# 資料探勘
# (Data Mining)

# 遞歸神經網絡
# (Recurrent Neural Networks)

1092DM10
MBA, IM, NTPU (M5026) (Spring 2021)
Tue 2, 3, 4 (9:10-12:00) (B8F40)

**Min-Yuh Day**
戴敏育
**Associate Professor**
副教授
**Institute of Information Management**, **National Taipei University**
國立臺北大學 資訊管理研究所

https://web.ntpu.edu.tw/~myday

2021-05-25

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

1  2021/02/23  資料探勘介紹 (Introduction to data mining)

2  2021/03/02  ABC：人工智慧，大數據，雲端運算
(ABC: AI, Big Data, Cloud Computing)

3  2021/03/09  Python資料探勘的基礎
(Foundations of Data Mining in Python)

4  2021/03/16  資料科學與資料探勘：發現，分析，可視化和呈現數據
(Data Science and Data Mining:
Discovering, Analyzing, Visualizing and Presenting Data)

5  2021/03/23  非監督學習：關聯分析，購物籃分析
(Unsupervised Learning: Association Analysis,
Market Basket Analysis)

6  2021/03/30  資料探勘個案研究 I
(Case Study on Data Mining I)

# 課程大綱 (Syllabus)

週次 (Week)　日期 (Date)　內容 (Subject/Topics)

7  2021/04/06 放假一天 (Day off)

8  2021/04/13 非監督學習：集群分析，行銷市場區隔
(Unsupervised Learning: Cluster Analysis, Market Segmentation)

9  2021/04/20 期中報告 (Midterm Project Report)

10  2021/04/27 監督學習：分類和預測
(Supervised Learning: Classification and Prediction)

11  2021/05/04 機器學習和深度學習
(Machine Learning and Deep Learning)

12  2021/05/11 卷積神經網絡
(Convolutional Neural Networks)

# 課程大綱 (Syllabus)

週次 (Week)    日期 (Date)    內容 (Subject/Topics)

13  2021/05/18  資料探勘個案研究 II
(Case Study on Data Mining II)

14  2021/05/25  遞歸神經網絡
(Recurrent Neural Networks)

15  2021/06/01  強化學習
(Reinforcement Learning)

16  2021/06/08  社交網絡分析
(Social Network Analysis)

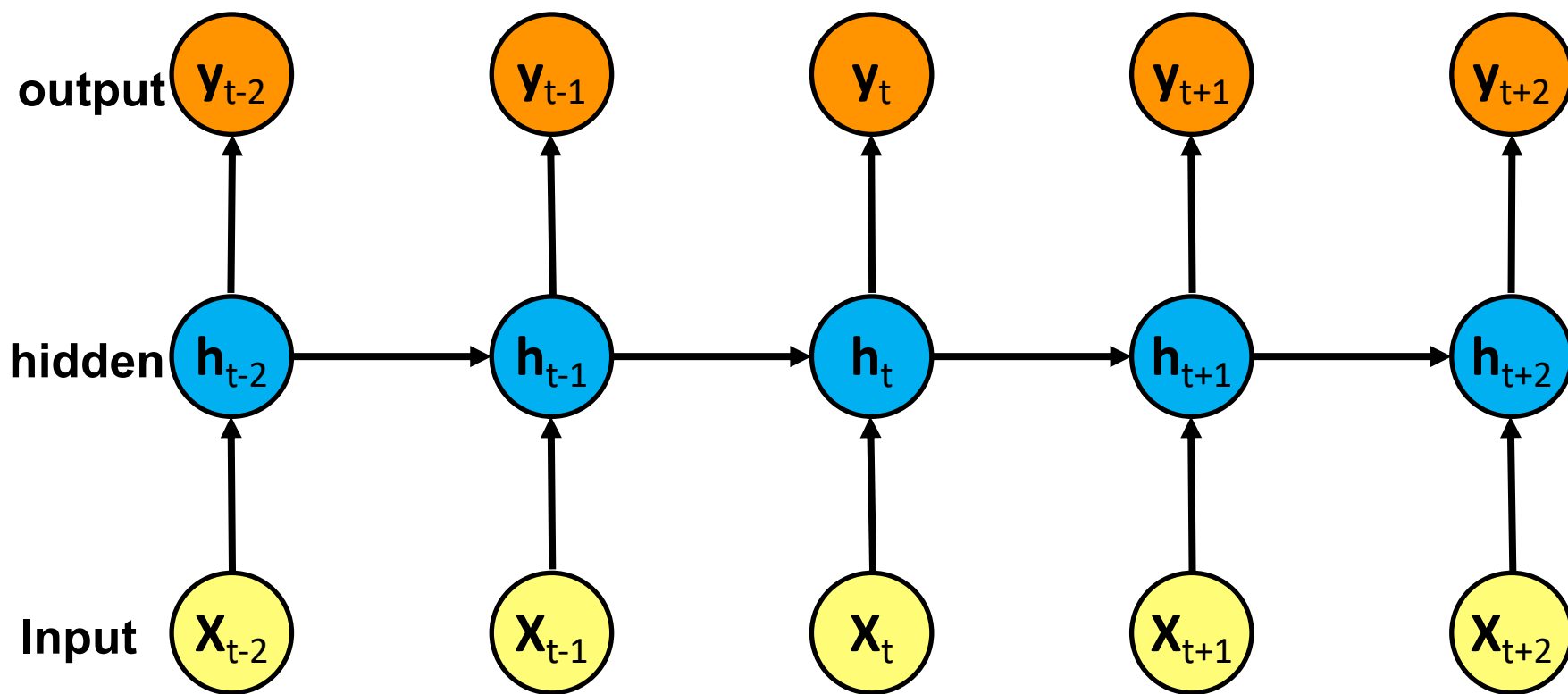17  2021/06/15  期末報告 I (Final Project Report I)

18  2021/06/22  期末報告 II (Final Project Report II)
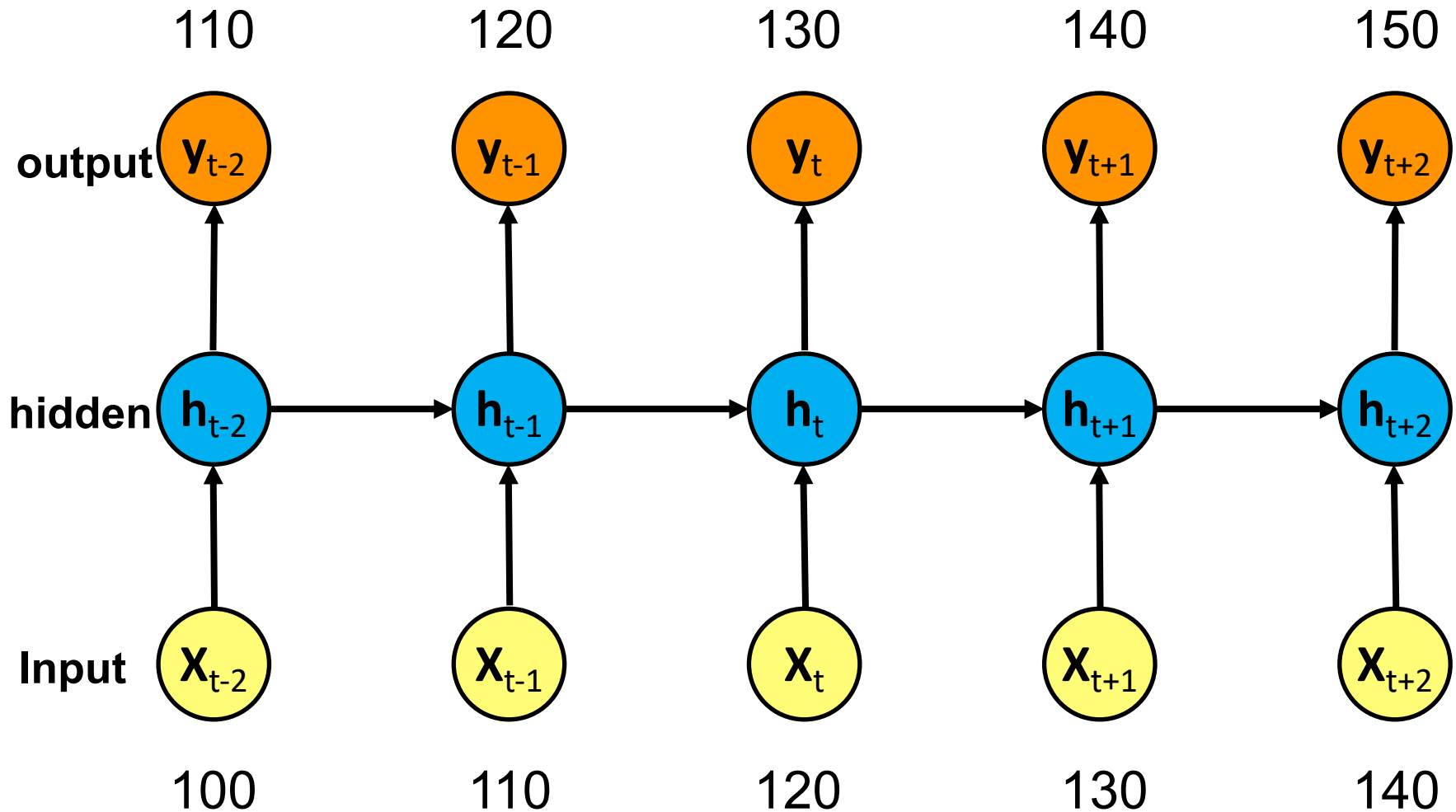
# Recurrent Neural Networks (RNN)

