Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00000

References
0000

## Attentions Mechanisms and Transformers

E. A. León-Gómez [1]

ealeon@unal.edu.co[1]
Universidad Nacional de Colombia

May 30, 2023

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00000

References
0000

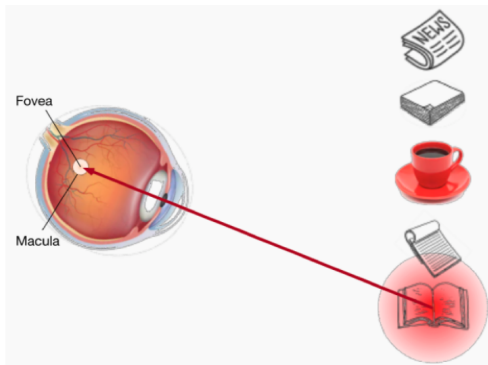## Overview

1 Introduction

2 Attentions mechanism
   - Attention

3 Transformers
   - Input embedding
   - Positional encoding
   - Encoder
   - Decoder

4 Random Feature Attention

**Introduction**
●○○

Attentions mechanism
○○○○○○○○

Transformers
○○○○○○○○○○○○○○

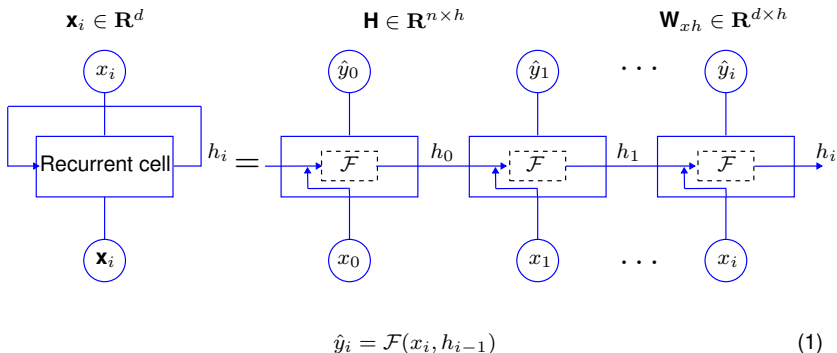Random Feature Attention
○○○○○

References
○○○○

# Introduction

## What is Attention?

Attention is a critical mechanism in deep learning that enables models to concentrate on the **most relevant parts** of input data for the given task at hand



Credit:http://d2l.ai/

**Introduction**
○○●
Attentions mechanism
○○○○○○○○
Transformers
○○○○○○○○○○○○○
Random Feature Attention
○○○○○
References
○○○○

## Recurrent Neural Networks (RNN)



$$\hat{y}_i = \mathcal{F}(x_i, h_{i-1}) \tag{1}$$

### Problems

- Distant positions in the sequence can be disregarded
- Parallelizing the work is challenging because it processes variables sequentially

Introduction
000

Attentions mechanism
●0000000

Transformers
0000000000000000

Random Feature Attention
00000

References
0000

# Attentions mechanism

| Introduction | Attentions mechanism | Transformers | Random Feature Attention | References |
|---|---|---|---|---|
| 000 | 0●000000 | 0000000000000 | 00000 | 0000 |

Attention

## Learning task

$$\{\mathbf{x}_i\}_{i=1}^n = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_d\} \; ; \; \mathbf{x}_i \in \mathcal{R}^d \qquad \mathbf{h} = \mathbf{x}_1\alpha_1 + \mathbf{x}_2\alpha_2 + \cdots + \mathbf{x}_d\alpha_d \; ; \mathbf{h} \in \mathcal{R}^n$$

$$\mathbf{X} \in \mathcal{R}^{n \times d} \qquad\qquad\qquad\qquad \mathbf{H} \in \mathcal{R}^{n \times d}$$

Self-attention computes the output sequence **H** from **X** as follows:

- Projection the input into diferrent subspaces
- Computing the output as a weighted average
- Multi-Head Attention

