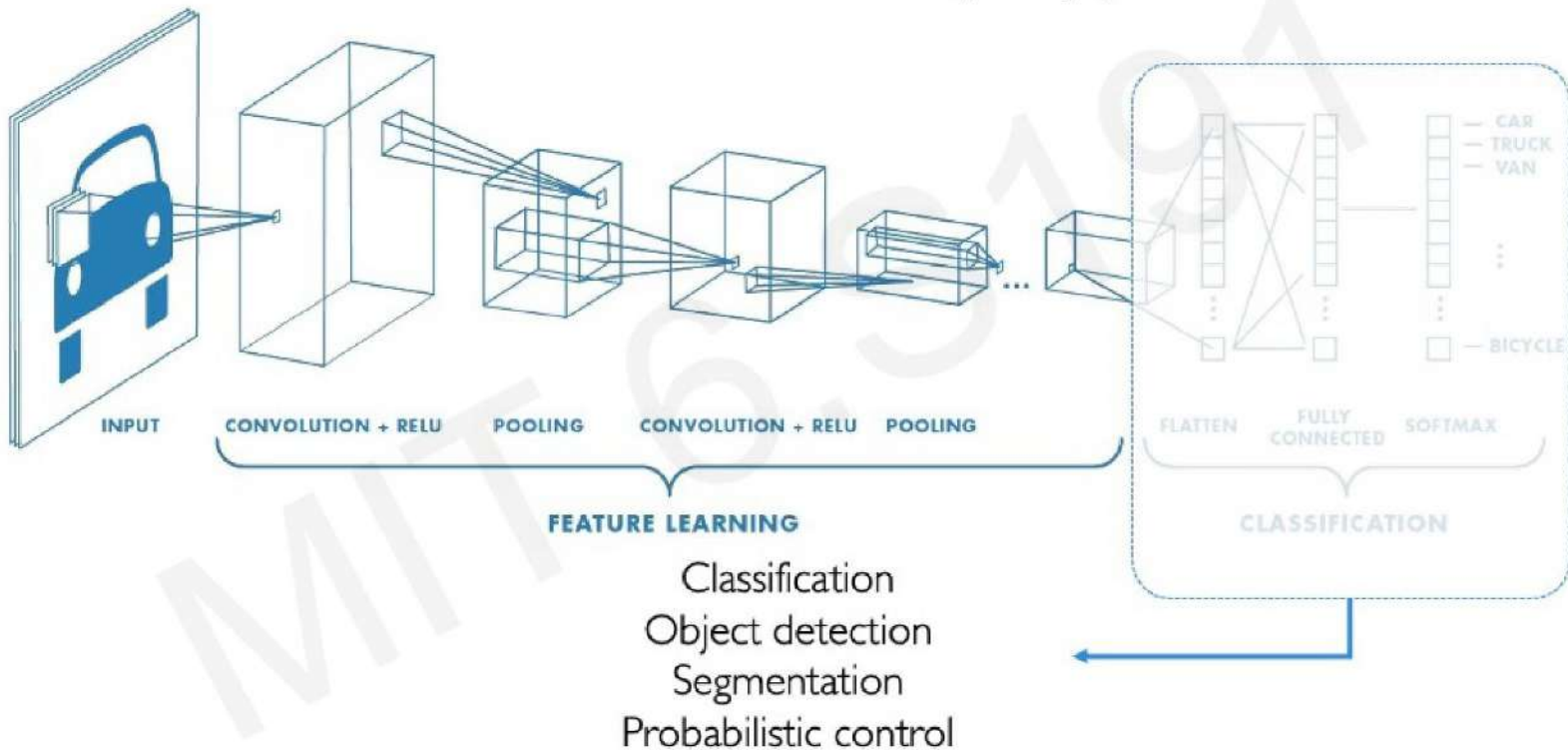


人工智能基本理论

An Architecture for Many Applications



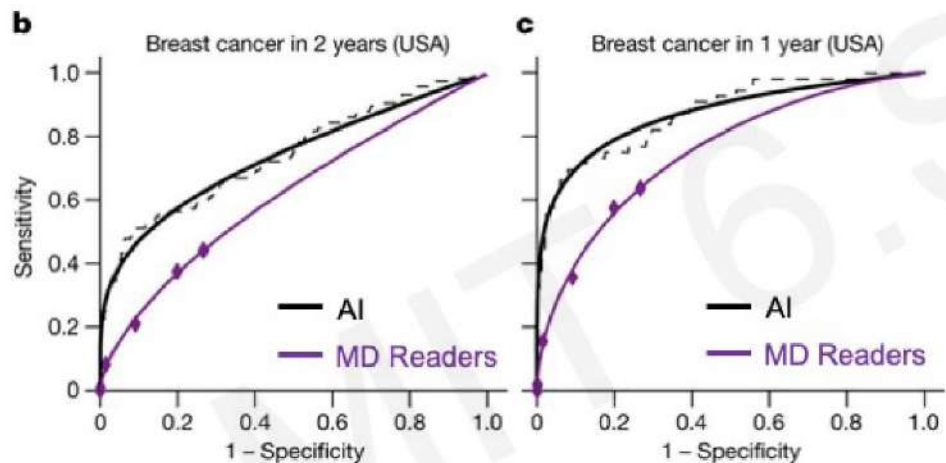
An Architecture for Many Applications



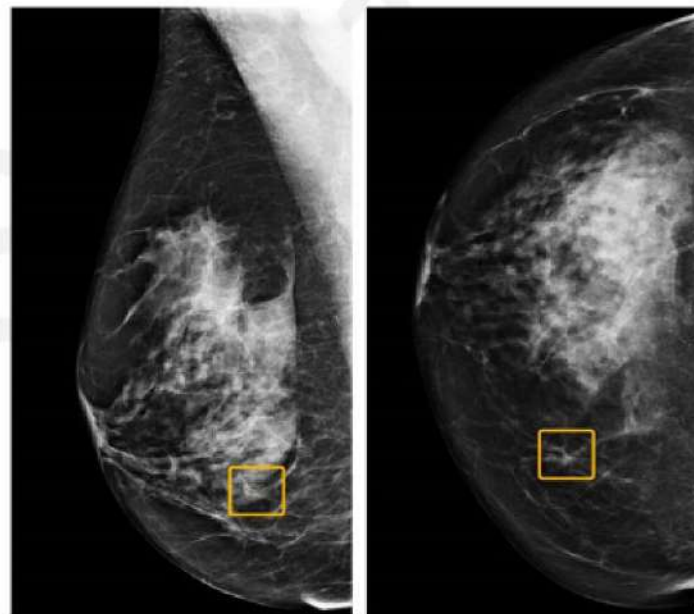


Classification: Breast Cancer Screening

International evaluation of an AI system for breast cancer screening nature



CNN-based system outperformed expert radiologists at detecting breast cancer from mammograms



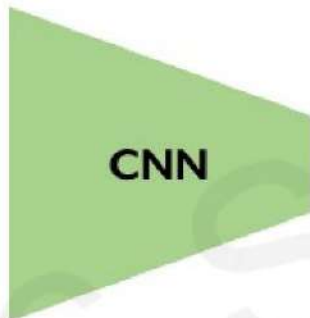