# Outline

- **Recurrent Neural Networks (RNN)**

  - **Long Short Term Memory (LSTM)**

  - **Gated Recurrent Unit (GRU)**

- **Deep Learning (RNN) for Time Series Prediction**

- **Deep Learning (RNN) for Text Analytics (NLP)**
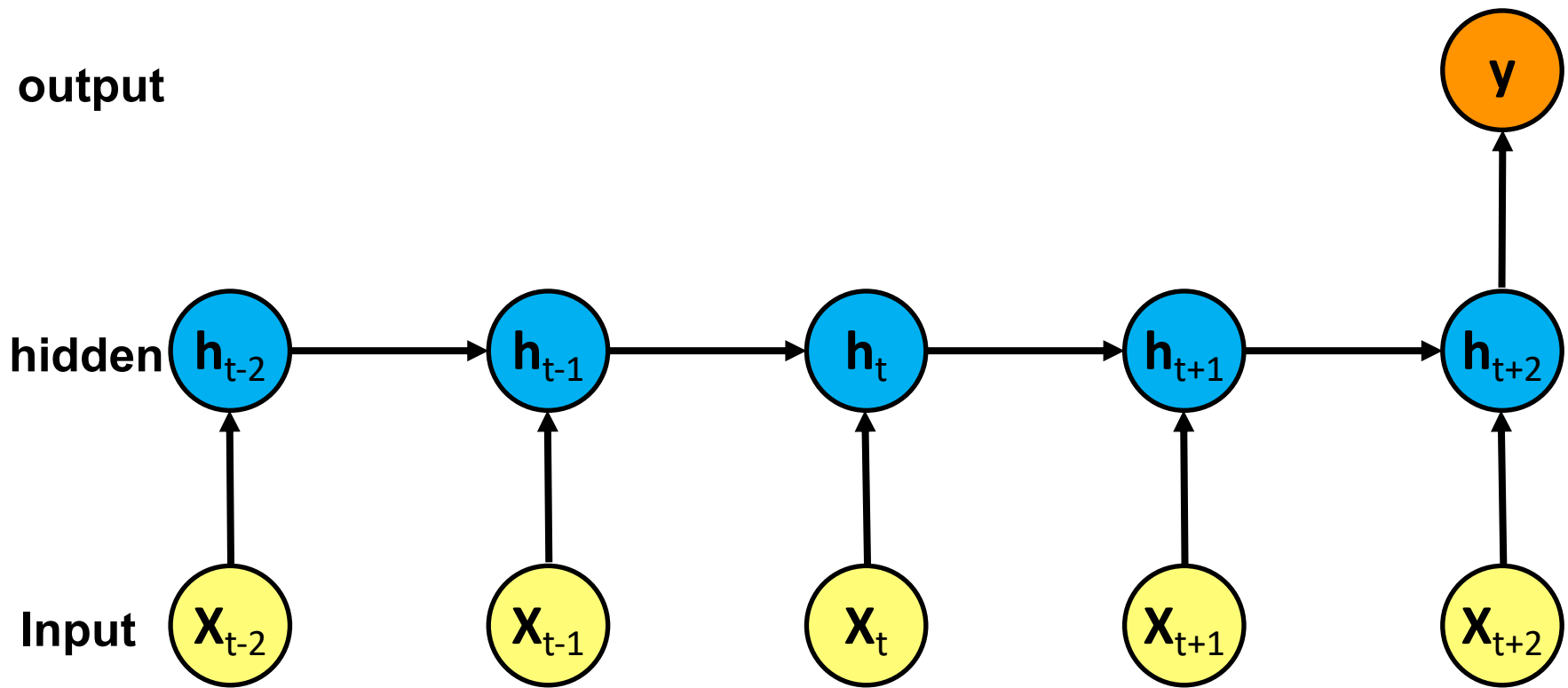
# Recurrent Neural Networks (RNN)

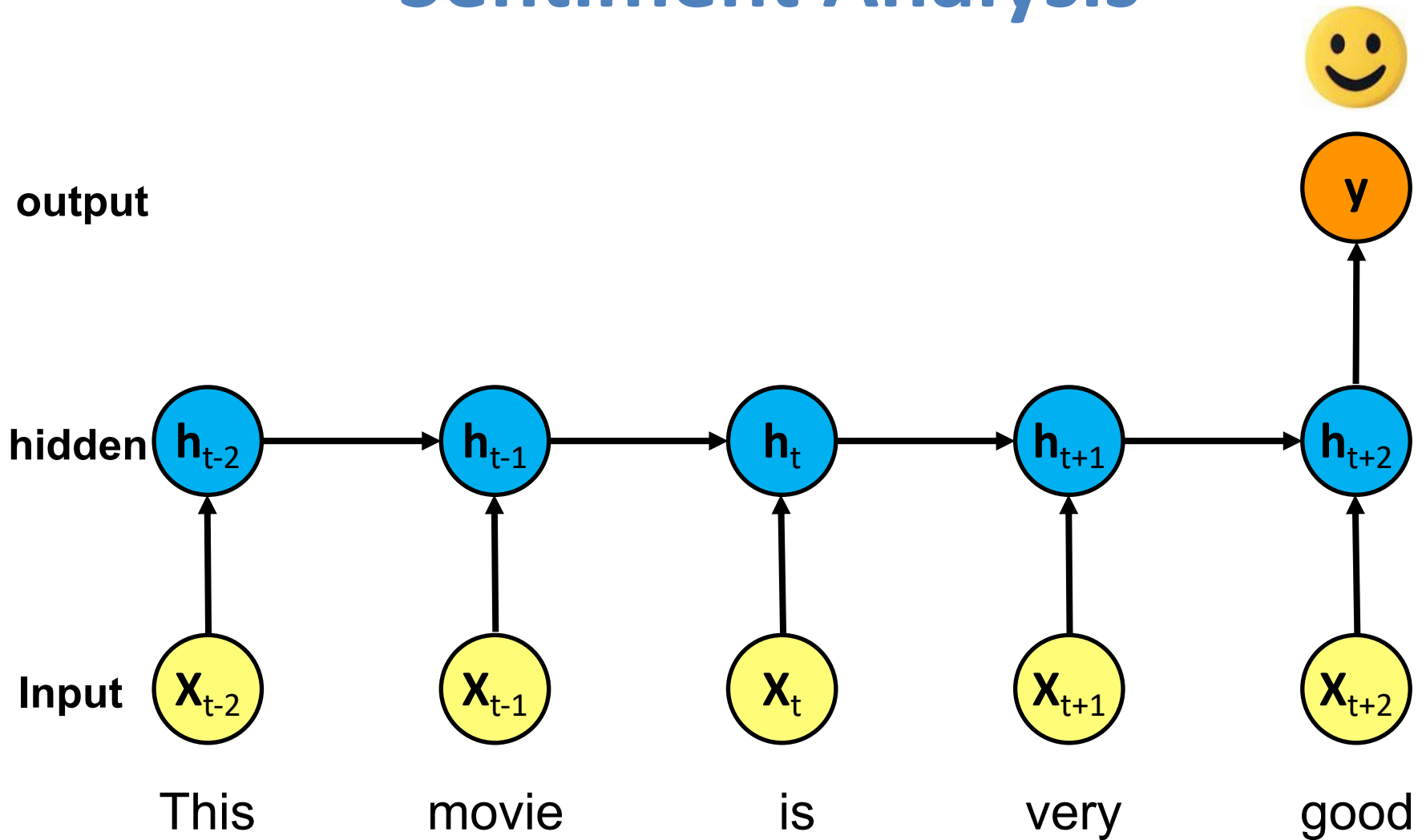# Recurrent Neural Networks (RNN) Time Series Forecasting
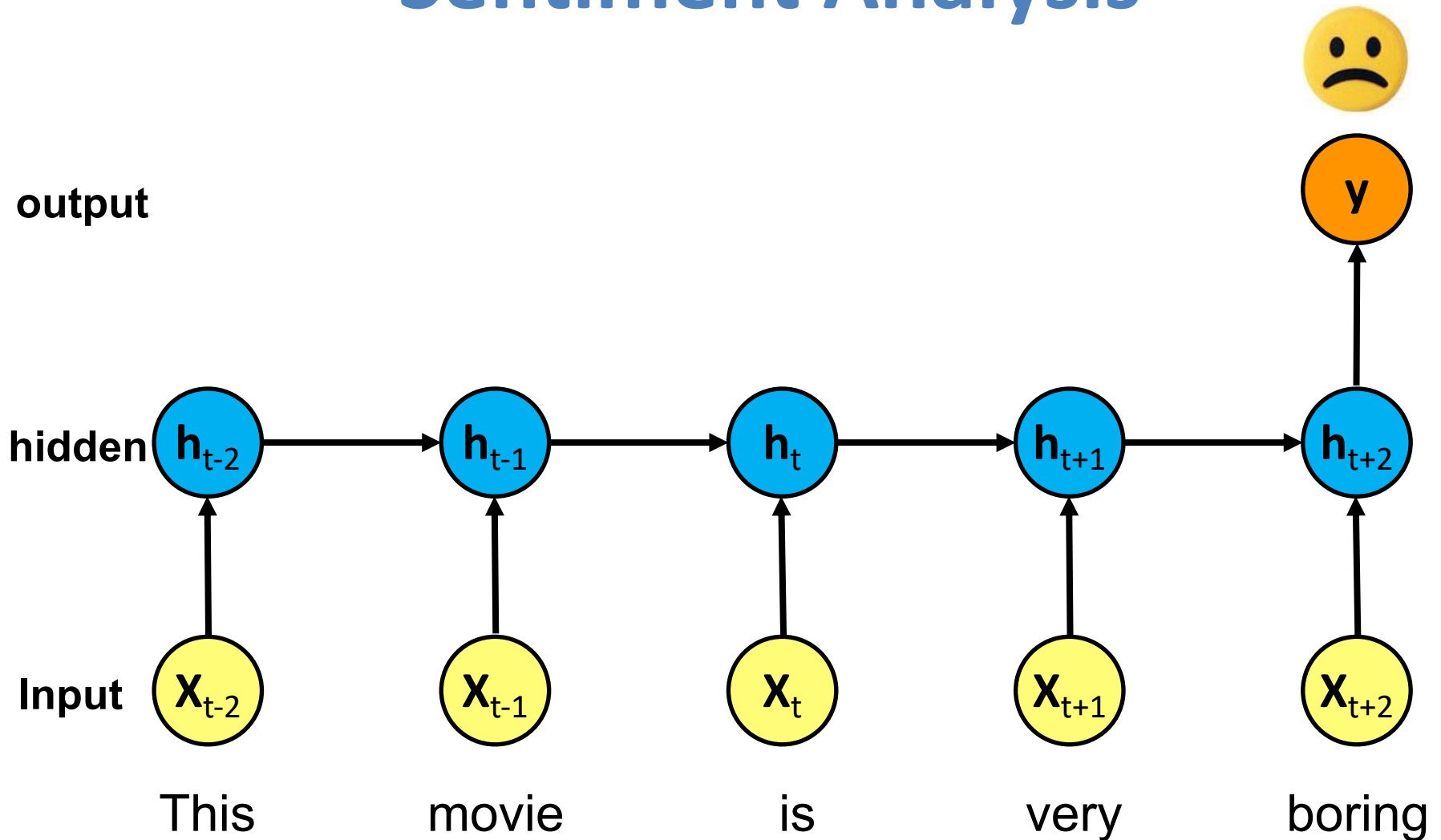
# Recurrent Neural Networks (RNN)
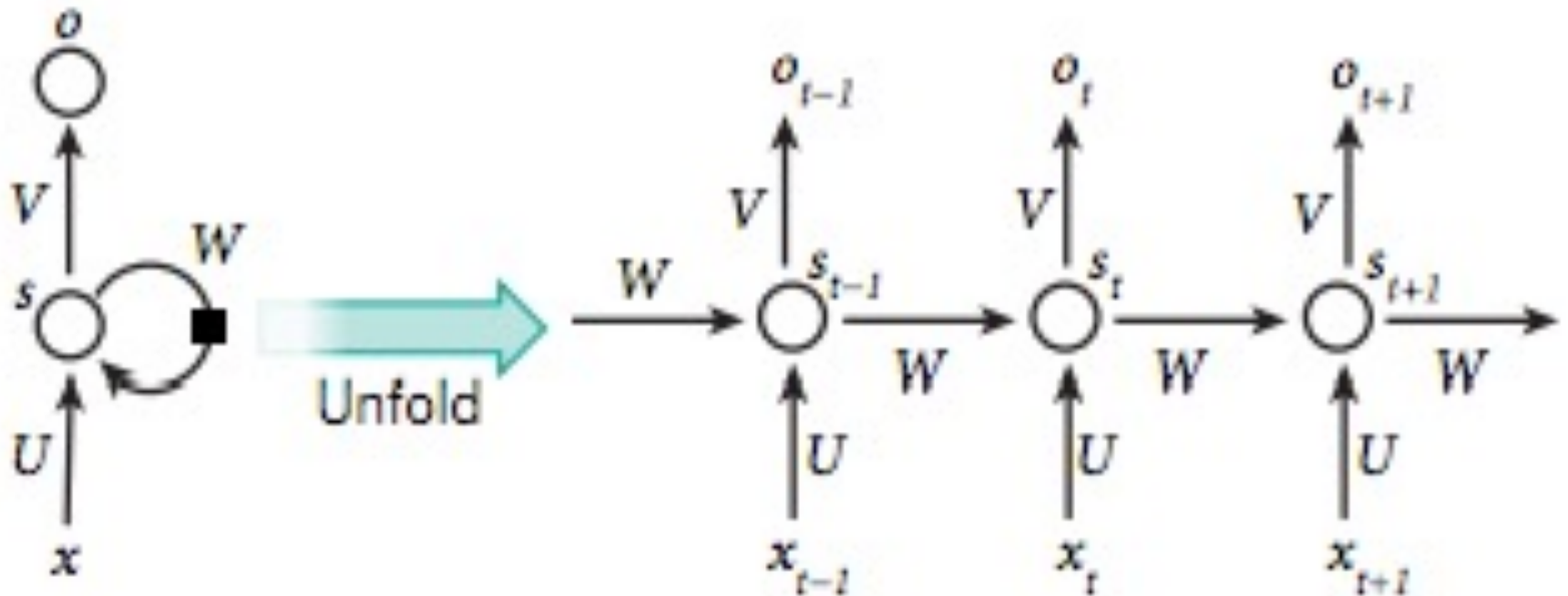
# Recurrent Neural Networks (RNN) Sentiment Analysis

