## Attention Meets Post-hoc Interpretability: A Mathematical Perspective

Gianluigi Lopardo,<sup>1</sup> Frederic Precioso,<sup>1</sup> Damien Garreau<sup>2</sup>

<sup>1</sup>Université Côte d'Azur - Inria <sup>2</sup>Julius-Maximilians Universität Würzburg - CAIDAS

January 8, 2025





## Where is Würzburg?





 Figure: (left panel) Würzburg is located in Northern Bavaria (right panel) Festung Marienberg (credits google maps / wikipedia)

## Outline

#### 1. Transformers

- 2. Post-hoc interpretability Gradient-based approaches Perturbation-based approaches Attention weights
- 3. Analysis on a simple model

#### 4. Conclusion

# 1. Transformers

### Introduction

▶ Context = natural language processing: from text input x ∈ X, predict y ∈ Y as f(x)
 ▶ Running example: sentiment analysis:

the selection on the menu is great and so is the food the service \_\_\_\_\_ **>** positive is not bad prices are fine

- ▶  $f = f_{\theta}$  corresponds to some architecture choice, and  $\theta \in \Theta$  to the parameters
- $\theta^{\star} = \text{good parameter learned from data}$
- State-of-the-art toady:  $f_{\theta}$  = attention<sup>1</sup>-based model (a transformer<sup>2</sup>)
- **Goal of this section:** describe a simple transformer architecture

<sup>&</sup>lt;sup>1</sup>Bahdanau, Cho, Bengio, *Neural machine translation by jointly learning to align and translate*, ICLR, 2025 <sup>2</sup>Vaswani et al., *Attention is all you need*, NeurIPS, 2017

## Tokens

- **Notation:**  $\xi \in \mathcal{X}$  is a document
- encoded for the computer as a sequence of tokens
- ▶ we identify tokens with elements of  $\{1, ..., D\} = [D]$
- Several possibilities in practice:
  - individual letters (D = 100)
  - words (D = 100.000)
  - ▶ in-between (*e.g.*, BERT<sup>3</sup> uses WordPiece,  $^4 D = 30.000$ )
- Example:

"DATAIA" 
$$\mapsto$$
 [4, 1, 20, 1, 9, 1]

**Special tokens:** <UNK>, <CLS>, etc.

<sup>&</sup>lt;sup>3</sup>Devlin et al., *BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding*, ACL Proc., 2019

<sup>&</sup>lt;sup>4</sup>Wu et al., *Google's Neural Machine Translation System: Bridging the Gap between Human and Machine Translation*, preprint, 2016

### Token embeddings

- Summary from last slide:  $\xi = \text{ ordered sequence of tokens } \xi_1, \dots, \xi_T \in [D]$
- Next step: vector representation of each token
- ▶ for each  $t \in [T]$ , token  $\xi_t = j$  is embedded as

$$e_t := (W_e)_{:,j} + W_p(t) \in \mathbb{R}^{d_e},$$

where:

▶ 
$$W_e \in \mathbb{R}^{d_e imes D}$$
 matrix containing embeddings of all tokens

•  $W_p : \mathbb{N} \to \mathbb{R}^{d_e}$  positional embedding

**Typically,**  $W_e$  is learned and  $W_p$  is set to something arbitrary

**Example:** 

$$\begin{cases} W_p(t)_{2i} &= \cos(t/T_{\max}^{2i/d_e}) \\ W_p(t)_{2i-1} &= \sin(t/T_{\max}^{2i/d_e}) \,. \end{cases}$$

### Keys, queries, values

Note: I am following

Phuong and Hutter, Formal Algorithms for Transformers, preprint, 2022 Padding until  $T_{max}$  with <UNK> token, embedded to  $h + W_p(t)$ 

**Next step:** for all  $t \in [t]$ , map embeddings to

$$\begin{cases} k_t &:= W_k e_t + b_k \in \mathbb{R}^{d_{\text{att}}} & (\text{key}) \\ q_t &:= W_q e_t + b_q \in \mathbb{R}^{d_{\text{att}}} & (\text{query}) \\ v_t &:= W_v e_t + b_v \in \mathbb{R}^{d_{\text{out}}} & (\text{value}) \end{cases}$$

▶ matrices  $W_k, W_q \in \mathbb{R}^{d_{\mathsf{att}} \times d_e}$ ,  $W_v \in \mathbb{R}^{d_{\mathsf{out}} \times d_e}$  learned parameters