Introduction
000

Attentions mechanism
○○●○○○○○○

Transformers
○○○○○○○○○○○○○○○

Random Feature Attention
○○○○○

References
○○○○

Attention

# Attention as search

## Projection the input into diferrent subspaces

The input **X** is transformed into the query matrix **Q**, the key matrix **K**, and the value matrix **V** via three linear transformations:

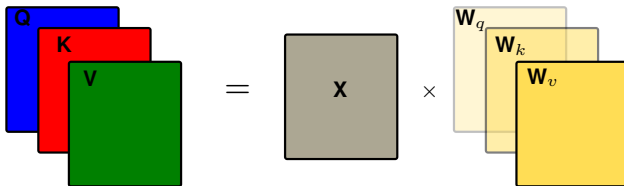$$\mathbf{Q} = \mathbf{X}\mathbf{W}_Q^\top \qquad\qquad \mathbf{K} = \mathbf{X}\mathbf{W}_K^\top \qquad\qquad \mathbf{V} = \mathbf{X}\mathbf{W}_V\top$$

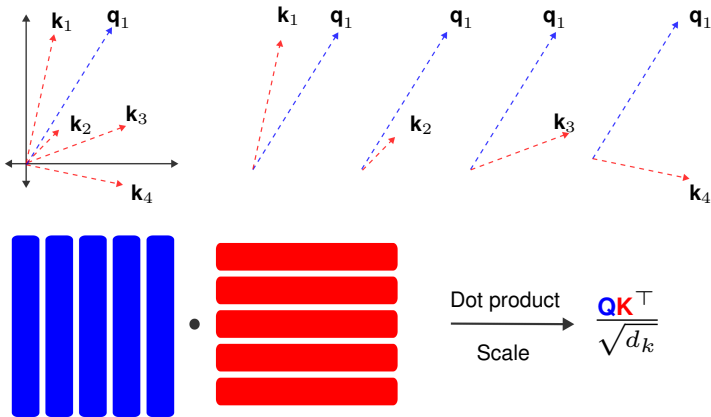$$\mathbf{Q} \in \mathcal{R}^{n \times d_k}, \mathbf{W}_Q \in \mathcal{R}^{d \times d_k} \qquad \mathbf{K} \in \mathcal{R}^{n \times d_k}, \mathbf{W}_K \in \mathcal{R}^{d \times d_k} \qquad \mathbf{V} \in \mathcal{R}^{n \times d_v}, \mathbf{W}_V \in \mathcal{R}^{d \times d_v}$$

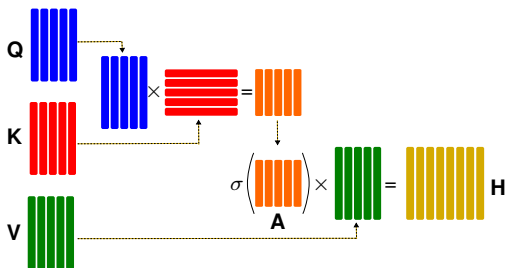| Introduction | Attentions mechanism | Transformers | Random Feature Attention | References |
|---|---|---|---|---|
| 000 | 00000●000 | 0000000000000 | 00000 | 0000 |

Attention

# Interpretability of query arrays, keys, and values

Query arrays, keys, and values can be considered an "information retrieval" system



Dot product
Scale
$$\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}$$

Introduction
000

Attentions mechanism
00000●00

Transformers
0000000000000

Random Feature Attention
00000

References
0000

Attention

# Computing the output as a weighted average



$$\mathbf{H} = \underbrace{softmax\left(\frac{\mathbf{QK}^\top}{\sqrt{d_k}}\right)}_{\mathbf{A}}\mathbf{V}$$
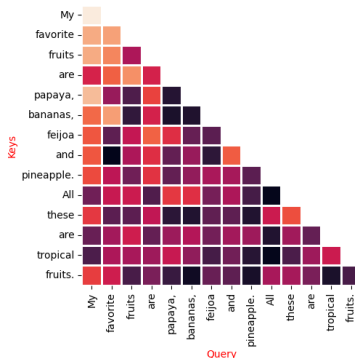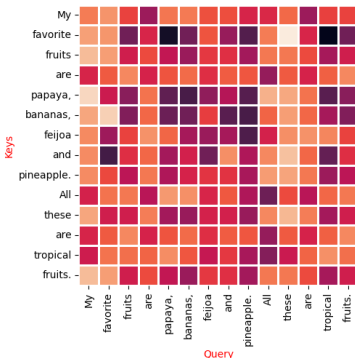
- The self-attention capture the intra-correlation of a given input **X**
- Where $\mathbf{A} \in \mathcal{R}^{n \times n}$ is a probability distribution over the element of **K**

