# Feature Attribution:

## Introduction to the Classic Methodologies

2024-07-19

Li Peng-Hsuan 李朋軒

# Agenda

- Why feature attribution

- Classic methodologies—the general, the good and the sound
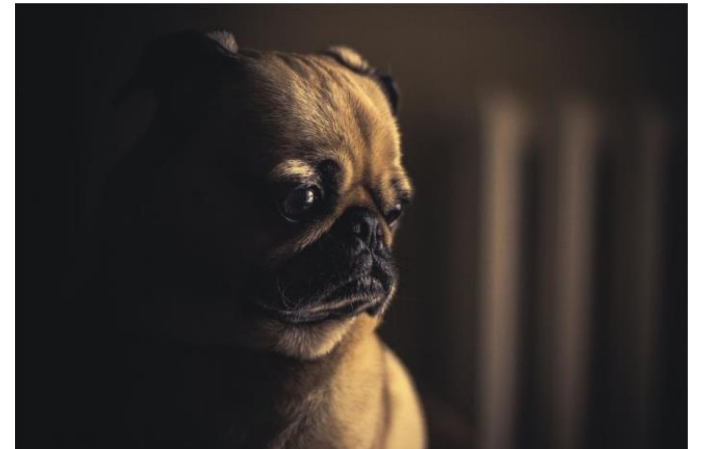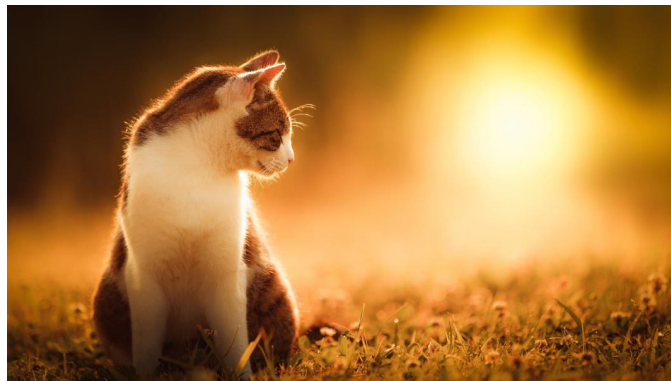
- The SOTA and the limit

# Agenda

- Why feature attribution

- Classic methodologies—the general, the good and the sound

- The SOTA and the limit

# Why Feature Attribution

- Model validation
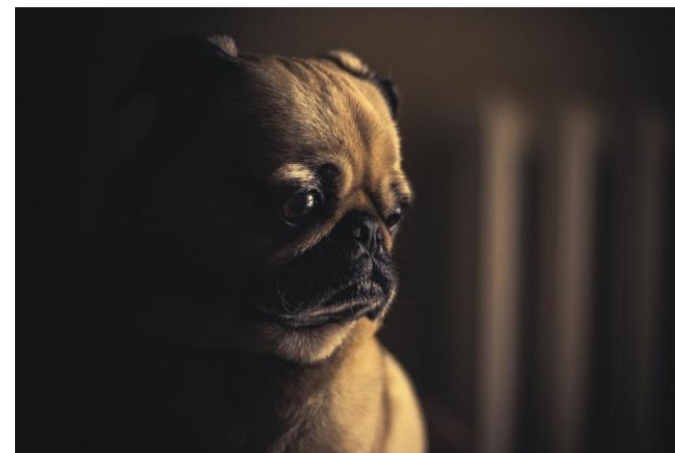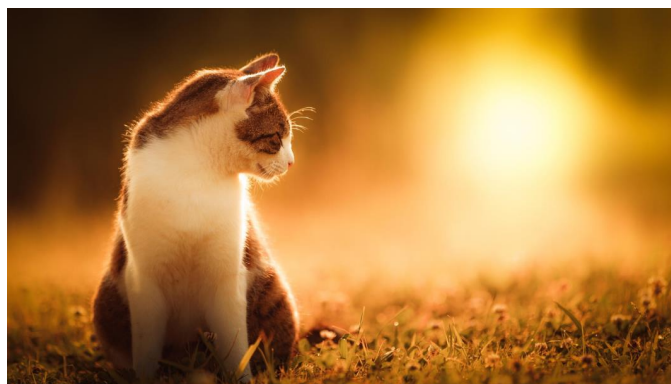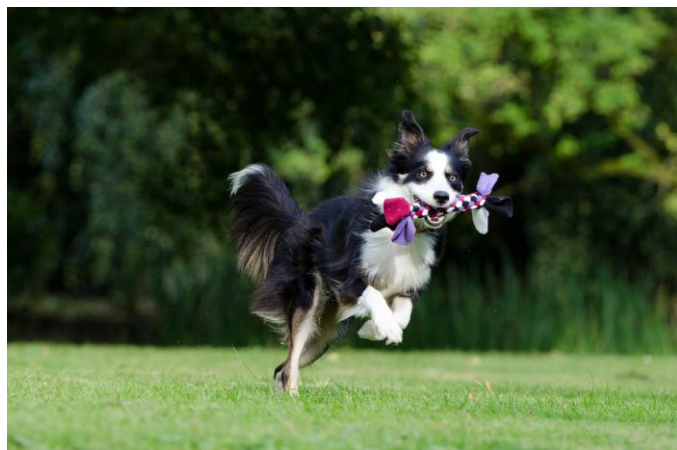
- Knowledge discovery