**output**

**hidden**

**Input**

| $h_{t-2}$ | $h_{t-1}$ | $h_t$ | $h_{t+1}$ | $h_{t+2}$ |

| $X_{t-2}$ | $X_{t-1}$ | $X_t$ | $X_{t+1}$ | $X_{t+2}$ |

This       movie       is       very       good

**y**

# Recurrent Neural Networks (RNN) Sentiment Analysis
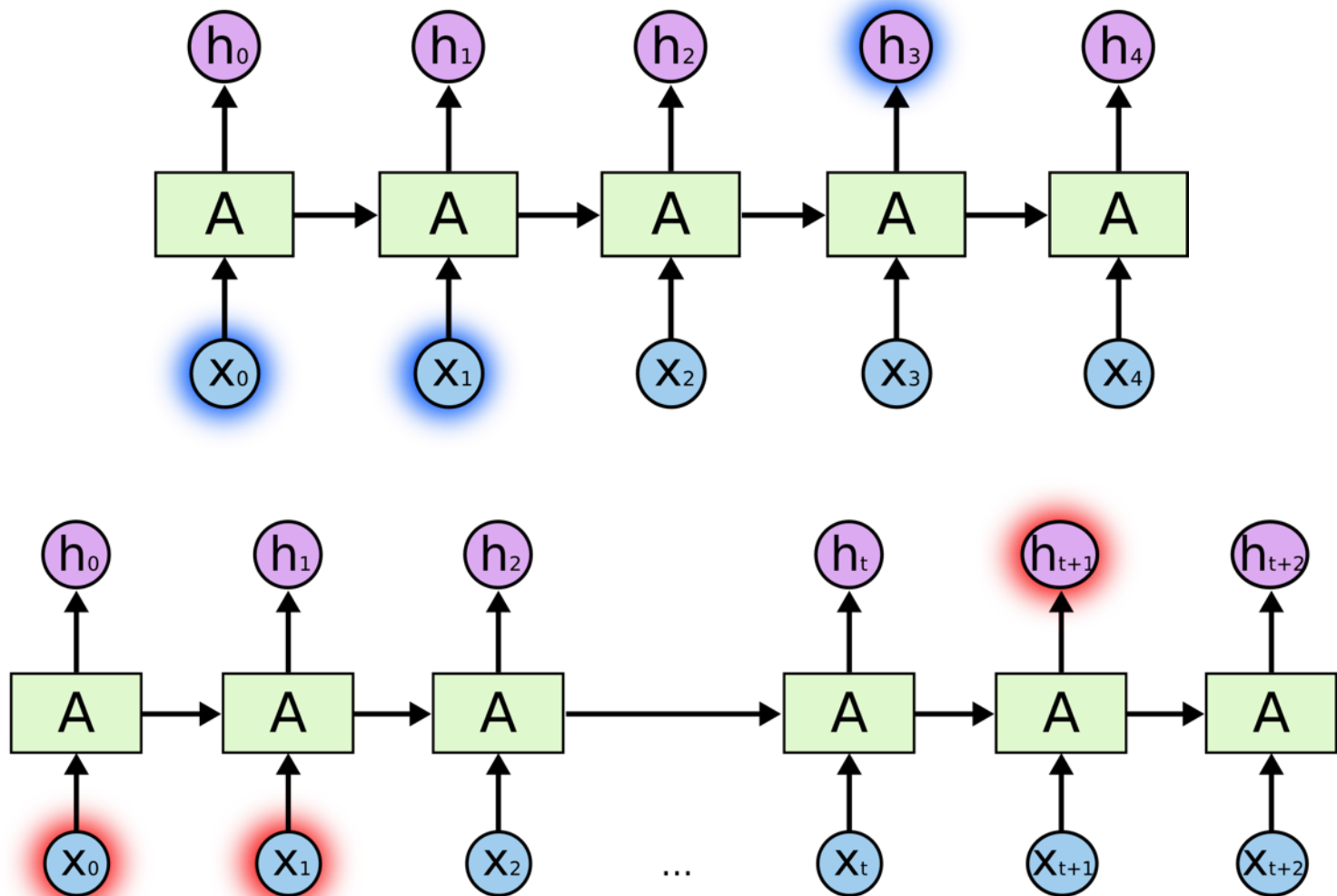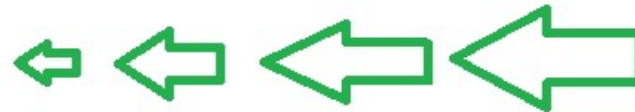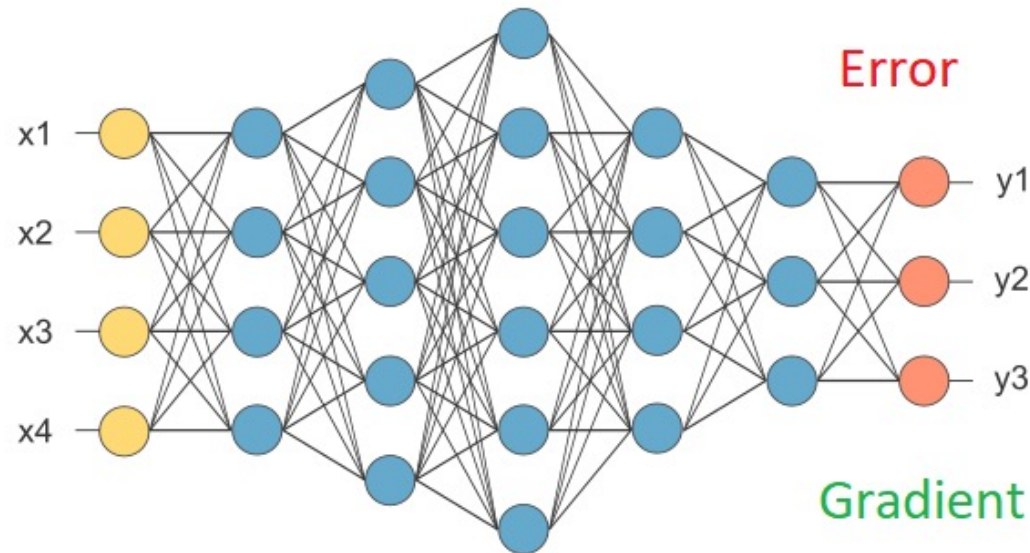
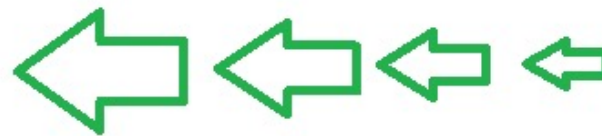# Recurrent Neural Network (RNN)