Introduction
○○○

Attentions mechanism
○○○○○○○●○○

Transformers
○○○○○○○○○○○○○○○

Random Feature Attention
○○○○○

References
○○○○

Attention

# Masked attention

$$\mathbf{H} = softmax\left(\frac{\mathbf{QK}^\top + \mathbf{M}}{\sqrt{d_k}}\right)\mathbf{V}$$
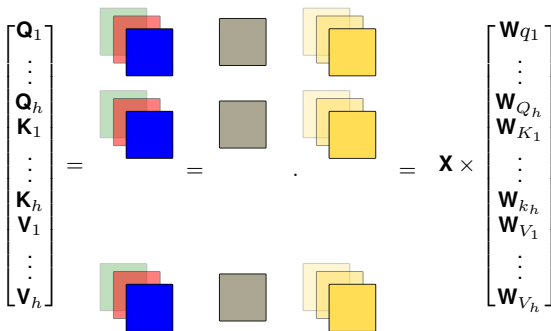
$$\mathbf{M} = \begin{bmatrix} 0 & -\infty & -\infty & \dots & -\infty \\ 0 & 0 & -\infty & \dots & -\infty \\ 0 & 0 & 0 & \dots & -\infty \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & 0 & \dots & 0 \end{bmatrix}$$

$$\mathbf{M} \in R^{n \times n} \; ; M_{i,j} = \{0 \; if \; i \le j, -\infty \, if \, i > j\}$$

| Introduction | Attentions mechanism | Transformers | Random Feature Attention | References |
|---|---|---|---|---|
| 000 | 0000000● | 0000000000000 | 00000 | 0000 |

Attention

# Multi-Head Attention

Each output sequence **H** forms an attention head. Multi-head attention concatenates multiple heads to compute the final output.

$$MultiHead\left(\{\mathbf{Q}, \mathbf{K}, \mathbf{V}\}_{i=1}^{H}\right) = Concat\left(\mathbf{H}_1, \mathbf{H}_2, \ldots, \mathbf{H}_h\right)\mathbf{W}^O \tag{2}$$

Introduction
ooo

Attentions mechanism
oooooooo

Transformers
●oooooooooooooo

Random Feature Attention
ooooo

References
oooo

# Transformers

Introduction
000

Attentions mechanism
00000000

Transformers
0●000000000000

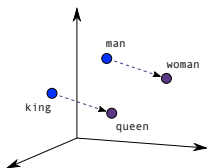Random Feature Attention
00000

References
0000

## How deep learning can have attention?



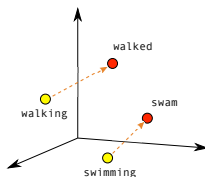### Why?

- Distinct parts of the input convey unique information
- Retain memories of specific, interconnected events from the past

Introduction    Attentions mechanism    **Transformers**    Random Feature Attention    References
○○○              ○○○○○○○○                 ○○●○○○○○○○○○○○      ○○○○○                      ○○○○

Input embedding

# Word Embeddings for NLP



Figure: Embeddings representation [1]

## Word embeddings models

- Bag of words(BOW)
- Word2Vec
- GloVe: Global Vector for word representation

---
1

https://developers.google.com/machine-learning/crash-course/embeddings/translating-to-a-lower-dimensional-space

