

## Interpretable Machine Learning

An overview-flavored talk containing a subset of work from the last year

#### Been Kim

with a lot of awesome people inside and outside of Google:

Marten Wattenberg, Julius Adebayo Justin Gilmer, Carrie Cai, Fernanda Viegas, Rory Sayres, Mortiz Hardt, Ian Goodfellow











## ML community is responding



#### Year

## This is not a new problem. Why now?

Complexity and prevalence!





I heard you can just use decision trees...

#### Can we go home now?

http://www.ogroup.com.au/raise-your-hand-when-you-should-and-why-you-should/

# Experiment.

Data = [Sunny, 200]

#people

• I will show you a decision tree. Follow the right path given a data point, and you do:



# Experiment.

Clap!

• I will show you a decision tree. Follow the right path given a data point, and you do:





 As soon as you know the answer, do the action!









Weather =

And can you explain what the overall logic of the system was?

If I give you a lot of data points, Year can you guess which feature was most / 'important' (i.e., used in more number of examples)?

Clap!

Stomp

Right

Left



ree e here No

Left

Left

### Common misunderstanding: Decision trees and linear models are always interpretable.

#### Do we need a different model? How about rule lists?

If ( sunny and hot )	then	go swim
Else if ( sunny and cold )	then	go ski
Else	then	go work

#### Do we need a different model? How about rule lists?

If ( sunny and hot )
Else if ( sunny and cold )
Else if ( wet and weekday )
Else if ( free coffee )
Else if ( cloudy and hot )
Else if ( snowing )
Else if ( New Rick and Morty)
Else if ( paper deadline )
Else if ( hungry )
Else if ( tired )
Else if ( advisor might come )
Else if ( code running )
Else

then	go swim
then	go ski
then	go work
then	attend tutorial
then	go swim
then	go ski
then	watch TV
then	go work
then	go eat
then	watch TV
then	go work
then	watch TV
then	go work

### Maybe rule sets are better?

IF ( sunny and hot ) OR ( cloudy and hot ) OR ( sunny and thirsty and bored ) THEN go to beach ELSE work

### Maybe rule sets are better?

IF (sunny and hot) OR (cloudy and hot) OR (sunny and thirsty and bored) OR (bored and tired) OR (thirty and tired) OR (code running) OR (friends away and bored) OR (sunny and want to swim) OR (sunny and friends visiting) OR (need exercise) OR (want to build castles) OR (sunny and bored) OR (done with deadline and hot) OR ( need vitamin D and sunny ) OR (just feel like it) THEN go to beach **ELSE** work

Are you saying decision trees, rule lists and rule sets don't work?!



Decision trees, rule lists or rule sets may work for your case!

The point here is that there is no one-size-fits-all method.

http://blog.xfree.hu/myblog.tvn?SID=&from=20&pid=&pev=2016&pho=02&pnap=&kat=1083&searchkey=&hol=&n=sarkadykati

## Is interpretability possible at all?

#### # WIRED

Our Machines Now Have Knowledge We'll Never Understand

SUBSCRIBE

#### DAVID WEINBERGER BACKCHANNEL 04.18.17 08:22 PM

#### OUR MACHINES NOW HAVE KNOWLEDGE WE'LL NEVER UNDERSTAND

#### SHARE



The new availability of huge amounts of data, along with the statistical tools to crunch these numbers, offers a whole new way of understanding the world. Correlation supersedes causation, and science can advance even without coherent models, unified theories, or really any mechanistic explanation at all.



TWEET

So wrote Wired's Chris Anderson in 2008. It kicked up a

## Is interpretability possible at all?

WIRED

TWEET

Our Machines Now Have Knowledge We'll Never Understand

SUBSCRIBE

DAVID WEINBERGER BACKCHANNEL 04.18.17 08:22 PM

OIR MACHINES NOW HAVE ENOWLEDCE WF'LL Common misunderstanding: We need to understand every single thing about the model.

and science can advance even without coherent models, unified theories, or really any mechanistic explanation at all.

Key Point:

Interpretability is NOT about understanding all bits and bytes of the model for all data points.