# RNN

# RNN long-term dependencies



I grew up in France… I speak fluent French.
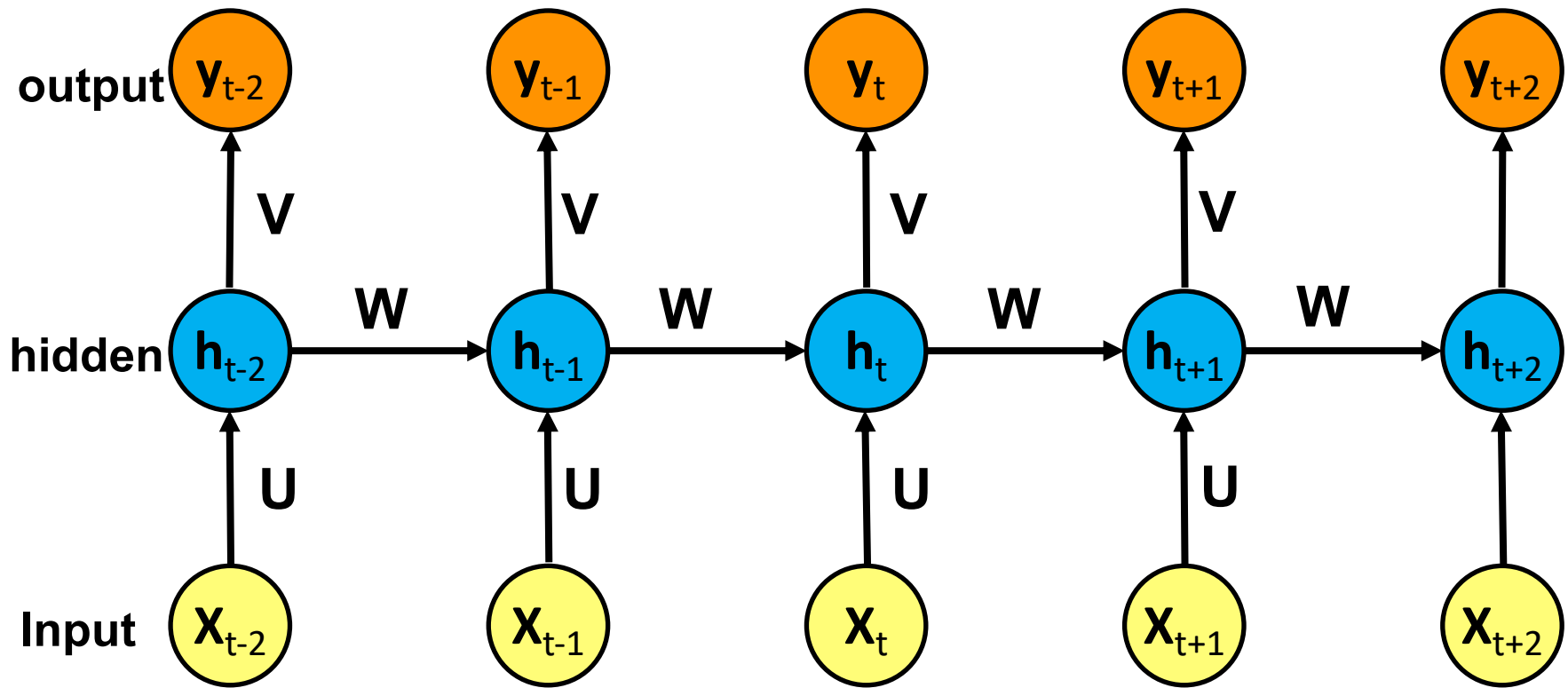
# Vanishing Gradient Exploding Gradient

# Recurrent Neural Networks (RNN)



output   $y_{t-2}$   $y_{t-1}$   $y_t$   $y_{t+1}$   $y_{t+2}$

V   V   V   V

hidden   $h_{t-2}$ —W→ $h_{t-1}$ —W→ $h_t$ —W→ $h_{t+1}$ —W→ $h_{t+2}$

U   U   U   U

Input   $X_{t-2}$   $X_{t-1}$   $X_t$   $X_{t+1}$   $X_{t+2}$

16

# RNN
# Vanishing Gradient problem
# Exploding Gradient problem



**Error**

output: $y_{t-2}$, $y_{t-1}$, $y_t$, $y_{t+1}$, $y_{t+2}$

V

hidden: $h_{t-2}$ —W→ $h_{t-1}$ —W→ $h_t$ —W→ $h_{t+1}$ —W→ $h_{t+2}$

U

Input: $X_{t-2}$, $X_{t-1}$, $X_t$, $X_{t+1}$, $X_{t+2}$

**if |W| < 1 (Vanishing)**
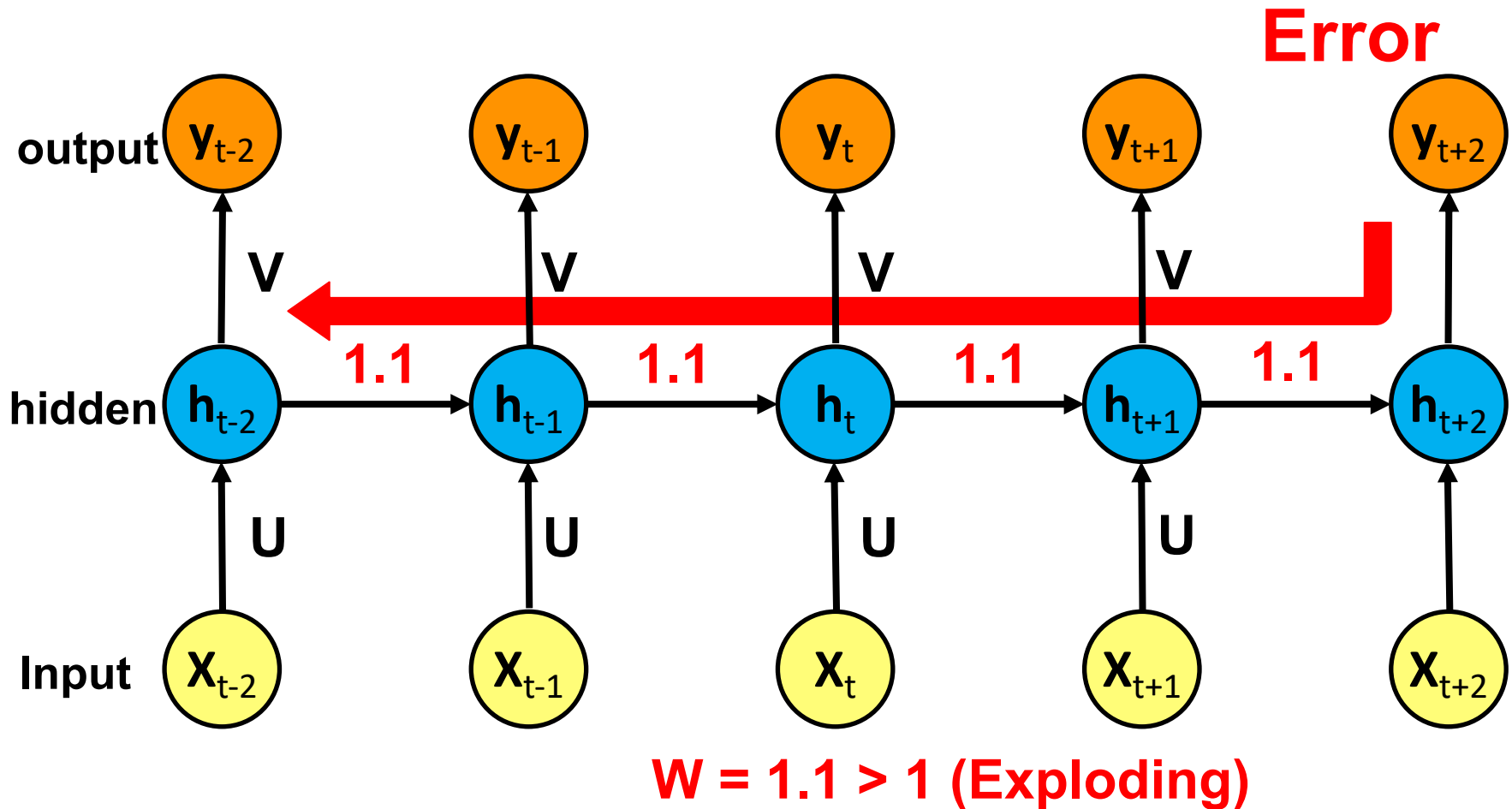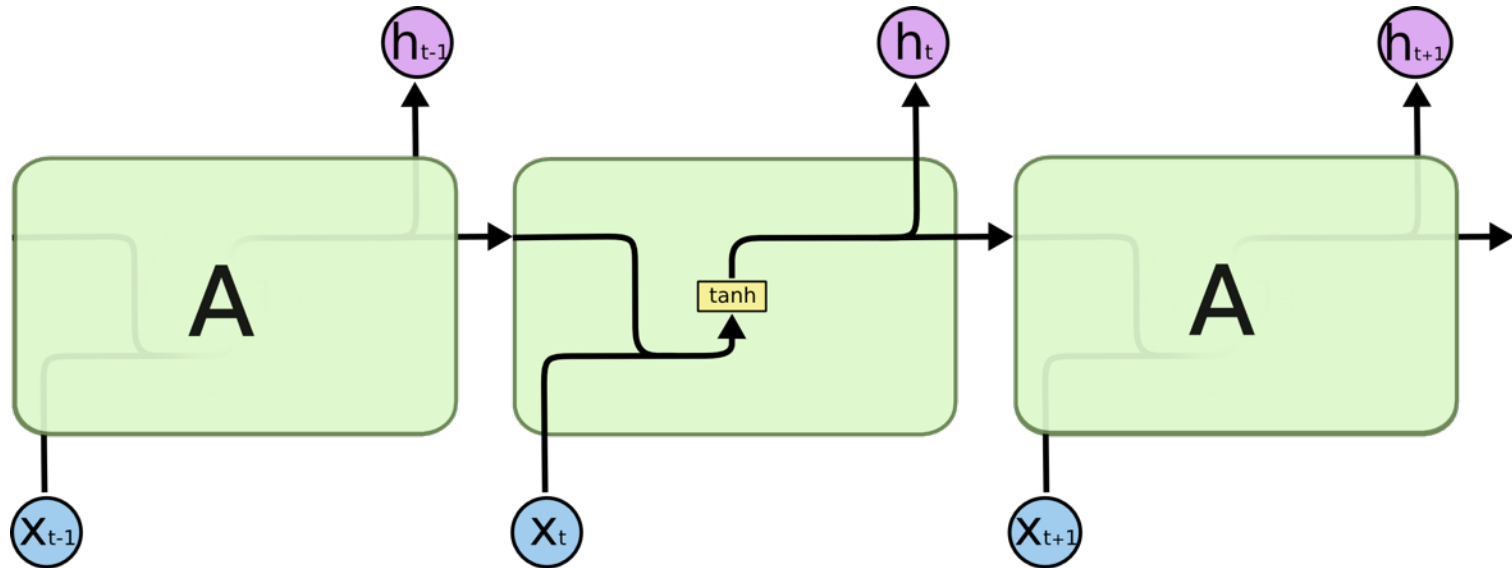**if |W| > 1 (Exploding)**