Introduction
000

Attentions mechanism
00000000

**Transformers**
0000●000000000

Random Feature Attention
00000

References
0000

Positional encoding

# Intuition

Transformers use wave frequency to find positional information



$$\mathbf{p}_{pos,t} = \begin{cases} \sin(w_k t), & \text{if } pos = 2k \\ \cos(w_k t), & \text{if } pos = 2k + 1 \end{cases}$$

$$w_k = \frac{pos}{10000^{\frac{2k}{d}}}$$

- $w_k$:
- $pos$: Position of input in the sequence
- $d$: size of input/word/token
- $i$: Individual dimension of the embedding

| Introduction | Attentions mechanism | **Transformers** | Random Feature Attention | References |
| 000 | 00000000 | 0000●00000000 | 00000 | 0000 |

Positional encoding

# Example

$$
\mathbf{p}_t = \begin{bmatrix} \sin(w_1 t) \\ \cos(w_1 t) \\ \\ \sin(w_2 t) \\ \cos(w_2 t) \\ \\ \vdots \\ \\ \sin\left(w_{\frac{d}{2}} t\right) \\ \cos\left(w_{\frac{d}{2}} t\right) \end{bmatrix}_{d \times 1}
$$

# Positional encoding [2]

Positional encoding for a sentence of 100 words and dimension embedding of 500

Introduction
000

Attentions mechanism
00000000
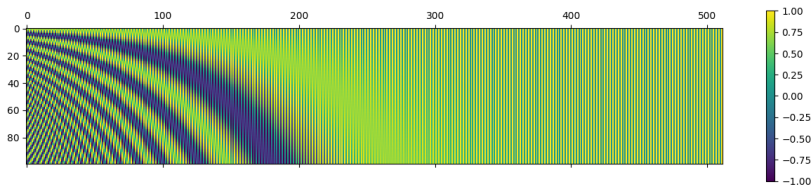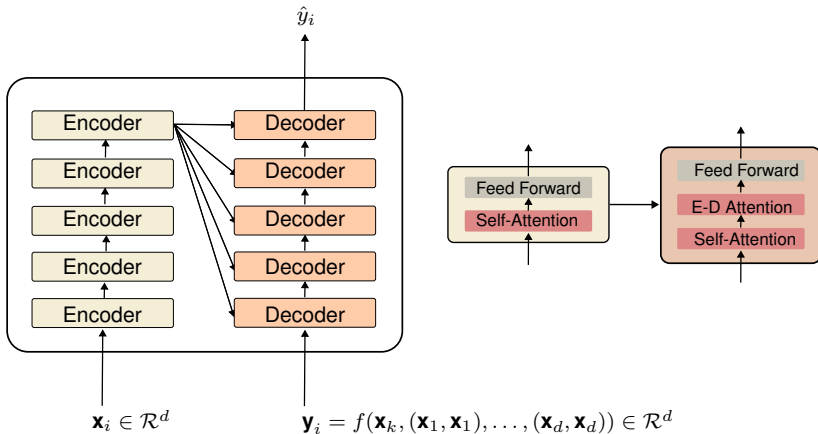
**Transformers**
000000●0000000

Random Feature Attention
00000

References
0000

Encoder

## Learning Self-Attention with Neural Network



$$\hat{y}_i$$

| | |
|---|---|
| Encoder | Decoder |
| Encoder | Decoder |
| Encoder | Decoder |
| Encoder | Decoder |
| Encoder | Decoder |

Feed Forward
Self-Attention

Feed Forward
E-D Attention
Self-Attention

$$\mathbf{x}_i \in \mathcal{R}^d \qquad \mathbf{y}_i = f(\mathbf{x}_k, (\mathbf{x}_1, \mathbf{x}_1), \ldots, (\mathbf{x}_d, \mathbf{x}_d)) \in \mathcal{R}^d$$

### Features

- The encoders are equal but with different weights
- **E-D Attention**: Allows you to focus on relevant parts of the input