▶ bias vectors  $b_k, b_q \in \mathbb{R}^{d_{\text{att}}}, b_v \in \mathbb{R}^{d_{\text{out}}}$  also learnable, set to zero for simplicity

#### Attention mechanism

▶ for a given query  $q \in \mathbb{R}^{d_{\text{att}}}$ , index *t* receives *attention* 

$$\forall t \in T_{\max}, \qquad \alpha_t := \frac{\exp\left(q^\top k_t / \sqrt{d_{\mathsf{att}}}\right)}{\sum_{u=1}^{T_{\max}} \exp\left(q^\top k_u / \sqrt{d_{\mathsf{att}}}\right)},$$

softmax of the vector  $(q^{ op}k_1,\ldots,q^{ op}k_{\mathcal{T}_{\max}})^{ op}$ 

- **Intuition:** if query "matches" with  $k_t$ , then  $\alpha_t$  large
- ▶ this mechanism is at the core of the transformer architecture

	[CLS]	ai	iņ	paris	is	great	but	the	weather	is	bad
[CLS]·	0.11	0.19	0.00	0.04	0.14	0.16	0.00	0.00	0.00	0.04	0.17

### Final output

 $\blacktriangleright$  for a given query q, output

$$\widetilde{\mathbf{v}} := \sum_{s=1}^{T_{\max}} lpha_s \mathbf{v}_s \in \mathbb{R}^{d_{\mathrm{out}}}$$
 .

▶ In our simplified setting, q corresponds to the <CLS> token,  $W_{\ell} \in \mathbb{R}^{1 \times d_{\text{out}}}$ , and

$$f(x) := W_\ell \tilde{v} \in \mathbb{R}$$
.

**Remark:** we can deal with several heads:

$$\forall i \in [\mathcal{K}], \qquad \tilde{v}^{(i)} := \sum_{s=1}^{T_{\max}} \alpha_s^{(i)} v_s^{(i)} \in \mathbb{R}^{d_{\text{out}}},$$

and

$$f(x) := rac{1}{K} \sum_{i=1}^K W^{(i)}_\ell ilde{\mathbf{v}}^{(i)} \in \mathbb{R}$$
 .

### Our model, in pictures



### In practice

In reality, self-attention, meaning each token associated to value v

 *v*<sub>t</sub>
 then layer-norm + small perceptron, several layers



Figure: transformer architecture from Vaswani et al. (2017)

# 2. Post-hoc interpretability

## A (very brief) introduction to interpretability

- **Context:** automated systems are reaching human-level in many applications
- Problem: sometime performance is not the only metric (especially in critical applications)
- Interpretability methods: give insights to why a specific decision was taken
- this talk = local, post-hoc (single example, model already trained)
- **Example:** sentiment analysis, outline words that are important for the decision



Figure: Anchors outlining words supporting a positive prediction

# 2.1. Gradient-based approaches

### Gradients

- **Simple idea:** take the gradient of prediction with respect to input
- Intuition: if feature is (positively) important, positive partial derivative
- Problem: documents are discrete objects
- **Solution:** rewrite *f* as a function of the embeddings, *i.e.*,

$$f(x) = g(e_1(x), e_2(x), \ldots, e_{\mathcal{T}_{\max}}(x))$$

▶ gradient-based approaches compute, for each token,  $abla_{e_t} g \in \mathbb{R}^{d_e}$ 



### Token-level measure of importance

> Problem: still complicated object, need to come back to token level

- several competing approaches:
  - ▶ take the mean (G-AVG)<sup>5</sup>
  - take the  $L^1$  norm  $(G-L1)^6$
  - take the  $L^2$  norm (G-L2)<sup>7</sup>
  - ▶ take the dot-product with  $e_t$  (G×I)<sup>8</sup>
  - ▶ ...

**Example:** *L*<sup>2</sup> norm of the token gradients:

attention based explanations are popular but questionable

<sup>&</sup>lt;sup>5</sup>Atanasova et al., *A diagnostic study of explainability techniques for text classification*, EMNLP, 2020 <sup>6</sup>Li et al., *Visualizing and understanding neural models in NLP*, ACL Proc., 2016

<sup>&</sup>lt;sup>7</sup>Poerner et al., *Evaluating neural network explanation methods using hybrid documents and morphosyntactic agreement*, ACL Proc., 2018

<sup>&</sup>lt;sup>8</sup>Denil et al., *Extraction of salient sentences from labelled documents*, preprint, 2014