# Cats and Dogs
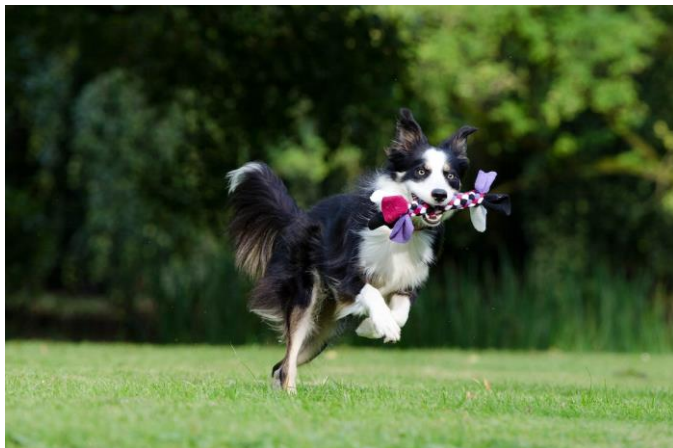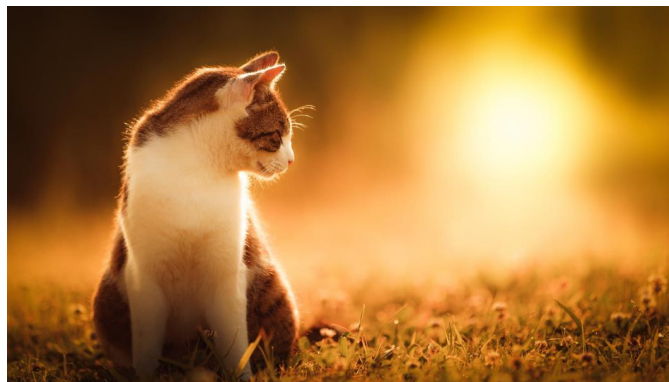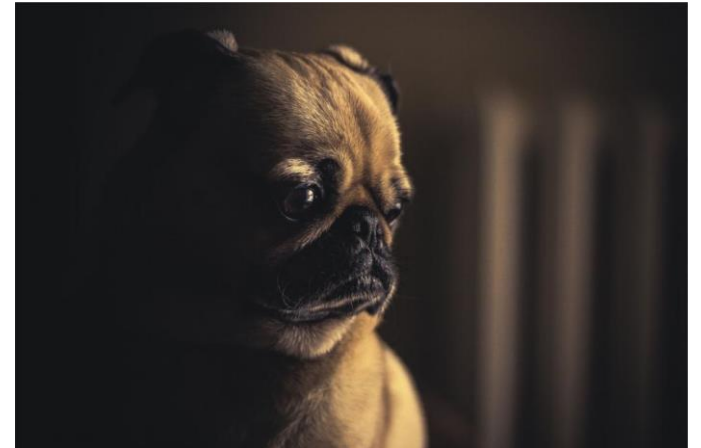
cat / dog

cat / dog

cat / dog

cat / dog

cat
cat



dog   dog



cat   cat



dog
dog



cat   cat



dog
dog

important feature to me

cat / dog

cat

dog

cat

dog

cat

dog

# Sequence & Motif



BPNet architecture

Output: predicted ChIP–nexus profile and counts for each TF

Regression of bias track

Convolutional layers with dilations and residual skip connections

Input: ~150,000 sequences of 1 kb

Function

Generation of ChIP–nexus footprints

Detection of motifs in context

Motif detection

Sequence scanning

CAGATGCATAACAAAGGTGC

10

Input: trained BPnet model

Backtracking
of signal
through
network

... $^A_e$ATGCATAAcAAAGG$_⊥$e$^c$ ...