Breast cancer case missed by radiologist but detected by AI



Object Detection



Image x



CNN

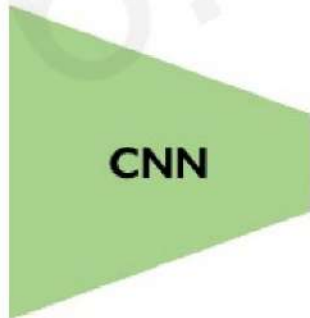


Taxi

Class label y



Image x



CNN

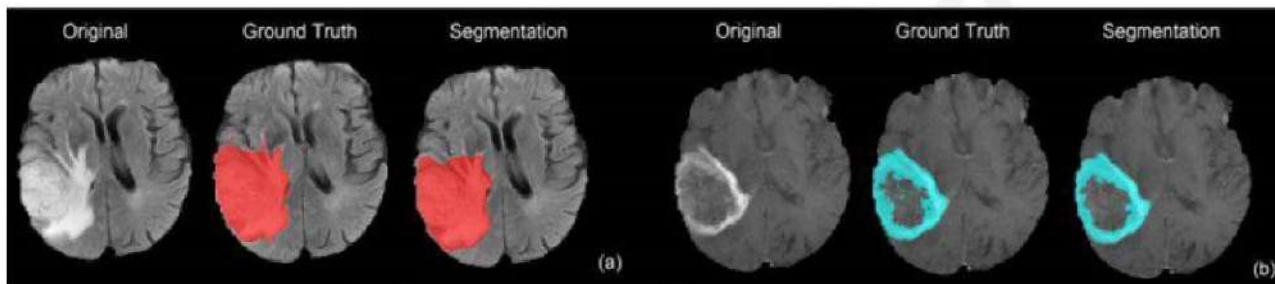


Label (x, y, w, h)

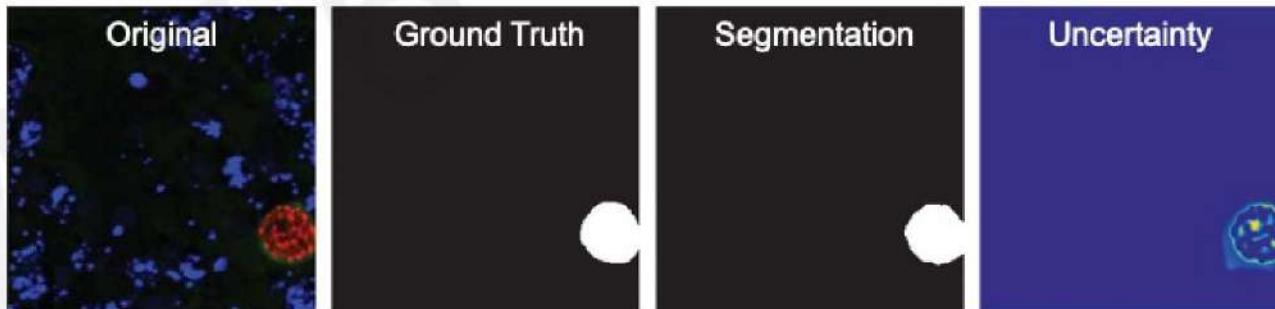


Semantic Segmentation: Biomedical Image Analysis

Brain Tumors
Dong+ *MIUA* 2017.



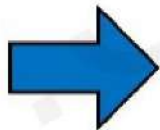
Malaria Infection
Soleimany+ *arXiv* 2019.