17

# RNN
# Vanishing Gradient problem
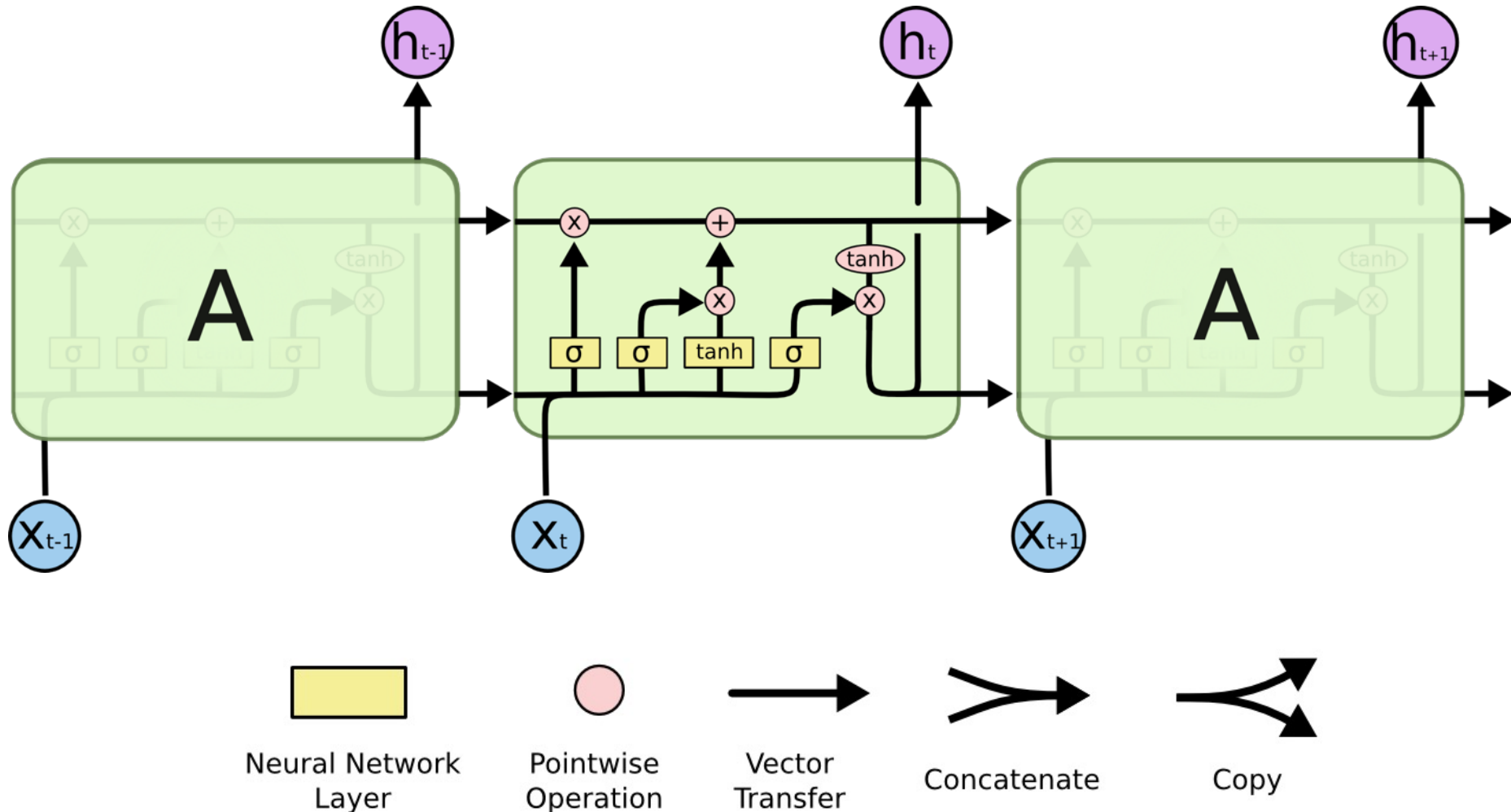
# RNN
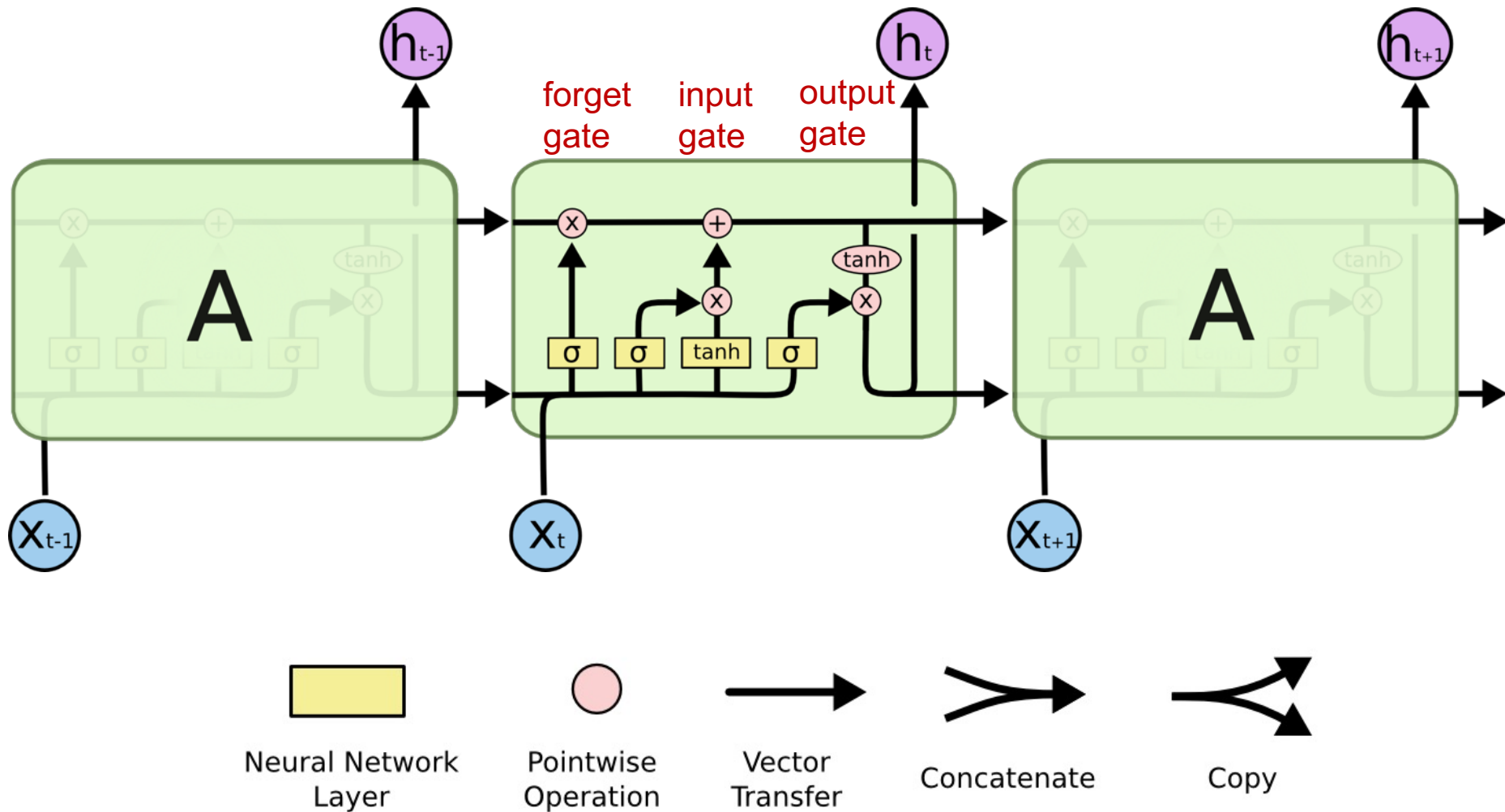# Exploding Gradient problem

# RNN LSTM

**RNN**



**LSTM**

# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)



forget gate    input gate    output gate

Neural Network Layer    Pointwise Operation    Vector Transfer    Concatenate    Copy

# Gated Recurrent Unit (GRU)

# Gated Recurrent Unit (GRU)

# LSTM

# LSTM vs GRU



## LSTM

## GRU

i, f and o are the input, forget and output gates, respectively.
c and c˜ denote the memory cell and the new memory cell content.

r and z are the reset and update gates, and h and h˜ are the activation and the candidate activation.