# 2.2. Perturbation-based approaches

### Perturbation-based approaches

- Idea: remove words at random and look at how the prediction varies
- Example: LIME<sup>9</sup>
- ▶ recall  $\xi = (\xi_1, ..., \xi_T) \in [D]^T$  document, f our model
- ▶ local dictionary [d] with d < D
- **Step 1:** create *n* perturbed samples  $X_1, \ldots, X_n$  by removing *s* words at random
- s follows uniform distribution on [d]
- the subset S of words to be removed is chosen uniformly
- words are removed with repetition

<sup>9</sup>Ribeiro et al., "Why should i trust you?" Explaining the predictions of any classifier, SIGKDD, 2016

## Sampling



**Figure:** LIME sampling on a small example

## Weights

- **Step 2:** give positive weights to the samples
- define  $Z_i \in \{0,1\}^T$  as presence / absence of words
- 1 corresponds to the original document
- weights are defined by

$$\forall i \in [n], \qquad \pi_i := \exp\left(\frac{-\delta(\mathbbm{1}, Z_i)^2}{2\nu^2}\right) \,,$$

with  $\nu > 0$  bandwidth parameter and  $\delta$  the cosine distance

$$\delta(\pmb{a},\pmb{b}):=1-rac{\pmb{a}^{ op}\pmb{b}}{\|\pmb{a}\|\cdot\|\pmb{b}\|}\,.$$

**Intuition:** if perturbed sample far from  $\xi$ ,  $\delta(\mathbb{1}, X_i) \approx 1$ , assign small weight

### Local surrogate model

**Step 3:** local surrogate model

fit linear model on absence / presence of words:

$$\hat{eta}_n \in \operatorname*{arg\,min}_{eta \in \mathbb{R}^d} \left\{ \sum_{i=1}^n \pi_i (Y_i - eta^ op Z_i)^2 + \lambda \left\|eta
ight\|^2 
ight\} \,,$$

where 
$$Y_i = f(X_i)$$
 and  $\lambda > 0$ 

#### Explaining a prediction with LIME





**Figure:** visualizing LIME output

# 2.3. Attention weights

### Attention weights

- Another idea: look directly at attention weights<sup>10</sup>
- ▶ In our setting, attention wrt <CLS> token  $\Rightarrow$  look at attention weight of individual tokens
- What happens with several heads?
  - either take the average

$$\alpha\text{-}\mathsf{avg}_t := \frac{1}{K}\sum_{i=1}^K \alpha_t^{(i)}\,,$$

or the max<sup>11</sup>

$$\alpha\operatorname{-max}_t := \max_{i \in [K]} \alpha_t^{(i)} \, .$$

▶ if several layers, further aggregation scheme required<sup>12</sup>

<sup>&</sup>lt;sup>10</sup>Clark et al., What does BERT look at? An analysis of BERT's attention, Blackbox @ EMNLP, 2019 <sup>11</sup>Schwenke and Atzmueller, Show me what you're looking for: visualizing abstracted transformer attention for enhancing their local interpretability on time series data, FLAIRS, 2021

<sup>&</sup>lt;sup>12</sup>Mylonas et al., An attention matrix for every decision: faithfulness-based arbitration among multiple attention-based interpretations of transformers in text classification, Data Mining and Knowledge Discovery, 2024

### Post-hoc interpretability III: attention, in pictures



Figure: attention patterns inside a single-layer transformer

### Is attention explanation?

- tempting to rely on these coefficients: they are really used by the model
- **But**, some dissident voices:<sup>13</sup>
  - ▶ if attention is explanation, attention coefs should correlate with feature importance
  - counterfactual attention weight configuration should change prediction
- the debate is not settled
- ▶ there are criticisms regarding experimental setting of Jain and Wallace<sup>14</sup>
- not many theoretical contributions
- related work show that single-layer attention models can get near-perfect accuracy with un-informative attention pattern<sup>15,16</sup>

<sup>&</sup>lt;sup>13</sup>Jain and Wallace, Attention is not explanation, NAACL Proc., 2019

<sup>&</sup>lt;sup>14</sup>Wiegreffe and Pinter, Attention is not not Explanation, EMNLP, 2019

<sup>&</sup>lt;sup>15</sup>Wen et al., *Transformers are uninterpretable with myopic methods: a case study with bounded Dyck grammars*, NeurIPS, 2024

<sup>&</sup>lt;sup>16</sup>Cui, Behrens, Krzakala, Zdeborová, A phase transition between positional and semantic learning in a solvable model of dot-product attention, preprint, 2024