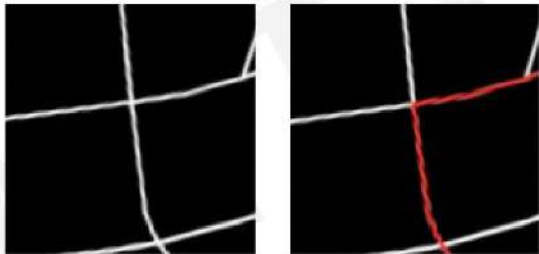


Continuous Control: Navigation from Vision

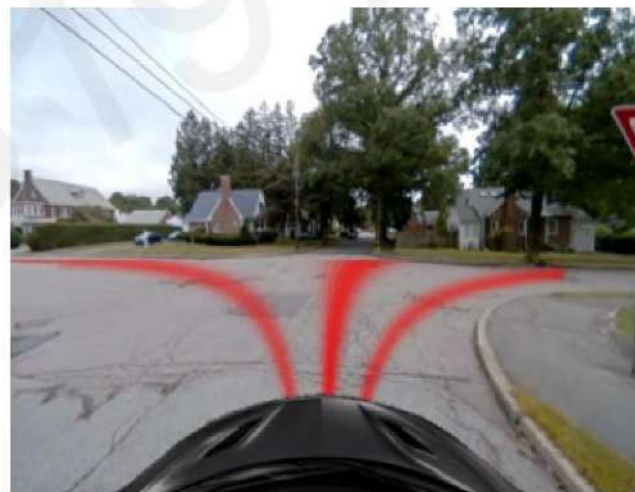
Raw Perception
 I
(ex. camera)



Coarse Maps
 M
(ex. GPS)



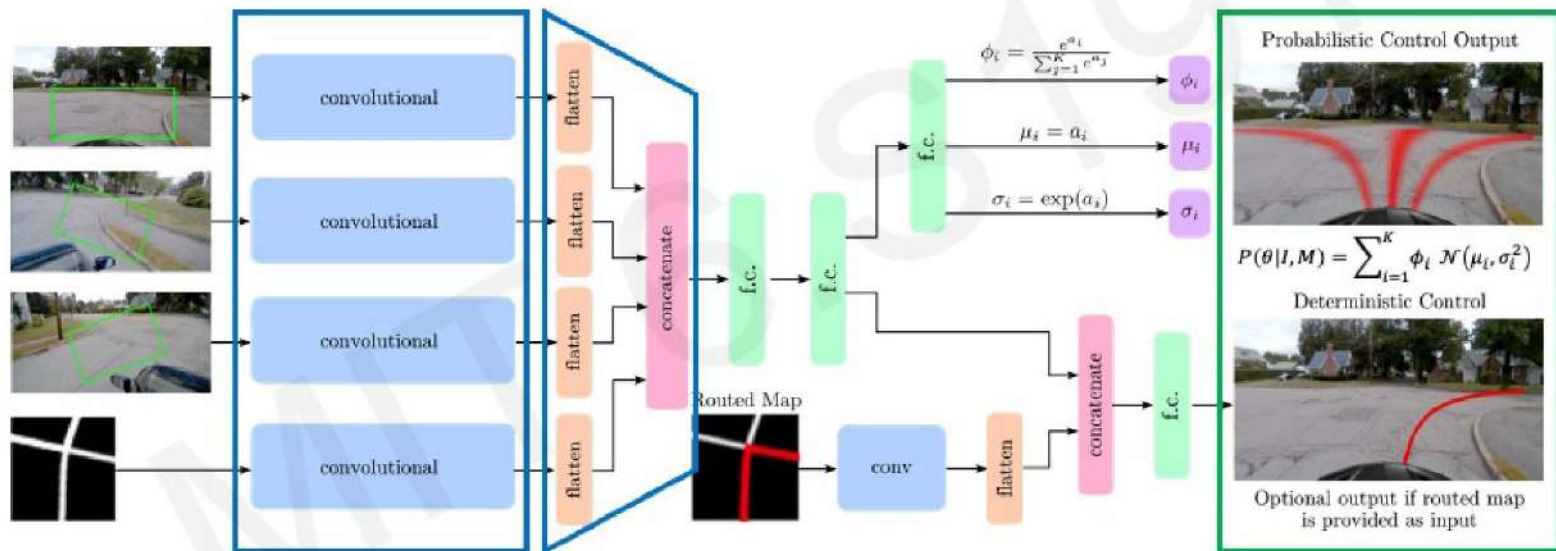
Possible Control Commands





End-to-End Framework for Autonomous Navigation

Entire model is trained end-to-end **without any human labelling or annotations**



$$L = -\log(P(\theta|I, M))$$



人工智能基本理论

ImageNet: 深度学习热潮的关键推动者之一

ImageNet 可以说是计算机视觉研究人员进行大规模物体识别和检测时, 最先想到的视觉大数据来源。ImageNet 数据集最初由**斯坦福大学李飞飞**等人在 CVPR 2009 的一篇论文中推出。

ImageNet Dataset

IMAGENET



Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., ... & Fei-Fei, L. (2015). [Imagenet large scale visual recognition challenge](#). *arXiv preprint arXiv:1409.0575*. [\[web\]](#)



人工智能基本理论

截至 2016 年，ImageNet 中含有超过 1500 万由人手工注释的图片网址，也就是带标签的图片，标签说明了图片中的内容，超过 2.2 万个类别。

English foxhound

An English breed slightly larger than the American foxhounds originally used to hunt in packs