Output: profile contribution scores for each TF

# Agenda

- Why feature attribution

- Classic methodologies—the general, the good and the sound

- The SOTA and the limit

# The General, the Good and the Sound

- The *shapley* methodology

- The *inversion* methodology

- The *gradient* methodology

# The General

*Awesome Cooperative Development Plan*

.........................................................................
.........................legalize marijuana trade.................
......drill oil from the Arctic Ocean...................................
.........................................................................
.....................send dissidents to the moon..............
.........................................................................

Signature:

$100B

*Awesome Cooperative Development Plan*

..................................................................
.............................legalize marijuana trade..............
......drill oil from the Arctic Ocean.............................
..................................................................
.......................send dissidents to the moon...............
..................................................................

Signature:

$100B

$25B    $25B    $25B    $25B

# Awesome Cooperative Development Plan

........................................................
....................legalize marijuana trade...............
......drill oil from the Arctic Ocean........................
........................................................
...................send dissidents to the moon...............
........................................................

Signature:

$100B

$25B    $25B    $25B    $25B

# Awesome Cooperative Development Plan

..................................................................................
..............................legalize marijuana trade................
......drill oil from the Arctic Ocean....................................
..................................................................................
........................send dissidents to the moon................
..................................................................................

Signature:



$\cancel{\$100B}$

$\$10B$



We contributed $90B!

## Order 1

Worth $10B

$10B

Worth $25B

$10B $15B

Worth $50B

$10B $15B $25B

Worth $100B

$10B $15B $25B $50B

## Order 2

Worth $10B

$10B

Worth $45B

$10B $35B

Worth $80B

$10B $35B $35B

Worth $100B

$10B $35B $35B $20B

# Economics vs. ML



Worth $100B

attribute

$$

vs.

$$y = f(\boldsymbol{x})$$

attribute

(feature) $x_i$

(contribution) $c_i$

# The *Shapley* Methodology

$$\$100B = \Box(\ \text{🔴🔵⬛🍁}\ )$$

$$\to$$

$$\text{🍁}$$
$$\$\$$$

$$y = f(\boldsymbol{x})$$

(feature) $\quad x_i$

(contribution) $\quad c_i$

attribute

(reference) $\quad r_i$

Sample $x$



Reference $r$

Sample $x$

Reference $r$

Population → Sequence → Assembled genomes → Calculate allele frequencies → Population-specific consensus genome ← Update allele frequencies ← Additional genomes