### Is attention explanation?, ctd.

- Histogram task: count number of times token appears in the sequence<sup>17</sup>
- **Example:** "DATAIAIAIA"  $\mapsto$  [1, 5, 1, 5, 3, 5, 3, 5, 3, 5]
- Architecture: single-layer with tied weights
- two vastly different local minima found, one with un-informative attention pattern





figure obtained running code from Cui et al., 2024

<sup>17</sup>Weiss, Goldberg, Yohav, *Thinking like transformers*, ICML, 2021

# 3. Analysis on a simple model

## Summary so far

- Many different methods providing explanations
- different results even on single-layer transformer:

G-avg:	attention	based	explanations	are	popular	but	questionable
G-I1:	attention	based	explanations	are	popular	but	questionable
G-I2:	attention	based	explanations	are	popular	but	questionable
G×I:	attention	based	explanations	are	popular	but	questionable
lime:	attention	based	explanations	are	popular	but	questionable
lpha-avg:	attention	based	explanations	are	popular	but	questionable
$\alpha$ -max:	attention	based	explanations	are	popular	but	questionable

- **Figure:** different explainers yield different explanations
- Our work: what should we use?
- starting point = attention coefficients =  $\alpha_t$
- Problem: no dependency in the linear layer / values (!)

### Gradient-based: closed-form expression

- **Recall:** we are looking at  $\nabla_{e_t} g$
- straightforward computations yield:

**Theorem (Lopardo, Precioso, G., '24):** The gradient of our simple model with respect to token  $e_t$  is given by

$$abla_{e_t} g(x) = lpha_t W_{\mathsf{v}}^\top W_{\ell}^\top + rac{lpha_t}{\sqrt{d_{\mathsf{att}}}} W_{\ell} \left( \mathsf{v}_t - \sum_{s=1}^{T_{\mathsf{max}}} lpha_s \mathsf{v}_s 
ight) W_k^\top q \in \mathbb{R}^{d_e}$$

- **•** Main insight: token contributes if  $\alpha_t \neq 0$  and  $v_t$  deviates from average
- **Remark (i):** easy corollary for K heads by linearity
- **Remark (ii):** straightforward derivations for average,  $L^1$  and  $L^2$  norms, etc.

### Additional notation

- much more challenging analysis for LIME
- set h the index for the <UNK> token
- $\blacktriangleright$  corresponding key / value vectors for <UNK> token at position t are

$$egin{cases} k_{h,t} &:= W_k h + W_k W_
ho(t) \in \mathbb{R}^{d_{ ext{att}}} \ v_{h,t} &:= W_
u(h + W_
ho(t)) \in \mathbb{R}^{d_{ ext{out}}} \,. \end{cases}$$

define further

$$g_{h,t} := \exp\left(q^{ op} k_{h,t}/\sqrt{d_{\mathsf{att}}}
ight) \,,$$

and

$$\alpha_{h,t} := \frac{g_{h,t}}{\sum_{u} g_{h,u}} \, .$$

Intuition: attention coefficient for all <UNK> tokens

### Perturbation-based, main result

with these notation:

**Theorem (Lopardo, Precioso, G., '24):** Assume that  $d = T = T_{\max}^{\varepsilon}$ , with  $\varepsilon \in (0, 1)$ . Assume further that there exist positive constants 0 < c < C such that, as  $T \to +\infty$ , for all  $t \in [T_{\max}]$ ,  $\max(|v_t|, |v_{h,t}|) \leq C$ , and  $c \leq \min(g_t, g_{h,t}) \leq C$ .

$$orall j \in [d], \qquad eta_j^\infty pprox rac{3}{2} \sum_{t=1}^{T_{ ext{max}}} W_\ell \left( lpha_t m{v}_t - lpha_{h,t} m{v}_{h,t} 
ight) \mathbbm{1}_{\xi_t=j}.$$

- approximate expression of LIME coefficients for a single-layer transformer
- **•** Main insight: token contributes if  $\alpha_t v_t$  deviates from "average"
- Remark: straightforward extension to several heads

### Experimental check



Figure: boxplots = 5 runs of LIME, red crosses = approximation. T = d = 99 in this example

## Sketch of proof (i)

we use previous work doing the analysis in the asymptotic setting:<sup>18</sup>

**Theorem (Mardaoui, G., '21):** take  $\lambda = 0$ , assume f is bounded, then  $\hat{\beta} \xrightarrow{\mathbb{P}} \beta^{f}$ , where  $\beta^{f}$  is defined as

$$\forall j \in [d], \quad \beta_j^f = c_d^{-1} \left\{ \sigma_1 \mathbb{E} \left[ \pi f(X) \right] + \sigma_2 \mathbb{E} \left[ \pi Z_j f(X) \right] + \sigma_3 \sum_{\substack{k=1\\k \neq j}}^d \mathbb{E} \left[ \pi Z_k f(X) \right] \right\}.$$

Here, X has the distribution of the perturbed document described previously, and  $c_d$ ,  $\sigma_1, \sigma_2, \sigma_3$  have explicit expressions.