454
pictures

37.57%
Popularity
Percentile

Wordnet
IDs

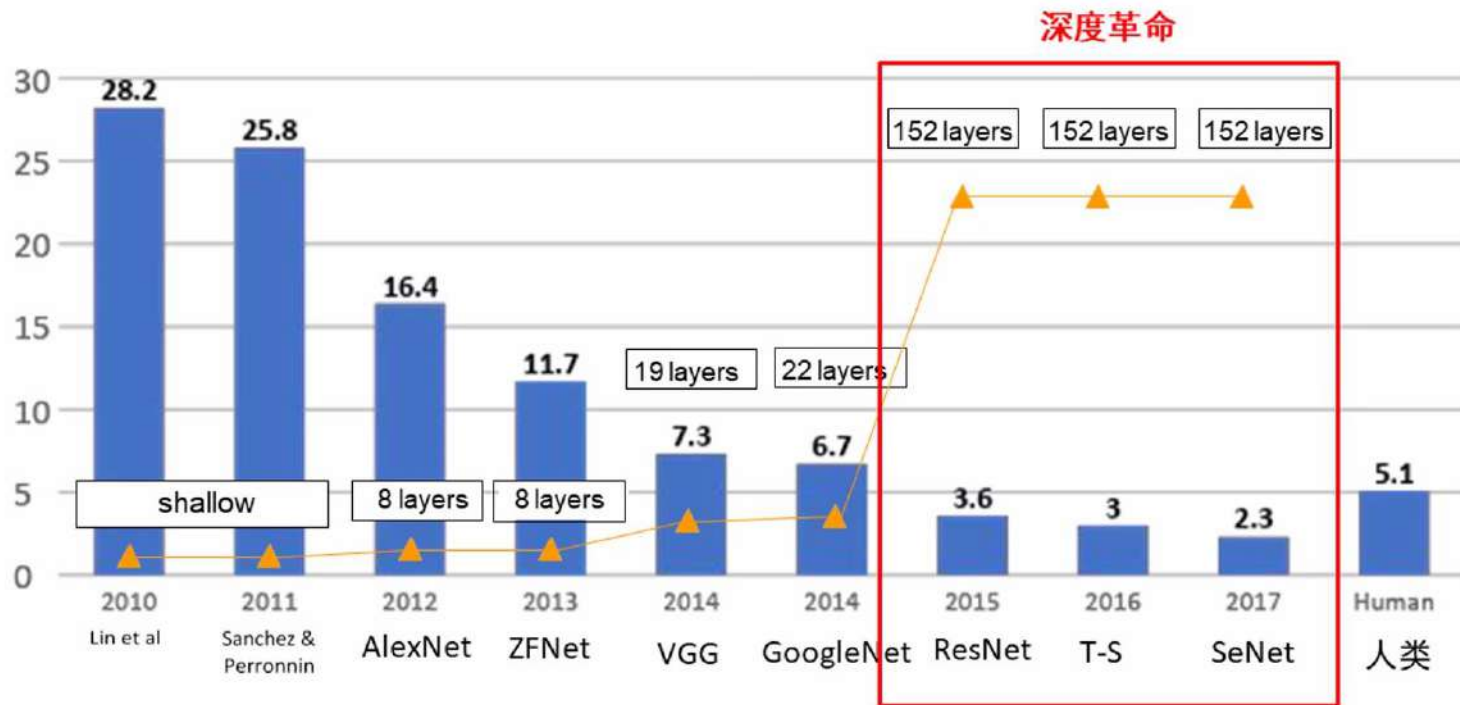
English foxhound

- range animal (0)
- creepy-crawly (0)
- domestic animal, domesticated animal (213)
- domestic cat, house cat, Felis domesticus, Felis catus (18)
- dog, domestic dog, Canis familiaris (189)
 - pooch, doggie, doggy, barker, bow-wow (0)
 - hunting dog (101)
 - sporting dog, gun dog (28)
 - dachshund, dachsie, tsadger dog (1)
 - terrier (37)
 - courser (0)
 - hound, hound dog (29)
 - Plott hound (0)
 - wolfhound (2)
 - Scottish deerhound, deerhound (0)
 - coonhound (2)
 - foxhound (3)
 - Walker hound, Walker foxhound (0)
 - American foxhound (0)
 - English foxhound (0)
 - Weimaraner (0)
 - otterhound, otter hound (0)
 - bloodhound, sleuthhound (0)
 - Norwegian elkhound, elkhound (0)
 - Seluki, gazelle hound (0)
 - Afghan hound, Afghan (0)
 - staghound (0)
 - grayhound (2)
 - beagle (0)
 - harrier (0)
 - basset, basset hound (0)
 - bluetick (0)
 - redbone (0)





ImageNet大赛历年冠军





人工智能基本理论

经典的卷积网络

AlexNet:

使用relu激活函数，提升训练速度；
使用Dropout，缓解过拟合

InceptionNet:

一层内使用不同尺寸卷积核，提升感知力
使用批标准化，缓解梯度消失



LeNet:

卷积网络开篇之作，
共享卷积核，减少网络参数

VGGNet:

小尺寸卷积核减少参数，网络结构规整，适合并行加速

ResNet:

层间残差跳连，引入前方信息，缓解模型退化，使神经网络层数加深成为可能



人工智能基本理论

LeNet

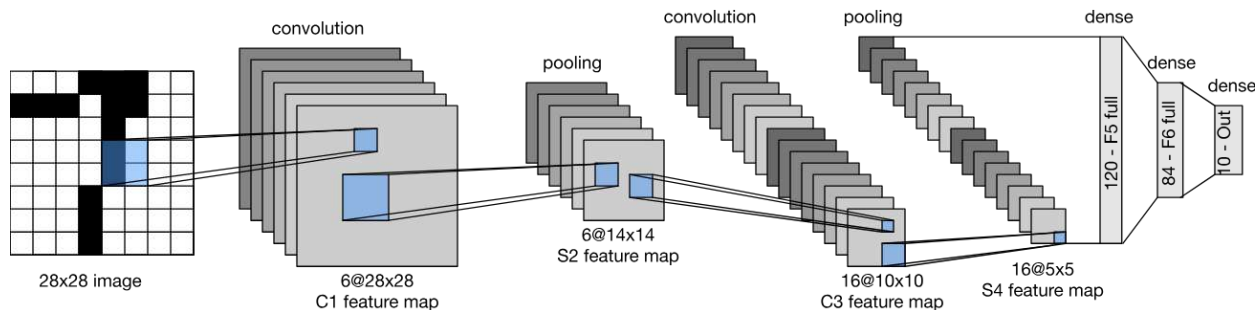
LeNet是最早的CNN架构之一,由Yann LeCun等人于1998年提出,用于手写数字识别。它的结构如下:

INPUT => CONV => POOL => CONV => POOL => FC => FC => OUTPUT

其中:

- CONV: 卷积层,提取局部特征
- POOL: 池化层,降低特征维度
- FC: 全连接层,进行分类

LeNet的创新之处在于引入了卷积层和池化层,极大地减少了网络参数,提高了模型的泛化能力。





人工智能基本理论

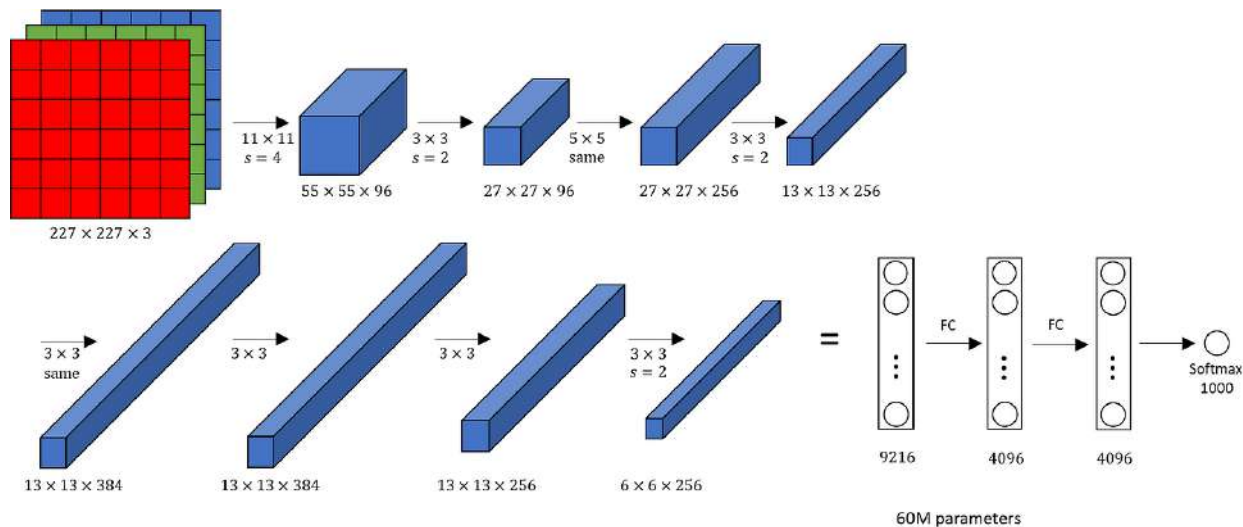
AlexNet

AlexNet是2012年ImageNet大赛的冠军模型,由Alex Krizhevsky等人提出,被认为是深度学习在计算机视觉领域的开端。它的结构如下:

INPUT => CONV => POOL => CONV => POOL => CONV => CONV => CONV => POOL => FC => FC => FC => OUTPUT

AlexNet的主要创新点包括:

- 1.使用了ReLU激活函数,提高了训练效率。
- 2.引入了Dropout技术,有效缓解了过拟合问题。
- 3.利用GPU并行计算,大大加快了训练速度。
- 4.数据增强技术,如翻转、裁剪等,增加了训练数据的多样性。