# The Good

# Attribution for a linear function

$$y = f(\boldsymbol{x}) = \boldsymbol{x} \cdot \boldsymbol{w} + b$$



$$c_i = x_i w_i$$

# Attribution for ReLU

$$y = f(x) = \max(x, 0)$$

$y$

$1[x \geq 0]$

$x$

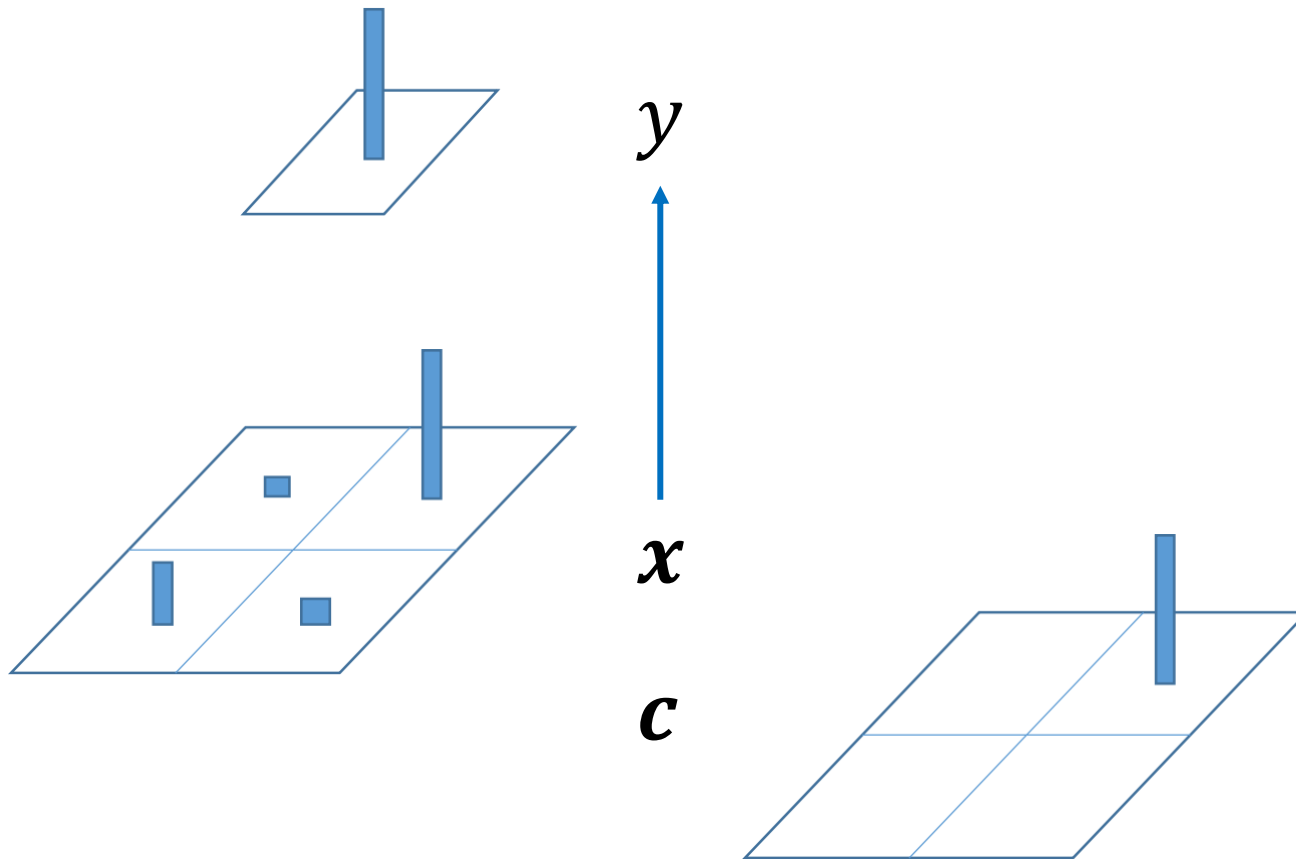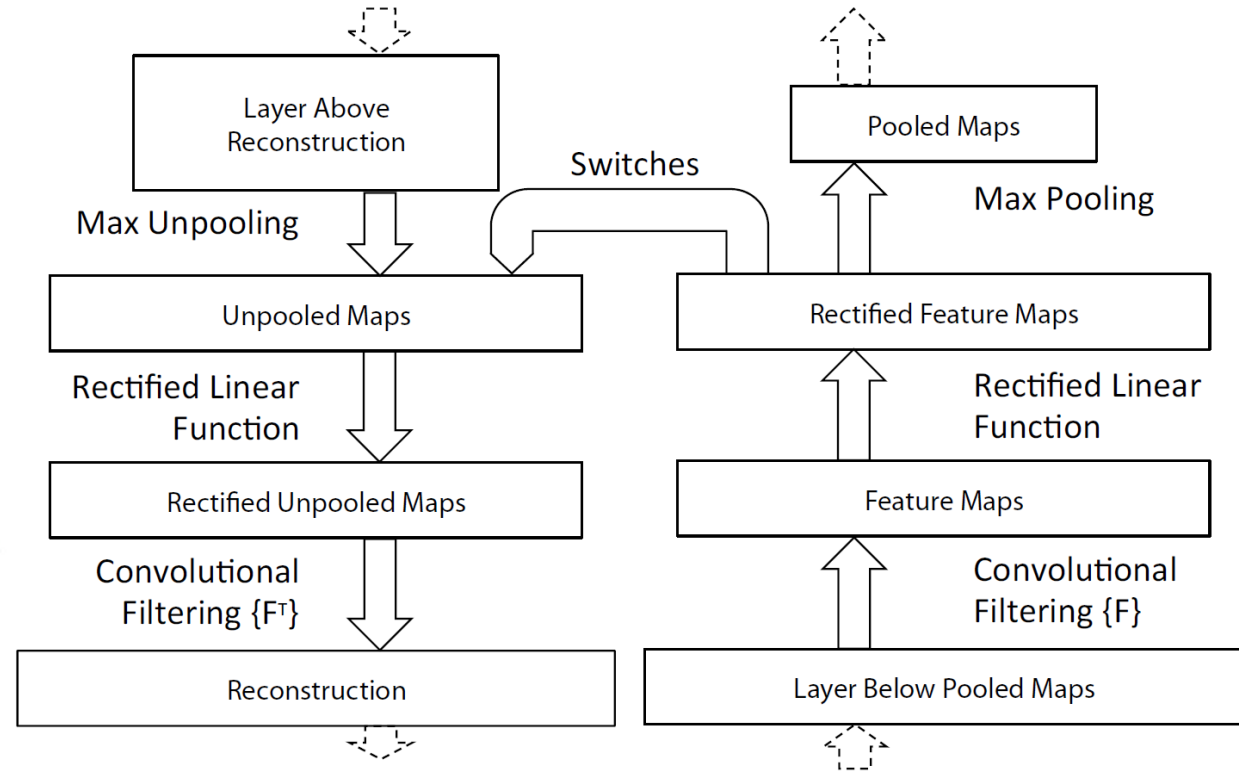$$c = x \cdot 1[x \geq 0]$$