# Long Short Term Memory (LSTM)

# Long Short Term Memory (LSTM)

# LSTM
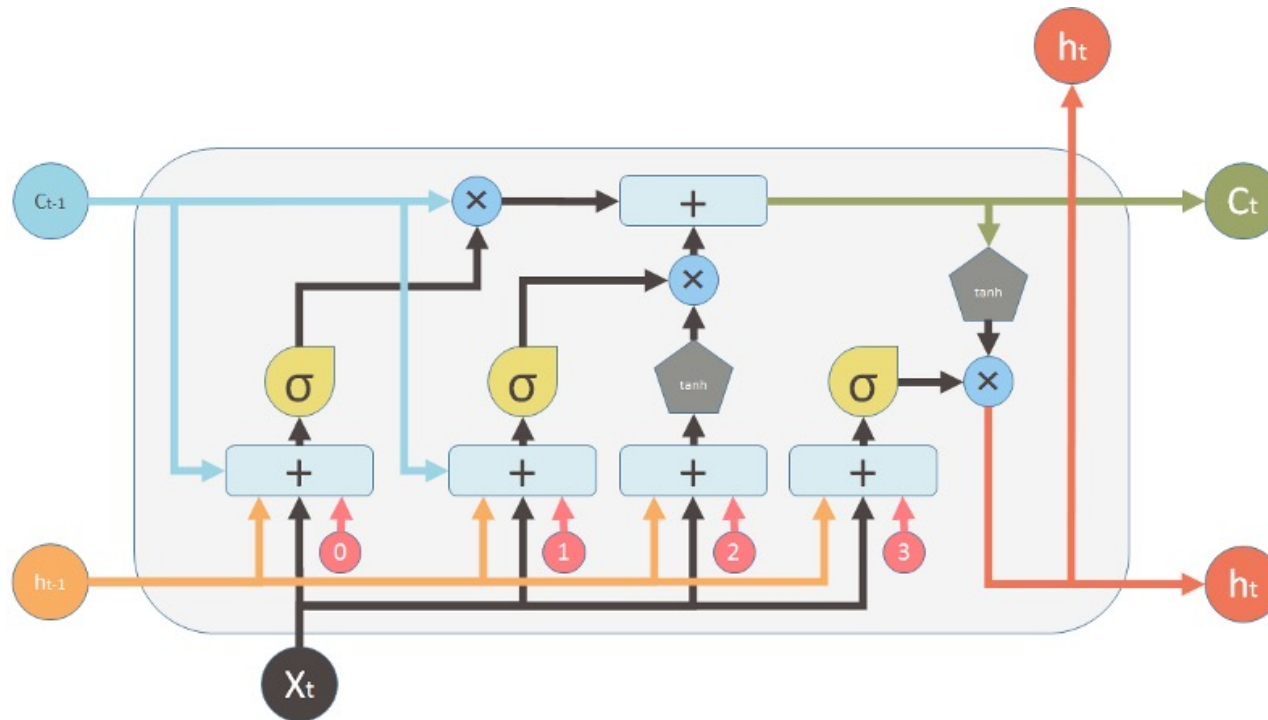# Memory state (C)

# LSTM
# *forget gate (f)*



$$f_t = \sigma \left( W_f \cdot [h_{t-1}, x_t] \; + \; b_f \right)$$

# LSTM
# input gate (i)



$$i_t = \sigma\left(W_i \cdot [h_{t-1}, x_t] \; + \; b_i\right)$$

$$\tilde{C}_t = \tanh(W_C \cdot [h_{t-1}, x_t] \; + \; b_C)$$

# LSTM
# Memory state (C)



$$C_t = f_t * C_{t-1} + i_t * \tilde{C}_t$$

# *LSTM output gate (o)*



$$o_t = \sigma \left( W_o \left[ h_{t-1}, x_t \right] + b_o \right)$$

$$h_t = o_t * \tanh \left( C_t \right)$$

# LSTM
# forget (f), input (i), output (o) gates



$$f_t = \sigma \left( W_f \cdot [C_{t-1}, h_{t-1}, x_t] \; + \; b_f \right)$$

$$i_t = \sigma \left( W_i \cdot [C_{t-1}, h_{t-1}, x_t] \; + \; b_i \right)$$

$$o_t = \sigma \left( W_o \cdot [C_t, h_{t-1}, x_t] \; + \; b_o \right)$$

# Gated Recurrent Unit (GRU)
## *update (z), reset (r) gates*



$$z_t = \sigma\left(W_z \cdot [h_{t-1}, x_t]\right)$$

$$r_t = \sigma\left(W_r \cdot [h_{t-1}, x_t]\right)$$

$$\tilde{h}_t = \tanh\left(W \cdot [r_t * h_{t-1}, x_t]\right)$$

$$h_t = (1 - z_t) * h_{t-1} + z_t * \tilde{h}_t$$

# Long Short Term Memory (LSTM) for Time Series Forecasting

# Deep Learning for Time Series Prediction

# Deep Learning
# for
# Time Series Prediction
# Financial Market Prediction
## Stock Market Prediction
## Stock Price Prediction

# Time Series Data

[100, 110, 120, 130, 140, 150]

X

Y

Y

[100  110  120  130  140]  150

$X_{t1}$  $X_{t2}$  $X_{t3}$  $X_{t4}$  $X_{t5}$

# Long Short Term Memory (LSTM) for Time Series Forecasting

# Time Series Data

[10, 20, 30, 40, 50, 60, 70, 80, 90]

|  X  |  Y  |
|-----|-----|
| [10  20  30] | 40 |
| [20  30  40] | 50 |
| [30  40  50] | 60 |
| [40  50  60] | 70 |
| [50  60  70] | 80 |
| [60  70  80] | 90 |

# Deep Learning for Text Analytics (NLP)

# LSTM Recurrent Neural Network

# The Sequence to Sequence model (seq2seq)

44

# Sequence to Sequence (Seq2Seq)

# NLP

## Classical NLP



## Deep Learning-based NLP

# Modern NLP Pipeline

# Modern NLP Pipeline

# Deep Learning NLP

# Transformer (Attention is All You Need)
## (Vaswani et al., 2017)

# BERT:
# Pre-training of Deep Bidirectional Transformers for Language Understanding

**BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding**

**Jacob Devlin**    **Ming-Wei Chang**    **Kenton Lee**    **Kristina Toutanova**

Google AI Language

{jacobdevlin,mingweichang,kentonl,kristout}@google.com

# BERT

## Bidirectional Encoder Representations from Transformers



## Pre-training model architectures

**BERT** uses a bidirectional Transformer.
**OpenAI GPT** uses a left-to-right Transformer.
**ELMo** uses the concatenation of independently trained left-to-right and right- to-left LSTM to generate features for downstream tasks.
Among three, only BERT representations are jointly conditioned on both left and right context in all layers.

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

## BERT (Bidirectional Encoder Representations from Transformers)

## Overall pre-training and fine-tuning procedures for BERT

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018).
"Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

53

# BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding

BERT (Bidirectional Encoder Representations from Transformers)

## BERT input representation



The input embeddings is the sum of the token embeddings, the segmentation embeddings and the position embeddings.

# Fine-tuning BERT on Different Tasks



(a) Sentence Pair Classification Tasks: MNLI, QQP, QNLI, STS-B, MRPC, RTE, SWAG

(b) Single Sentence Classification Tasks: SST-2, CoLA

(c) Question Answering Tasks: SQuAD v1.1

(d) Single Sentence Tagging Tasks: CoNLL-2003 NER

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Fine-tuning BERT on Question Answering (QA)



(c) Question Answering Tasks: SQuAD v1.1

Source: Devlin, Jacob, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.

# Fine-tuning BERT on Dialogue Intent Detection (ID; Classification)



(b) Single Sentence Classification Tasks: SST-2, CoLA

# Fine-tuning BERT on Dialogue Slot Filling (SF)



(d) Single Sentence Tagging Tasks:
CoNLL-2003 NER

# General Language Understanding Evaluation (GLUE) benchmark

# GLUE Test results

| System | MNLI-(m/mm) 392k | QQP 363k | QNLI 108k | SST-2 67k | CoLA 8.5k | STS-B 5.7k | MRPC 3.5k | RTE 2.5k | Average - |
|---|---|---|---|---|---|---|---|---|---|
| Pre-OpenAI SOTA | 80.6/80.1 | 66.1 | 82.3 | 93.2 | 35.0 | 81.0 | 86.0 | 61.7 | 74.0 |
| BiLSTM+ELMo+Attn | 76.4/76.1 | 64.8 | 79.9 | 90.4 | 36.0 | 73.3 | 84.9 | 56.8 | 71.0 |
| OpenAI GPT | 82.1/81.4 | 70.3 | 88.1 | 91.3 | 45.4 | 80.0 | 82.3 | 56.0 | 75.2 |
| BERT$_{BASE}$ | 84.6/83.4 | 71.2 | 90.1 | 93.5 | 52.1 | 85.8 | 88.9 | 66.4 | 79.6 |
| BERT$_{LARGE}$ | **86.7/85.9** | **72.1** | **91.1** | **94.9** | **60.5** | **86.5** | **89.3** | **70.1** | **81.9** |

**MNLI**: Multi-Genre Natural Language Inference
**QQP**: Quora Question Pairs
**QNLI**: Question Natural Language Inference
**SST-2**: The Stanford Sentiment Treebank
**CoLA**: The Corpus of Linguistic Acceptability
**STS-B**:The Semantic Textual Similarity Benchmark
**MRPC**: Microsoft Research Paraphrase Corpus
**RTE**: Recognizing Textual Entailment

# Pre-trained Language Model (PLM)



Semi-supervised Sequence Learning
context2Vec
Pre-trained seq2seq

ULMFiT — ELMo

Transformer

Multi-lingual

Bidirectional LM

GPT

MultiFiT

Larger model
More data

BERT

Cross-lingual

GPT-2 — Defense → Grover

Multi-task

+ Generation

XLM
UDify

MT-DNN

Knowledge distillation

MASS
UniLM

MT-DNN_KD

Span prediction
Remove NSP

Longer time
Remove NSP
More data

SpanBERT

RoBERTa

+Knowledge Graph

Cross-modal

Whole Word Masking

Permutation LM
Transformer-XL
More data

XLNet

ERNIE
(Tsinghua)

Neural entity linker

KnowBert

VideoBERT
CBT
ViLBERT
VisualBERT
B2T2
Unicoder-VL
LXMERT
VL-BERT
UNITER

ERNIE (Baidu)
BERT-wwm

By Xiaozhi Wang & Zhengyan Zhang @THUNLP

# Turing Natural Language Generation (T-NLG)



T-NLG
17b

MegatronLM
8.3b

GPT-2
1.5b

BERT-Large
340m

RoBERTa
355m

DistilBERT
66m

2018

2019

2020

# Transformers

## State-of-the-art Natural Language Processing for TensorFlow 2.0 and PyTorch

- Transformers
  - pytorch-transformers
  - pytorch-pretrained-bert
- provides state-of-the-art general-purpose architectures
  - (BERT, GPT-2, RoBERTa, XLM, DistilBert, XLNet, CTRL...)
  - for Natural Language Understanding (NLU) and
    Natural Language Generation (NLG)
    with over 32+ pretrained models
    in 100+ languages
    and deep interoperability between
    TensorFlow 2.0 and
    PyTorch.

