Representations: Looking Inside Networks COMP6258



What are "learned representations"?

Learned Representations

- Informal: How the model encodes the input
- At a particular point in the network: collection of outputs for all possible inputs (often implied: belonging to your data distribution)
- In papers often people refer to feature maps of input images as the "learned" representations"







Why should we care about them?



Test Accuracy Is Not Everything (Reminder)



Test Accuracy Is Not Everything (Reminder)

e.g. Use distance in representational space to gauge performance on unseen samples?

Quick Recap: FID



 $\operatorname{FID}\left(\boldsymbol{\mu}_{r},\boldsymbol{\Sigma}_{r},\boldsymbol{\mu}_{g},\boldsymbol{\Sigma}_{g}\right)=\left\|\boldsymbol{\mu}_{r}-\boldsymbol{\mu}_{g}\right\|_{2}^{2}+\operatorname{Tr}\left(\boldsymbol{\Sigma}_{r}+\boldsymbol{\Sigma}_{g}-2\left(\boldsymbol{\Sigma}_{r}\boldsymbol{\Sigma}_{g}\right)^{\frac{1}{2}}\right),$

Quick Recap: FID



We need to understand the representational spaces

Let's look at some feature maps

Let's look at some feature maps (First feature map for a few different samples)

ResNet-18 trained on ImageNet Feature maps after the first convolutional layer



Reference feature map



A













Let's look at some feature maps (Second feature map for a few different samples)

Reference feature map





R

A







Now let's look at a later point in the network

Reference feature map





B

A









Reference feature map





B

A









Why is it hard to identify similarity? (MANY reasons, but let's give some examples)

functional

functional

structural

functional

structural

• Intuition: Using certain feature maps, how well can a model **perform**?





VS















Can this model perform well? If yes, representations A and B are "compatible"



Reminder: How can we find out if similar input patterns are encoded by a collection of feature maps? What does compatibility of embeddings tell us about input patterns represented?



Reminder

class A



class C



class B







Digit 2

VS





Digit 0










Platonic Representations?

How can we find out if similar input patterns are encoded by a collection of feature maps?

functional

structural

Most common method

- Most common method
- Statistical perspective

- Most common method
- Statistical perspective
- Intuition:



- Most common method
- Statistical perspective
- Intuition:



- Most common method
- Statistical perspective
- Intuition:





- Most common method
- Statistical perspective
- Intuition:









- Most common method
- Statistical perspective
- Intuition:





OR



• A and B are mapped to the X and Y space $X = A A^{\mathsf{T}}$





• A and B are mapped to the X and Y space $X = A A^{\mathsf{T}}$

M - number of samples





M x M tensors

• A and B are mapped to the X and Y space $X = A A^{\mathsf{T}}$

M x M tensors

M - number of samples

Similarity between each pair of samples according to all the feature maps



• A and B are mapped to the X and Y space $X = A A^{\mathsf{T}}$

• Hilbert-Schmidt Independence Criterion (HSIC):

HSIC(X, Y) =

$$= \frac{1}{(M-1)^2} \operatorname{tr}(X^*, Y^*)$$

- A and B are mapped to the X and Y space $X = A A^{\mathsf{T}}$
- Hilbert-Schmidt Independence Criterion (HSIC):
 - HSIC(X, Y) =
- Centered Kernel Alignement (CKA):
 - $\mathsf{CKA}(X,Y) =$ -

$$= \frac{1}{(M-1)^2} \operatorname{tr}(X^*, Y^*)$$

$$\frac{\mathsf{HSIC}(X,Y)}{\sqrt{\mathsf{HSIC}(X,X)\,\mathsf{HSIC}(Y,Y)}}$$

What did the functional perspective (try to) capture that CKA doesn't?

CKA take-away

 But not accounting for the functional use of the feature maps is a major drawback of structural similarity

Compression is problematic for both perspectives

Structural similarity is more nuanced and useful than functional similarity

Do these feature maps of the

Do these feature maps capture the same patterns in the input?















How can we find out what input patterns a feature map encodes?

Saliency Map Comparison GradCAM

Intent: What input regions are "important" when making a prediction?



Saliency Map Comparison GradCAM

Intent: What input regions are "important" when making a prediction?



• If I make a small change to the input, how much does it affect the prediction?



Saliency Map - CIFAR10 Example





-3.4393 -2.7674 -4.2603 6.4203 -0.5697 10.8899 -3.1707 0.5069 -1.9005 -1.7105

Saliency Map - CIFAR10 Example



-3.4393 -2.7674 -4.2603 6.4203 -0.5697 10.8899 -3.1707 0.5069 -1.9005 -1.7105

Saliency Map - CIFAR10 Example



-3.4393 -2.7674 -4.2603 6.4203 -0.5697 10.8899 -3.1707 0.5069 -1.9005 -1.7105

Saliency Map Comparison



Cutout

CAM for 'Poodle'





Yun et al. (2019). CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features

Saliency Map Comparison



CAM for

'Poodle'

Cutout

Yun et al. (2019). CutMix: Regularization Strategy to Train Strong Classifiers with Localizable Features









 $\lim_{h \to 0} \frac{f(x) - f(x+h)}{h}$

 $\lim_{h \to 0} \frac{f(x) - f(x+h)}{h}$

to the input

 $\lim_{h \to 0} \frac{f(x) - f(x+h)}{h}$

How much the output changes when I make an infinitesimally small change
Back to Basics: Derivative

 $\lim_{h\to 0}$

- to the input
- Robust feature prediction doesn't change when I make a perturbation

)
$$\frac{f(x) - f(x+h)}{h}$$

How much the output changes when I make an infinitesimally small change

Back to Basics: Derivative

 $\lim_{h\to 0}$

- to the input
- Robust feature prediction doesn't change when I make a perturbation
- Non-robust feature prediction is sensitive to perturbations

)
$$\frac{f(x) - f(x+h)}{h}$$

How much the output changes when I make an infinitesimally small change

Back to Basics: Derivative

 $\lim_{h\to 0}$

- How much the output changes whe to the input
- Importance or sensitivity?

)
$$\frac{f(x) - f(x+h)}{h}$$

• How much the output changes when I make an infinitesimally small change

Saliency Map Comparison

• Looking at sensitivity could be insightful



Saliency Map Comparison

Looking at sensitivity could be insightful

But is it only capturing a small part of what a model learns



Saliency Map Comparison

Looking at sensitivity could be insightful

But is it only capturing a small part of what a model learns

 Sensitive features can be important, but not all important features are sensitive





• We don't have good ways of interpreting what is going on inside DNNs

Take-aways

- existing techniques

 We don't have good ways of interpreting what is going on inside DNNs But we can use basic DL concepts to construct correct interpretations of

Take-aways

- existing techniques
- And hopefully build better ones

 We don't have good ways of interpreting what is going on inside DNNs But we can use basic DL concepts to construct correct interpretations of

Take-aways

- existing techniques
- And hopefully build better ones
- Understanding representations is very much work in progress

• We don't have good ways of interpreting what is going on inside DNNs But we can use basic DL concepts to construct correct interpretations of