#### **Intuition:** weighted least squares $\rightarrow$ closed-form

<sup>18</sup>Mardaoui and Garreau, An analysis of LIME for text data, AISTATS, 2021

## Sketch of proof (ii)

▶ in the large bandwidth regime, expression simplifies somewhat:

**Corollary (Mardaoui and G., '21):** same assumptions,  $\nu \to +\infty$ , then  $\beta_i^f$  converges to

$$\beta_j^{\infty} = 3\mathbb{E}\left[f(X) \mid j \notin S\right] - \frac{3}{d} \sum_k \mathbb{E}\left[f(X) \mid k \notin S\right],$$

where S is the random set defined in the sampling scheme.

very challenging to deal with this expectation (f non-linear and complicated distribution)

we resort to approximations

## Sketch of proof (iii)

Crux of the proof: approximate

$$\mathbb{E}\left[f(X) \mid j \notin S\right] = \mathbb{E}\left[\sum_{t=1}^{T_{\max}} A_t V_t \mid \ell \notin S\right] = \sum_{t=1}^{T_{\max}} \mathbb{E}\left[\frac{G_t V_t}{\sum_{u=1}^{T_{\max}} G_u} \mid \ell \notin S\right],$$

where  $A_t$  and  $V_t$  are the random version of attention / values

Proof technique: split expectation according to |S| = s, then approximate each piece using the following:

**Lemma:** Let X and Y be two random variables with finite variance. Assume that there exist two positive constants c and C such that  $|X| \le C$  and  $cn \le Y \le Cn$  a.s. Then

$$\mathbb{E}\left[\frac{X}{Y}\right] - \frac{\mathbb{E}\left[X\right]}{\mathbb{E}\left[Y\right]} \le \frac{C \operatorname{Var}\left(Y\right)}{c^3 n^3} + \frac{C^2 \sqrt{\operatorname{Var}\left(Y\right)}}{c^2 n^2}$$

## Sketch of proof (iv)

▶ finally, computation of expectation and variance of

$$\mathcal{H}_{\mathcal{S}} := \sum_{i} \left\{ a_{i} \mathbb{1}_{i \notin \mathcal{S}} + b_{i} \mathbb{1}_{i \in \mathcal{S}} \right\} \,,$$

conditionally to  $|\mathcal{S}| = s$  and  $\ell \notin \mathcal{S}$ 

**Lemma:** Let  $H_s$  be as before, then

$$\mathbb{E}_s\left[\mathcal{H}_S|\ell 
otin S
ight] = rac{n-1-s}{n-1}\sum_i \mathsf{a}_i + rac{s}{n-1}\sum_i b_i + rac{s}{n-1}(\mathsf{a}_\ell - b_\ell)\,,$$

and

$$\mathsf{Var}_{s}(\mathsf{H}_{\mathsf{S}} \mid \ell \notin \mathsf{S}) = \frac{ns(n-s-1)}{(n-1)(n-2)} \left[ \widehat{\mathsf{Var}} \left( \mathsf{a} - \mathsf{b} \right) - \frac{1}{n-1} \left( \mathsf{a}_{\ell} - \mathsf{b}_{\ell} - \left( \overline{\mathsf{a}} - \overline{\mathsf{b}} \right) \right)^{2} \right] \,.$$

# 4. Conclusion

## Conclusion

#### Summary:

- single-layer attention-based model
- closed-form or exact approximations for token-importance measure
- methods are very un-alike no clear recommended choice

#### Main reference:

Lopardo, Precioso, Garreau, Attention Meets Post-hoc Interpretability: A Mathematical Perspective, ICML, 2024

#### Future directions:

- Anchors<sup>19</sup> meets attention (existing theoretical framework<sup>20</sup>)
- more realistic architecture (skip connection, non-linearities, more layers)

<sup>19</sup>Ribeiro, Singh, Guestrin, Anchors: High-precision model-agnostic explanations, AAAI, 2018
 <sup>20</sup>Lopardo, Precioso, Garreau, A sea of words: an in-depth analysis of anchors for text data, AISTATS, 2023

## Thank you for your attention!