Source: https://github.com/huggingface/transformers

# Question Answering
# (QA)
# SQuAD

**Stanford Question Answering Dataset**

# SQuAD

## SQuAD 2.0
### The Stanford Question Answering Dataset

### What is SQuAD?

**S**tanford **Qu**estion **A**nswering **D**ataset (SQuAD) is a reading comprehension dataset, consisting of questions posed by crowdworkers on a set of Wikipedia articles, where the answer to every question is a segment of text, or *span*, from the corresponding reading passage, or the question might be unanswerable.

**SQuAD2.0** combines the 100,000 questions in SQuAD1.1 with over 50,000 unanswerable questions written adversarially by crowdworkers to look similar to answerable ones. To do well on SQuAD2.0, systems must not only answer questions when possible, but also determine when no answer is supported by the paragraph and abstain from answering.

### Leaderboard

SQuAD2.0 tests the ability of a system to not only answer reading comprehension questions, but also abstain when presented with a question that cannot be answered based on the provided paragraph.

| Rank | Model | EM | F1 |
|---|---|---|---|
| | Human Performance<br>*Stanford University*<br>(Rajpurkar & Jia et al. '18) | 86.831 | 89.452 |
| 1<br>Apr 06, 2020 | SA-Net on Albert (ensemble)<br>*QIANXIN* | **90.724** | **93.011** |
| 2<br>May 05, 2020 | SA-Net-V2 (ensemble)<br>*QIANXIN* | 90.679 | 92.948 |
| 2 | Retro-Reader (ensemble) | 90.578 | 92.978 |

https://rajpurkar.github.io/SQuAD-explorer/

64

# SQuAD

## SQuAD: 100,000+ Questions for Machine Comprehension of Text

**Pranav Rajpurkar** and **Jian Zhang** and **Konstantin Lopyrev** and **Percy Liang**

{pranavsr,zjian,klopyrev,pliang}@cs.stanford.edu

Computer Science Department

Stanford University

### Abstract

We present the Stanford Question Answering Dataset (SQuAD), a new reading comprehension dataset consisting of 100,000+ questions posed by crowdworkers on a set of Wikipedia articles, where the answer to each question is a segment of text from the corresponding reading passage. We analyze the dataset to understand the types of reasoning required to answer the questions, leaning heavily on dependency and constituency trees. We build a strong logistic regression model, which achieves an F1 score of 51.0%, a significant improvement over a simple baseline (20%). However, human performance (86.8%) is much higher, indicating that the dataset presents a good challenge problem for future research. The dataset is freely available at https://stanford-qa.com.

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under **gravity**. The main forms of precipitation include drizzle, rain, sleet, snow, **graupel** and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals **within a cloud**. Short, intense periods of rain in scattered locations are called "showers".

What causes precipitation to fall?
**gravity**

What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?
**graupel**

Where do water droplets collide with ice crystals to form precipitation?
**within a cloud**

**Figure 1:** Question-answer pairs for a sample passage in the

# SQuAD (Question Answering)

**Q:** What causes precipitation to fall?

## Precipitation

From Wikipedia, the free encyclopedia

*For other uses, see Precipitation (disambiguation).*

In meteorology, **precipitation** is any product of the condensation of atmospheric water vapor that falls under gravity from clouds.[2] The main forms of precipitation include drizzle, rain, sleet, snow, ice pellets, graupel and hail. Precipitation occurs when a portion of the atmosphere becomes saturated with water vapor (reaching 100% relative humidity), so that the water condenses and "precipitates". Thus, fog and mist are not precipitation but suspensions, because the water vapor does not condense sufficiently to precipitate. Two processes, possibly acting together, can lead to air becoming saturated: cooling the air or adding water vapor to the air. Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers."[3]

https://en.wikipedia.org/wiki/Precipitation

# SQuAD (Question Answering)

Paragraph

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail... Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail… Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

**A:** gravity

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail… Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

**A:** graupel

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail… Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** Where do water droplets collide with ice crystals to form precipitation?

**A:** within a cloud

# SQuAD (Question Answering)

In meteorology, precipitation is any product of the condensation of atmospheric water vapor that falls under gravity. The main forms of precipitation include drizzle, rain, sleet, snow, graupel and hail… Precipitation forms as smaller droplets coalesce via collision with other rain drops or ice crystals within a cloud. Short, intense periods of rain in scattered locations are called "showers".

**Q:** What causes precipitation to fall?

**A:** gravity

**Q:** What is another main form of precipitation besides drizzle, rain, snow, sleet and hail?

**A:** graupel

**Q:** Where do water droplets collide with ice crystals to form precipitation?

**A:** within a cloud

# Natural Language Processing with Python
## – Analyzing Text with the Natural Language Toolkit

www.nltk.org/book/

## Natural Language Processing with Python

## – Analyzing Text with the Natural Language Toolkit

**Steven Bird, Ewan Klein, and Edward Loper**

*This version of the NLTK book is updated for Python 3 and NLTK 3. The first edition of the book, published by O'Reilly, is available at http://nltk.org/book_1ed/. (There are currently no plans for a second edition of the book.)*

http://www.nltk.org/book/

# spaCy



https://spacy.io/

# gensim

# TextBlob

## TextBlob: Simplified Text Processing

Release v0.12.0. (Changelog)

*TextBlob* is a Python (2 and 3) library for processing textual data. It provides a simple API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.

○ Star    3,777

TextBlob is a Python (2 and 3) library for processing textual data. It provides a consistent API for diving into common natural language processing (NLP) tasks such as part-of-speech tagging, noun phrase extraction, sentiment analysis, and more.

## Useful Links

TextBlob @ PyPI
TextBlob @ GitHub
Issue Tracker

## Stay Informed

○ Follow @sloria

## Donate

If you find TextBlob useful,

```python
from textblob import TextBlob

text = '''
The titular threat of The Blob has always struck me as the ultimate movie
monster: an insatiably hungry, amoeba-like mass able to penetrate
virtually any safeguard, capable of--as a doomed doctor chillingly
describes it--"assimilating flesh on contact.
Snide comparisons to gelatin be damned, it's a concept with the most
devastating of potential consequences, not unlike the grey goo scenario
proposed by technological theorists fearful of
artificial intelligence run rampant.
'''

blob = TextBlob(text)
blob.tags            # [('The', 'DT'), ('titular', 'JJ'),
                     #  ('threat', 'NN'), ('of', 'IN'), ...]

blob.noun_phrases    # WordList(['titular threat', 'blob',
                     #           'ultimate movie monster',
                     #           'amoeba-like mass', ...])

for sentence in blob.sentences:
    print(sentence.sentiment.polarity)
# 0.060
```

https://textblob.readthedocs.io

# Polyglot

Search docs

Installation

Language Detection

Tokenization

Command Line Interface

Downloading Models

Word Embeddings

Part of Speech Tagging

Named Entity Extraction

Morphological Analysis

Transliteration

Sentiment

polyglot

Docs » Welcome to polyglot's documentation!     ○ Edit on GitHub

## Welcome to polyglot's documentation!

## polyglot

`downloads` `17k/month`  `pypi package` `16.7.4`  `build` `passing`  `docs` `passing`

Polyglot is a natural language pipeline that supports massive multilingual applications.

- Free software: GPLv3 license
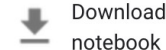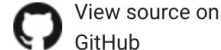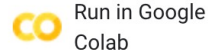- Documentation: http://polyglot.readthedocs.org.

## Features

- Tokenization (165 Languages)
- Language detection (196 Languages)
- Named Entity Recognition (40 Languages)
- Part of Speech Tagging (16 Languages)
- Sentiment Analysis (136 Languages)
- Word Embeddings (137 Languages)
- Morphological analysis (135 Languages)
- Transliteration (69 Languages)

https://polyglot.readthedocs.io/

# Text Classification with TF Hub

Overview    Tutorials    Guide    TF 1

TensorFlow tutorials
Quickstart for beginners
Quickstart for experts

**BEGINNER**

**ML basics with Keras** ^

  Basic image classification
  **Text classification with TF Hub**
  Text classification with preprocessed text
  Regression
  Overfit and underfit
  Save and load

**Load and preprocess data** ^

  CSV
  NumPy
  pandas.DataFrame
  Images
  Text
  Unicode
  TF.Text
  TFRecord and tf.Example
  Additional formats with tf.io ⬏

TensorFlow > Learn > TensorFlow Core > Tutorials         ☆ ☆ ☆ ☆ ☆

## Text classification with TensorFlow Hub: Movie reviews

[ CO  Run in Google Colab ]    [ ⚙ View source on GitHub ]    [ ⬇ Download notebook ]

This notebook classifies movie reviews as *positive* or *negative* using the text of the review. This is an example of *binary*—or two-class—classification, an important and widely applicable kind of machine learning problem.

The tutorial demonstrates the basic application of transfer learning with TensorFlow Hub and Keras.

We'll use the IMDB dataset that contains the text of 50,000 movie reviews from the Internet Movie Database. These are split into 25,000 reviews for training and 25,000 reviews for testing. The training and testing sets are *balanced,* meaning they contain an equal number of positive and negative reviews.