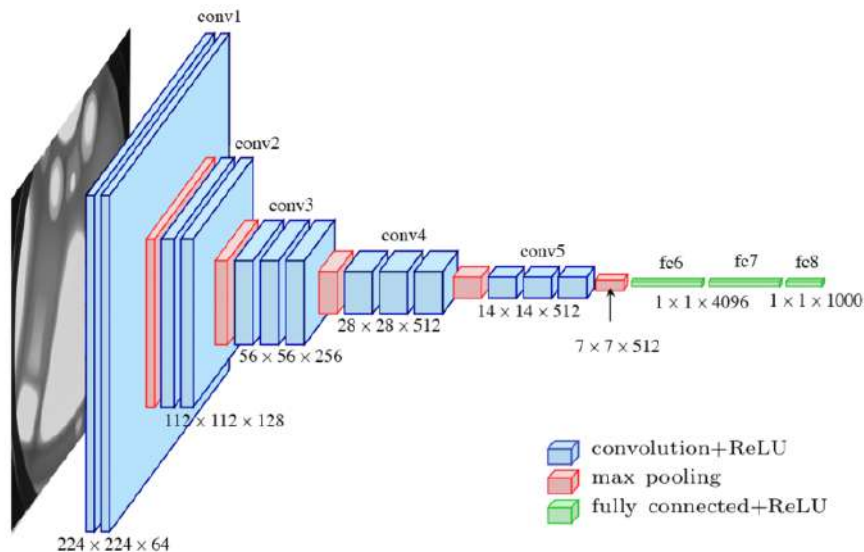
人工智能基本理论

VGGNet

VGGNet是2014年ImageNet大赛的亚军模型,由Karen Simonyan和Andrew Zisserman提出。它的结构如下:

INPUT => CONV => CONV => POOL => CONV => CONV => POOL => CONV => CONV => CONV => POOL => CONV => CONV => CONV => POOL => CONV => CONV => CONV => POOL => FC => FC => FC => OUTPUT

VGGNet的特点是使用了连续的 3×3 小卷积核,代替了 AlexNet 中的大卷积核,提高了非线性表达能力。同时, VGGNet 也探索了不同深度的网络结构,发现深度对提高性能至关重要。





人工智能基本理论

VGGNet

```
INPUT: [224x224x3]      memory: 224*224*3=150K  weights: 0
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M  weights: (3*3*3)*64 = 1,728
CONV3-64: [224x224x64]  memory: 224*224*64=3.2M  weights: (3*3*64)*64 = 36,864
POOL2: [112x112x64]    memory: 112*112*64=800K  weights: 0
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  weights: (3*3*64)*128 = 73,728
CONV3-128: [112x112x128] memory: 112*112*128=1.6M  weights: (3*3*128)*128 = 147,456
POOL2: [56x56x128]     memory: 56*56*128=400K  weights: 0
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*128)*256 = 294,912
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*256)*256 = 589,824
CONV3-256: [56x56x256] memory: 56*56*256=800K  weights: (3*3*256)*256 = 589,824
POOL2: [28x28x256]     memory: 28*28*256=200K  weights: 0
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*256)*512 = 1,179,648
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [28x28x512] memory: 28*28*512=400K  weights: (3*3*512)*512 = 2,359,296
POOL2: [14x14x512]     memory: 14*14*512=100K  weights: 0
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
CONV3-512: [14x14x512] memory: 14*14*512=100K  weights: (3*3*512)*512 = 2,359,296
POOL2: [7x7x512]       memory: 7*7*512=25K    weights: 0
FC: [1x1x4096]         memory: 4096           weights: 7*7*512*4096 = 102,760,448
FC: [1x1x4096]         memory: 4096           weights: 4096*4096 = 16,777,216
FC: [1x1x1000]         memory: 1000           weights: 4096*1000 = 4,096,000
```

TOTAL memory: 24M * 4 bytes ≈ 93MB / image (only forward! ≈*2 for bwd)

TOTAL params: 138M parameters

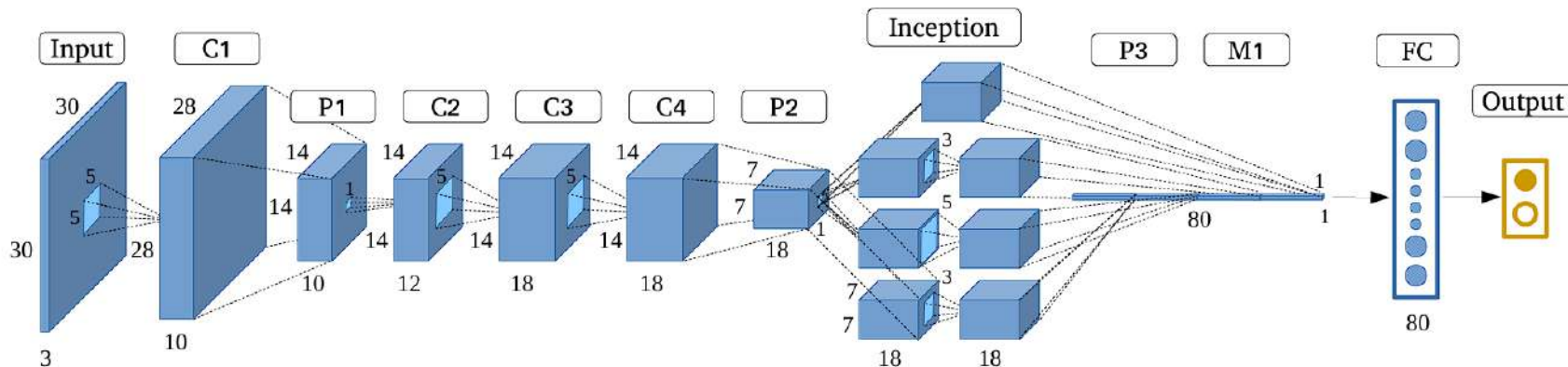
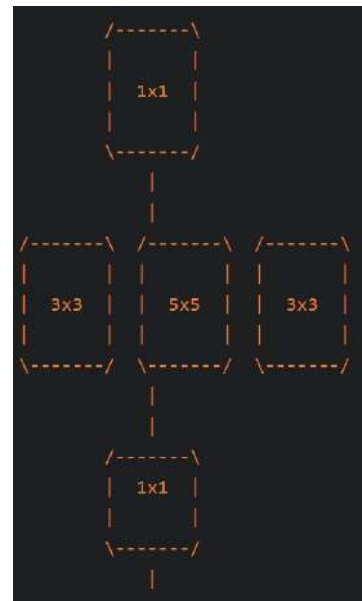