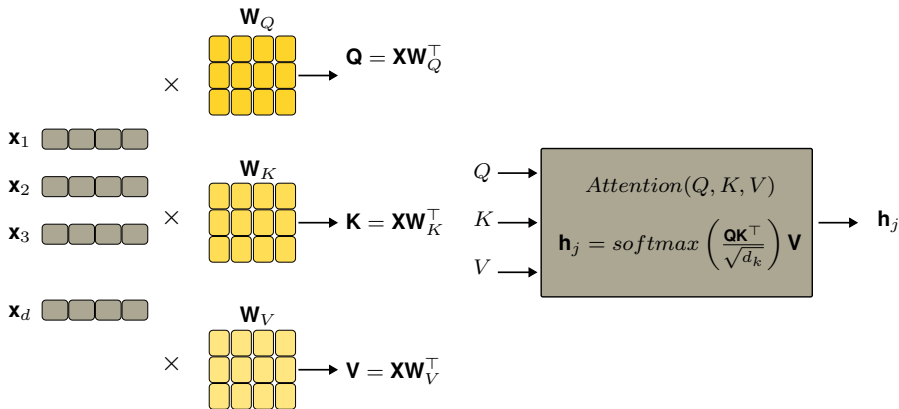| Introduction | Attentions mechanism | **Transformers** | Random Feature Attention | References |
|---|---|---|---|---|
| ○○○ | ○○○○○○○○ | ○○○○○○○●○○○○○○ | ○○○○○ | ○○○○ |

Encoder

# The Encoder Block

| Introduction | Attentions mechanism | **Transformers** | Random Feature Attention | References |
| 000 | 00000000 | 000000000000000 | 00000 | 0000 |

Encoder

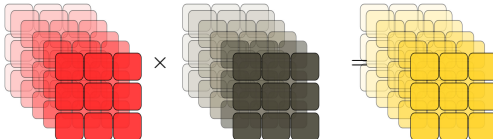# Self-Attention: Query, keys, and values

The dot-product attention fail to capture the correlations between these features



Where $\mathbf{W}_Q, \mathbf{W}_K \in \mathcal{R}^{d_k \times n}$ and $\mathbf{W}_Q \in \mathcal{R}^{d_v \times n}$ are represent learnable weight matrices.

| Introduction | Attentions mechanism | **Transformers** | Random Feature Attention | References |
| 000 | 00000000 | 000000000●0000 | 00000 | 0000 |

Encoder

# Multi-Head Self-Attention

The multi-head attention Self-Attention consists of multiple attention operators with different similarity function determined by different groups of weight matrices.

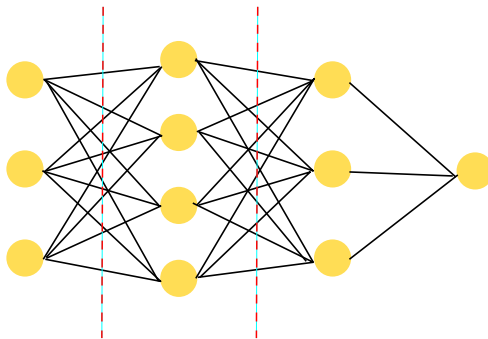$$\mathbf{z}_j = \mathbf{h}_j \mathbf{w}_j^o \qquad (3)$$



$$\mathbf{h}_j \in \mathcal{R}^{t \times d} \qquad \mathbf{w}_j^o \in \mathcal{R}^{d \times d} \qquad \mathbf{z}_j \in \mathcal{R}^{t \times d}$$

The multi-head attention allows each head to attend different locations based on the similarity in different representation subspaces.

Introduction
○○○

Attentions mechanism
○○○○○○○○○

**Transformers**
○○○○○○○○○○●○○○

Random Feature Attention
○○○○○

References
○○○○

Encoder

# Layer normalization



$$\mathbf{z}'_j = LayerNorm(\mathbf{z}_j + \mathbf{x}_j)$$

$$\mathbf{x}' = \delta(W_1^\top + b_1)$$

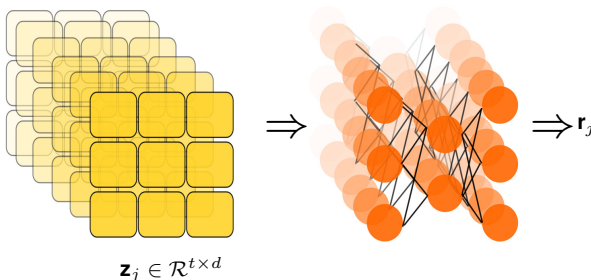$$\mathbf{z}'_j = \gamma_1 \left( \frac{\mathbf{x} - \mu_1}{\sigma} \right) + \beta_1$$