This notebook uses tf.keras, a high-level API to build and train models in TensorFlow, and TensorFlow Hub, a library and platform for transfer learning. For a more advanced text classification tutorial using `tf.keras`, see the MLCC Text Classification Guide.
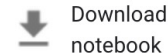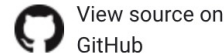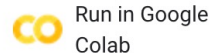
```
from __future__ import absolute_import, division, print_function, unicode_literals
```

**Contents**
Download the IMDB dataset
Explore the data
Build the model
  Loss function and optimizer
Train the model
Evaluate the model
Further reading

<u>https://www.tensorflow.org/tutorials/keras/text_classification_with_hub</u>

# Text Classification with Pre Text

TensorFlow tutorials
Quickstart for beginners
Quickstart for experts

BEGINNER

ML basics with Keras
  Basic image classification
  Text classification with TF Hub
  **Text classification with preprocessed text**
  Regression
  Overfit and underfit
  Save and load

Load and preprocess data
  CSV
  NumPy
  pandas.DataFrame
  Images
  Text
  Unicode
  TF.Text
  TFRecord and tf.Example
  Additional formats with tf.io

Estimator

☆ ☆ ☆ ☆ ☆

## Text classification with preprocessed text: Movie reviews

**Contents**
Setup
Download the IMDB dataset
Try the encoder
Explore the data
Prepare the data for training
Build the model
  Hidden units
  Loss function and optimizer
Train the model
Evaluate the model
Create a graph of accuracy and loss over time

[ CO ] Run in Google Colab    [ ] View source on GitHub    [ ↓ ] Download notebook

This notebook classifies movie reviews as *positive* or *negative* using the text of the review. This is an example of *binary*—or two-class—classification, an important and widely applicable kind of machine learning problem.

We'll use the IMDB dataset that contains the text of 50,000 movie reviews from the Internet Movie Database. These are split into 25,000 reviews for training and 25,000 reviews for testing. The training and testing sets are *balanced*, meaning they contain an equal number of positive and negative reviews.

This notebook uses tf.keras, a high-level API to build and train models in TensorFlow. For a more advanced text classification tutorial using `tf.keras`, see the MLCC Text Classification Guide.

## Setup

```
from __future__ import absolute_import, division, print_function, unicode_literals
```

https://www.tensorflow.org/tutorials/keras/text_classification

# Text Classification
# IMDB Movie Reviews

https://colab.research.google.com/drive/1x16h1GhHsLIrLYtPCvCHaoO1W-i_gror

Source: https://colab.research.google.com/github/tensorflow/docs/blob/master/site/en/tutorials/keras/basic_text_classification.ipynb

# Basic Regression
# Predict House Prices

📁 tf03_basic-regression.ipynb ☆

File   Edit   View   Insert   Runtime   Tools   Help

💬 COMMENT     👥 SHARE     Ⓐ

⊞ CODE   ⊞ TEXT   ⬆ CELL   ⬇ CELL          CONNECT ▾    ✏ EDITING   ⌃

▶ Copyright 2018 The TensorFlow Authors.

↳ 2 cells hidden

## ⌄ Predict house prices: regression

🔶 View on TensorFlow.org     CO Run in Google Colab     ⬛ View source on GitHub

In a *regression* problem, we aim to predict the output of a continuous value, like a price or a probability. Contrast this with a *classification* problem, where we aim to predict a discrete label (for example, where a picture contains an apple or an orange).

This notebook builds a model to predict the median price of homes in a Boston suburb during the mid-1970s. To do this, we'll provide the model with some data points about the suburb, such as the crime rate and the local property tax rate.

This example uses the `tf.keras` API, see this guide for details.

```
1  # memory footprint support libraries/code
2  !ln -sf /opt/bin/nvidia-smi /usr/bin/nvidia-smi
3  !pip install gputil
4  !pip install psutil
5  !pip install humanize
6  import psutil
7  import humanize
8  import os
9  import GPUtil as GPU
10 GPUs = GPU.getGPUs()
11 gpu = GPUs[0]
12 def printm():
13     process = psutil.Process(os.getpid())
14     print("Gen RAM Free: " + humanize.naturalsize( psutil.virtual_memory().available ), " | Proc size: "
15     print("GPU RAM Free: {0:.0f}MB | Used: {1:.0f}MB | Util {2:3.0f}% | Total {3:.0f}MB".format(gpu.memo
```

# Papers with Code
# State-of-the-Art (SOTA)

## Browse State-of-the-Art

1509 leaderboards • 1327 tasks • 1347 datasets • 17810 papers with code

Follow on Twitter for updates

## Computer Vision

**Semantic Segmentation**
33 leaderboards
667 papers with code

**Image Classification**
52 leaderboards
564 papers with code

**Object Detection**
54 leaderboards
467 papers with code

**Image Generation**
51 leaderboards
231 papers with code

**Pose Estimation**
40 leaderboards
231 papers with code

▸ See all 707 tasks

## Natural Language Processing

**Machine Translation**

**Language Modelling**

**Question Answering**

**Sentiment Analysis**

**Text Generation**

https://paperswithcode.com/sota

# Aurélien Géron (2019),
# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition
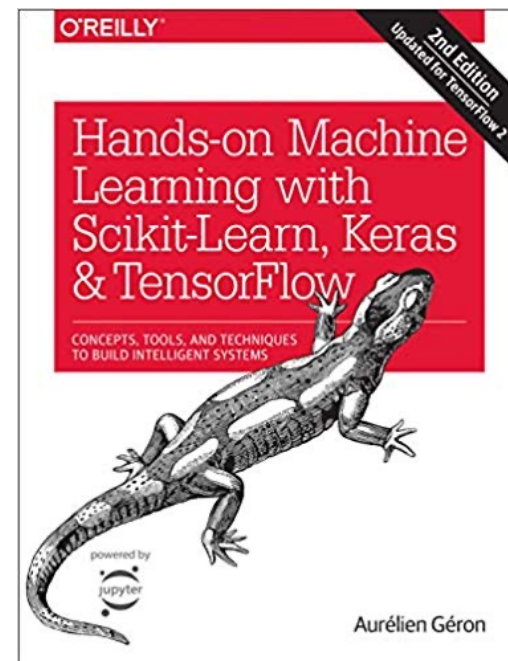# O'Reilly Media, 2019



https://github.com/ageron/handson-ml2

# Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow

**Notebooks**

1. The Machine Learning landscape
2. End-to-end Machine Learning project
3. Classification
4. Training Models
5. Support Vector Machines
6. Decision Trees
7. Ensemble Learning and Random Forests
8. Dimensionality Reduction
9. Unsupervised Learning Techniques
10. Artificial Neural Nets with Keras
11. Training Deep Neural Networks
12. Custom Models and Training with TensorFlow
13. Loading and Preprocessing Data
14. Deep Computer Vision Using Convolutional Neural Networks
15. Processing Sequences Using RNNs and CNNs
16. Natural Language Processing with RNNs and Attention
17. Representation Learning Using Autoencoders
18. Reinforcement Learning
19. Training and Deploying TensorFlow Models at Scale

https://github.com/ageron/handson-ml2

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

python101.ipynb

File   Edit   View   Insert   Runtime   Tools   Help      All changes saved

Comment      Share

+ Code      + Text

RAM / Disk      Editing

## Text Classification: IMDB Movie Review

- Source: https://www.tensorflow.org/tutorials/keras/text_classification_with_hub

```
[1]  1 !pip install -q tensorflow-hub
     2 !pip install -q tensorflow-datasets
```

```
[2]  1 import os
     2 import numpy as np
     3
     4 import tensorflow as tf
     5 import tensorflow_hub as hub
     6 import tensorflow_datasets as tfds
     7
     8 print("Version: ", tf.__version__)
     9 print("Eager mode: ", tf.executing_eagerly())
    10 print("Hub version: ", hub.__version__)
    11 print("GPU is", "available" if tf.config.list_physical_devices("GPU") else "NOT AVAILABLE")
```

```
Version:  2.4.1
Eager mode:  True
Hub version:  0.12.0
GPU is available
```

```
[3]  1 # Split the training set into 60% and 40% to end up with 15,000 examples
     2 # for training, 10,000 examples for validation and 25,000 examples for testing.
     3 train_data, validation_data, test_data = tfds.load(
     4     name="imdb_reviews",
```

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

https://colab.research.google.com/drive/1FEG6DnGvwfUbeo4zJ1zTunjMqf2RkCrT

python101.ipynb

File  Edit  View  Insert  Runtime  Tools  Help      All changes saved

Comment      Share

+ Code    + Text

RAM
Disk

Editing

- Huggingface Transformers: https://github.com/huggingface/transformers

```
[18]  1 !pip install transformers
```

```
1 from transformers import pipeline
2 classifier = pipeline('sentiment-analysis')
3 classifier('We are very happy to introduce pipeline to the transformers repository.')
```

Downloading: 100%  ████████████████  629/629 [00:00<00:00, 1.31kB/s]

Downloading: 100%  ████████████████  268M/268M [00:05<00:00, 46.9MB/s]

Downloading: 100%  ████████████████  232k/232k [00:01<00:00, 159kB/s]

Downloading: 100%  ████████████████  48.0/48.0 [00:00<00:00, 522B/s]

[{'label': 'POSITIVE', 'score': 0.9996980428695679}]

```
[11]  1 classifier('This movie is very good.')
```

[{'label': 'POSITIVE', 'score': 0.9998621940612793}]

```
[12]  1 classifier('This movie is very boring.')
```

[{'label': 'NEGATIVE', 'score': 0.999795138835907}]

https://tinyurl.com/aintpupython101

# Python in Google Colab (Python101)

# Summary

- **Recurrent Neural Networks (RNN)**

  - **Long Short Term Memory (LSTM)**

  - **Gated Recurrent Unit (GRU)**

- **Deep Learning (RNN) for Time Series Prediction**

- **Deep Learning (RNN) for Text Analytics (NLP)**

# References

- Aurélien Géron (2019), Hands-On Machine Learning with Scikit-Learn, Keras, and TensorFlow: Concepts, Tools, and Techniques to Build Intelligent Systems, 2nd Edition, O'Reilly Media. https://github.com/wesm/pydata-book
- Stuart Russell and Peter Norvig (2020), Artificial Intelligence: A Modern Approach, 4th Edition, Pearson.
- Christopher Olah, (2015) Understanding LSTM Networks, http://colah.github.io/posts/2015-08-Understanding-LSTMs/
- Martin Gorner (2017), Tensorflow, deep learning and modern RNN architectures, without a PhD, https://www.youtube.com/watch?v=pzOzmxCR37I
- Brandon Rohrer (2017), Recurrent Neural Networks (RNN) and Long Short-Term Memory (LSTM), https://www.youtube.com/watch?v=WCUNPb-5EYI
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova (2018). "Bert: Pre-training of deep bidirectional transformers for language understanding." arXiv preprint arXiv:1810.04805.
- Min-Yuh Day (2021), Python 101, https://tinyurl.com/aintpupython101