人工智能基本理论

GoogLeNet(Inception)

GoogLeNet是2014年ImageNet大赛的冠军模型,由Christian Szegedy等人提出。它的核心是Inception模块,如下所示:

Inception模块通过并行的卷积核组合,提高了网络的表达能力,同时控制了参数数量。GoogLeNet还引入了辅助分类器,缓解了梯度消失问题。



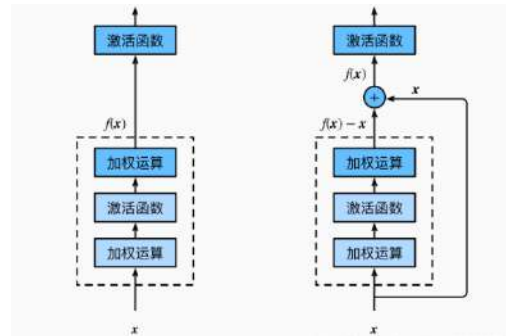


人工智能基本理论

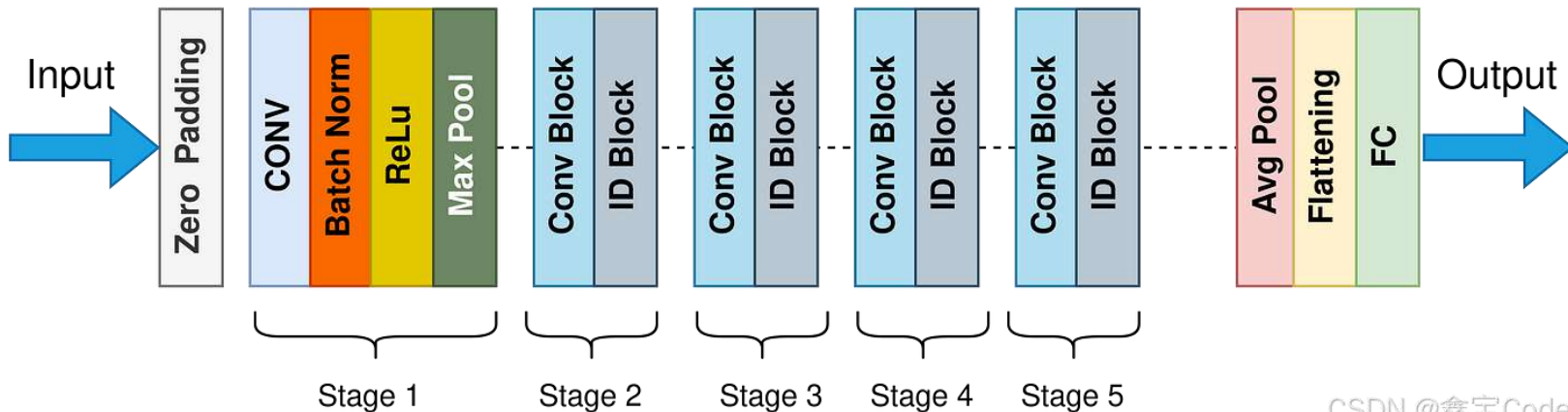
ResNet

ResNet是2015年ImageNet大赛的冠军模型,由Kaiming He等人提出。它的核心是残差模块(Residual Block),如下所示:

残差模块通过引入shortcut connection,使得梯度可以直接传递到较浅层,有效缓解了梯度消失/爆炸问题,从而可以训练出更深的网络。ResNet在ImageNet数据集上取得了极高的分类精度。