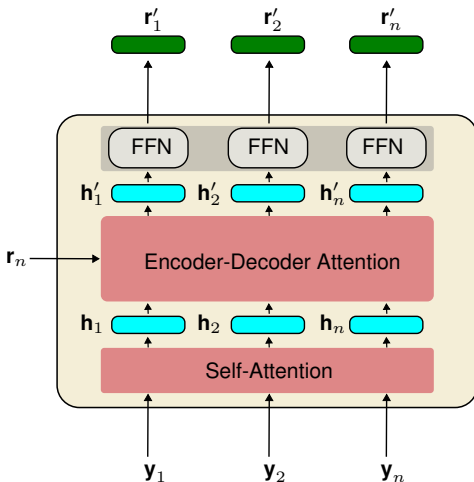| Introduction | Attentions mechanism | Transformers | Random Feature Attention | References |
| 000 | 00000000 | 000000●0000000 | 00000 | 0000 |

Encoder

# Full connected network

$$FFN(\mathbf{h}_j) = \mathbf{r}_j = ReLU(\mathbf{h}_j \mathbf{W}^1 + \mathbf{b}^1)\mathbf{W}^2 + \mathbf{b}^2$$
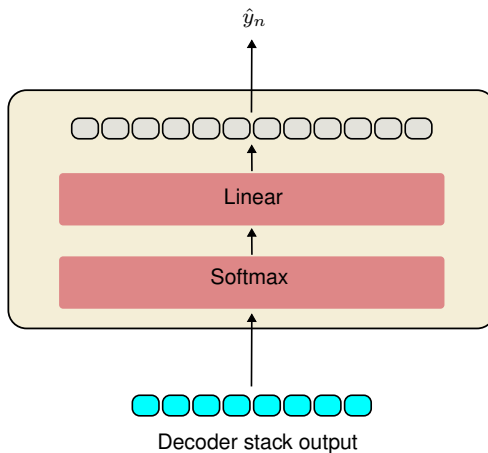


$\mathbf{z}_j \in \mathcal{R}^{t \times d}$

- Where $\mathbf{W}^1$, $\mathbf{W}^2$, $\mathbf{b}^1$ and $\mathbf{b}^2$ are parameters.
- The feed-forward layer is connected back-to-back with residual connections and normalization layers.
- The output of an encoder block is used as an input for the next encoder block.

| Introduction 000 | Attentions mechanism 00000000 | **Transformers** 00000000000●0 | Random Feature Attention 00000 | References 0000 |

Decoder

# The Decoder Block

Each decoder block consists of similar layers and operations as the encoder block.

Introduction
000

Attentions mechanism
00000000

Transformers
00000000000000●

Random Feature Attention
00000

References
0000

Decoder

# The Output

# Random Feature Attention

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00●000

References
0000
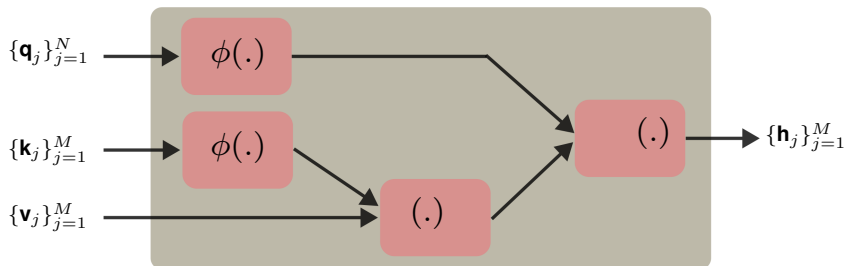
## Complexity computation for softmax attention



$$\mathbf{h} = \mathbf{D}^{-1} \exp \underbrace{\left( \frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right)}_{\hat{\mathbf{A}} \in \mathcal{R}^{n \times n}} \mathbf{V} = \mathbf{D}^{-1} \hat{\mathbf{A}} \mathbf{V} \; ; \mathbf{D} = diag\left( \hat{\mathbf{A}} \mathbf{1} \right)$$