# Attribution for Max-Pooling



$$y = f(\boldsymbol{x}) = \max(\boldsymbol{x})$$

# Attribution for CNN

Layer 4

Layer 5

$$y = f(x) = g_3\Big(g_2\big(g_1(x)\big)\Big) = g_3 \circ g_2 \circ g_1(x)$$

$$x = f^{-1}(y) = g_1^{-1} \circ g_2^{-1} \circ g_3^{-1}(y)$$

$$y = f(x) = g_3\left(g_2\big(g_1(x)\big)\right) = g_3 \circ g_2 \circ g_1(x)$$

$$x = f^{-1}(y) = g_1^{-1} \circ g_2^{-1} \circ g_3^{-1}(y)$$

🤔

Hey! But functions are not generally invertible!

# The *Inversion* Methodology

Given a model $f$

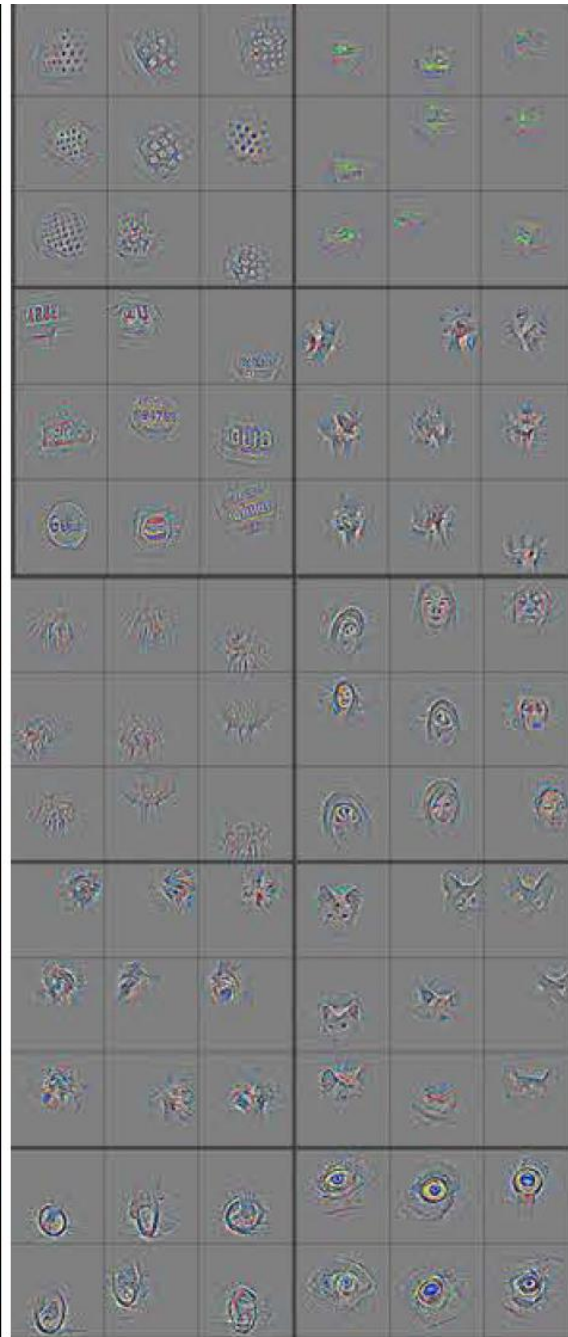1. Break down composition $\qquad$ $f = g_n \circ \cdots \circ g_2 \circ g_1$