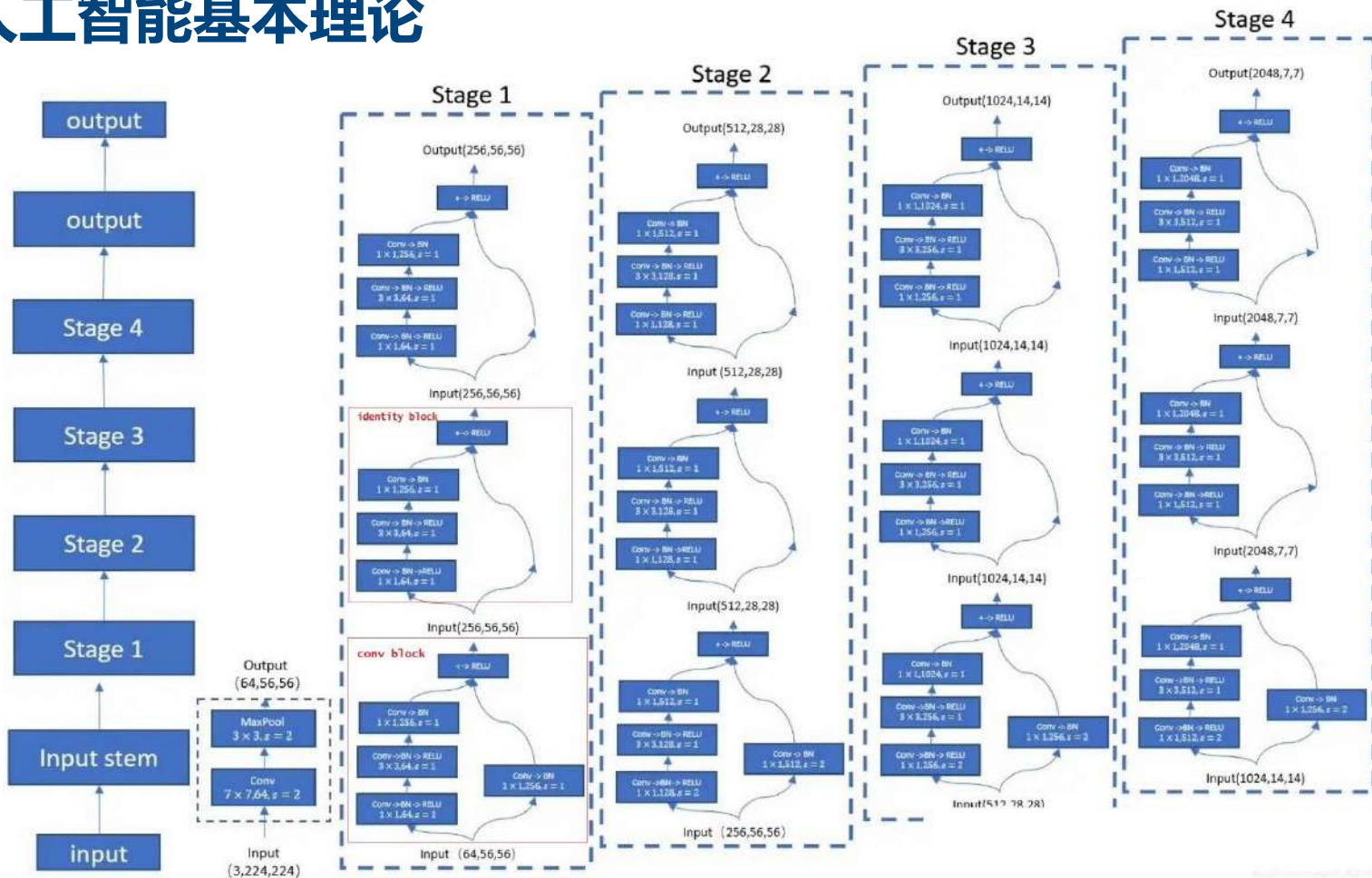


ResNet50 Model Architecture





人工智能基本理论





Deep Learning for Computer Vision: Summary

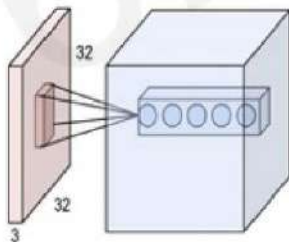
Foundations

- Why computer vision?
- Representing images
- Convolutions for feature extraction



CNNs

- CNN architecture
- Application to classification
- ImageNet



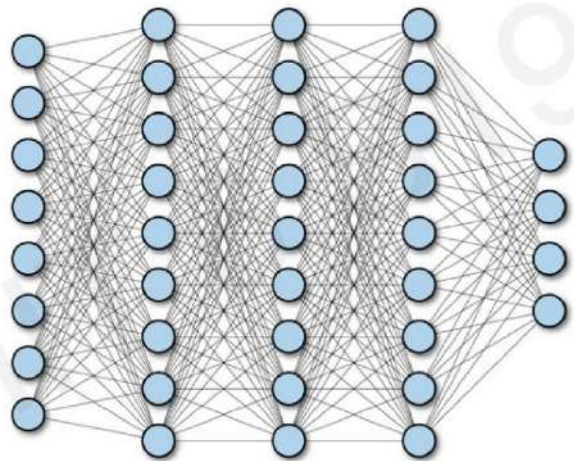
Applications

- Segmentation, image captioning, control
- Security, medicine, robotics

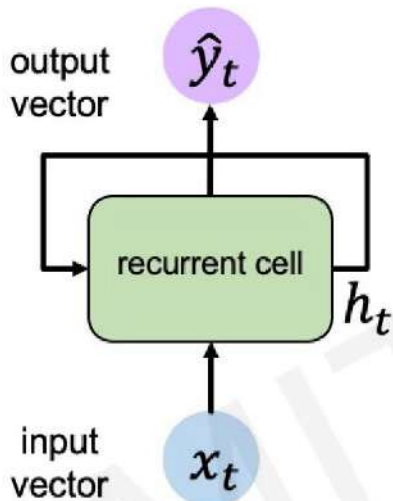




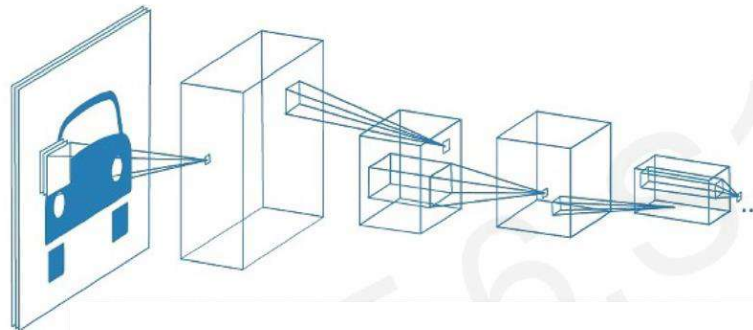
人工智能基本理论



前馈神经网络
MLP



循环神经网络
RNN



卷积神经网络
CNN



谢谢大家