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00●00

References
0000

## Linearized attention

Objetive

$$\exp\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \approx \phi(\mathbf{Q})\phi(\mathbf{K})^\top$$

Introduction
000

Attentions mechanism
00000000

Transformers
000000000000

Random Feature Attention
00000

References
0000

## Attention in kernel machine I

### Theorem: Random Fourier Features [4]

Let $\phi : \mathcal{R}^d \to \mathcal{R}^{2D}$ be transformation:

$$\phi(\mathbf{x}) = \sqrt{\frac{1}{D}} \left[ \sin(\mathbf{w}_1 \mathbf{x}), \dots, \sin(\mathbf{w}_D \mathbf{x}), \dots, \cos(\mathbf{w}_1 \mathbf{x}), \dots, \cos(\mathbf{w}_1 \mathbf{x}) \right]^\top$$

When d-dimensional random vectors $\mathbf{w}_i$ are independently sampled from $\mathcal{N}(0, \sigma^2 \mathbf{I}_d)$

$$\mathbf{h} = softmax \left( \frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right) \mathbf{V}; \ softmax(\mathbf{a}) = \frac{\exp(\mathbf{a}_i)}{\sum_k \exp(\mathbf{a}_k)}$$

$$\mathbf{h} = \frac{\exp \left( \frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right)}{\sum_{j=1}^n \exp \left( \frac{\mathbf{QK}^\top}{\sqrt{d_k}} \right)} \mathbf{V}$$

---

[4] https:
//proceedings.neurips.cc/paper/2007/hash/013a006f03dbc5392effeb8f18fda755-Abstract.html

## Attention in kernel machine II

$$\mathbf{D} = \sum_{j=1}^{n} \underbrace{\exp\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)}_{\hat{\mathbf{A}}} = \sum_{j=1}^{n} \hat{\mathbf{A}} = diag(\hat{\mathbf{A}}\mathbf{1}^\top)$$

$$\mathbf{h} = \frac{\exp\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right)}{\mathbf{D}}\mathbf{V}$$

$$\exp\left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}}\right) \approx \phi(\mathbf{Q})\phi^\top(\mathbf{K})$$

$$\boxed{\mathbf{h} = \mathbf{D}^{-1}\phi(\mathbf{Q})\phi^\top(\mathbf{K})\mathbf{V}}$$

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00000

References
●000

📄 Illustrated transformer/
*https: // jalammar. github. io/ illustrated- transformer/*.

📄 Rahimi, Ali and Recht, Benjamin
Attention is All You Need
*Advances in neural information processing systems (2007).*

📄 Vuckovic, James and Baratin, Aristide and Combes, Remi Tachet des
A mathematical theory of attention
*arXiv preprint arXiv (2020).*

📄 Ji, Shuiwang and Xie, Yaochen and Gao, Hongyang
A mathematical view of attention models in deep learning
*Texas A&M University: College Station, TX, USA (2019).*

📄 Nguyen, Tan and Pham, Minh and Nguyen, Tam and Nguyen, Khai and Osher,
Stanley J and Ho, Nhat
Transformer with Fourier Integral Attentions
*rXiv preprint arXiv (2022).*

📄 AAhmed, Sabeen and Nielsen, Ian E and Tripathi, Aakash and Siddiqui,
Shamoon and Rasool, Ghulam and Ramachandran, Ravi P
Transformers in time-series analysis: a tutorial
*rXiv preprint arXiv (2022).*

Rahimi, Ali and Recht, Benjamin
Random features for large-scale kernel machines
*Advances in neural information processing systems.*

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00000

References
0000

### Transformers

- Self-attention https://github.com/ealeongomez/Machine-Learning/blob/master/Attention-mechanisms/1_selfAttention.ipynb
- Multi Head Attention:
  https://github.com/ealeongomez/Machine-Learning/blob/master/Attention-mechanisms/2-TensorFlow_MultiHeadAttention.ipynb

Introduction
000

Attentions mechanism
00000000

Transformers
0000000000000

Random Feature Attention
00000

References
000●

Thank you!