2. Formulate inverse $\qquad\qquad$ $g_i^{-1}$

3. Reduce inverse to linear $\qquad$ $g_i^{-1}(y|\boldsymbol{x}) = g_i^L(y)$

4. Apply chain rule $\qquad\qquad$ $\boldsymbol{c} = f^L(y) = g_1^L \circ g_2^L \circ \cdots \circ g_n^L(y)$

# The Sound

Cat!

Still cat!

No cat 😓

How about using gradient as contribution?

$$c = \frac{\partial f(x)}{\partial x}$$

How about using gradient as contribution?

$$\boldsymbol{c} = \frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}}$$



$$\frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}_A} > \frac{\partial f(\boldsymbol{x})}{\partial \boldsymbol{x}_B}$$

# Sun

95%

Totally sun

# 80%

## Still sun

50%

Maybe not sun

😓

$$\frac{\partial f(\boldsymbol{x})}{\partial A} = 0$$

$A$

50%

😀

$$\frac{\partial f(\boldsymbol{x}^{50\%})}{\partial A} > 0$$

How about <span style="color:red">adding up gradients</span> at different percentages from reference?

$$c = \sum_{\alpha} \frac{\partial f(\boldsymbol{x}^{\alpha})}{\partial \boldsymbol{x}}$$

# The *Gradient* Methodology

$$c = (x - r) \cdot \int_{\alpha=0}^{1} \frac{\partial f(r + \alpha(x - r))}{\partial x} d\alpha$$

# Agenda

• Why feature attribution

• Classic methodologies—the general, the good and the sound

• **The SOTA and the limit**

# The SOTA and the Limit

- SOTA methods

- The general limit to feature attribution

- Remarks

| Shapley methodology | Inversion methodology | Gradient methodology |
|---|---|---|
| $SHAP^1$ | $DeepLIFT^2$ | Integrated Gradients[3] |
| For every function | For CNN | For deep networks |
| • Has variants with different niche<br>• Kernel SHAP → general functions<br>• Deep SHAP → CNN<br>• Tree SHAP → decision trees | • Tailored inverse for CNN, for which it is empirically powerful<br>• Not theoretically sound, no robustness guarantee for other model types | • Has soundness and robustness guarantee<br>• Powerful for deep differentiable functions<br>• Weak for non-differentiable functions, e.g., max-pooling<br>• Has a tunable parameter for tradeoff between speed and accuracy |

[1]
A Unified Approach to Interpreting Model Predictions.
https://arxiv.org/abs/1705.07874

[2]
Learning Important Features Through Propagating Activation Differences.
https://arxiv.org/abs/1704.02685

[3]
Axiomatic Attribution for Deep Networks.
https://arxiv.org/abs/1703.01365

$$y = f(\boldsymbol{x}) = x_1 + x_2 + x_3$$

$$c_1 = x_1$$

$$c_2 = x_2$$

$$c_3 = x_3$$

$$y = f(\boldsymbol{x}) = x_1 x_2 + x_3$$

$$c_1 = \ ?$$

$$c_2 = \ ?$$

$$c_3 = x_3$$

$$y = f(\boldsymbol{x}) = x_1 x_2 + x_3$$

$$c_1 = c_2 = \frac{x_1 x_2}{2} \quad 🤔$$

$$c_3 = x_3$$

$$y = f(\boldsymbol{x}) = {x_1}^{x_2} + x_3$$

$$c_1 = \;???$$

$$c_2 = \;???$$

$$c_3 = x_3$$

Various *sound, robust, powerful* methods exist for many popular model types.

A generally *correct* linear attribution doe not exist.

Large *generative* foundation models have started a new chapter.