It is about knowing enough for your goals/downstream tasks.

#### My goal

# interpretability

To use machine learning **responsibly** we need to ensure that 1. our **values** are aligned 2. our **knowledge** is reflected



Fundamental **underspecification** in the problem



#### Fundamental **underspecification** in the problem



example 2: Science



#### Fundamental **underspecification** in the problem



\_\_example 2: Science



#### Fundamental **underspecification** in the problem

example3: mismatched objectives





example 2: Science

Fundamental **underspecification** in the problem

Common misunderstanding: More data or more clever algorithm will solve interpretability.

29 drugs.com

# What is NOT underspecification?



# When we may **not** want interpretability

• No significant consequences or when predictions are all you need.

• Sufficiently well-studied problem

 Prevent gaming the system mismatched objectives.

https://cdn.theatlantic.com/assets/media/img/mt/2015/04/shutterstock\_11926084/lead\_large.jpg https://www.jal.com/assets/img/flight/safety/equipment/pic\_tcas\_001\_en.jpg

http://www.cinemablend.com/pop/Netflix-Using-Amazon-Cloud-Explore-Artificial-Intelligence-Movie-Recommendations-62248.html



# When we may **not** want interpretability

Prevent gaming the We always need interpretability.

Climb!

Common misunderstanding:

Descend

• No significant consequences or when predictions are all you need.

• Sufficiently well-studied problem

https://cdn.theatlantic.com/assets/media/img/mt/2015/04/shutterstock\_11926084/lead\_large.jpg

mismatched objectives

https://www.jal.com/assets/img/flight/safety/equipment/pic\_tcas\_001\_en.jpg

http://www.cinemablend.com/pop/Netflix-Using-Amazon-Cloud-Explore-Artificial-Intelligence-Movie-Recommer Agetions-62248.html



fairness accountability trust causality etc.



## Our cousins are not us



- Interpretability can help with them when we cannot formalize these ideas
- But once formalized, you may not need interpretability.

#### Let's build some interpretability methods.



www.memecenter.com
ingredients for interpretability methods.

# $\underset{E}{\operatorname{argmax}} Q(E|?)$







Data

X Class1



Data

X Class1







Post-training interpretability methods



Building inherently interpretable model





Post-training interpretability methods

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{E}_{x} planation | \mathbf{M} odel, \mathbf{H} uman, \mathbf{D} ata, \mathbf{T} ask)$ 

Building inherently interpretable model

 $\underset{E,M}{\operatorname{argmax}} Q(\operatorname{Explanation}, \operatorname{Model}| \operatorname{Human}, \operatorname{Data}, \operatorname{Task})$ 

Explaining data

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{Explanation} | \mathbf{H}\mathsf{uman}, \mathbf{D}\mathsf{ata}, \mathbf{T}\mathsf{ask})$ 

Post-training interpretability methods

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{E}_{x} planation | \mathbf{M} odel, \mathbf{H} uman, \mathbf{D} ata, \mathbf{T} ask)$ 

Building inherently interpretable model

 $\underset{E,M}{\operatorname{argmax}} Q(\mathbf{E} x planation, \mathbf{M} odel | \mathbf{H} uman, \mathbf{D} ata, \mathbf{T} ask)$ 

Explaining data

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{E} x p | \mathbf{a} nation | \mathbf{H} u m a n, \mathbf{D} a t a, \mathbf{T} a s k)$ 

Mv<sup>°</sup>ML

Post-training interpretability methods

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{E}_{x} planation | \mathbf{M} odel, \mathbf{H} uman, \mathbf{D} ata, \mathbf{T} ask)$ 

Building inherently interpretable model rule, feature, example-based and many more.  $argmax_{E,M} Q(\mathbf{E}xplanation, \mathbf{M}odel|\mathbf{H}uman, \mathbf{D}ata, \mathbf{T}ask)$ 

Explaining data

data visualization, exploratory data analysis

 $\underset{E}{\operatorname{argmax}} Q(\mathbf{Explanation} | \mathbf{H}\mathsf{uman}, \mathbf{D}\mathsf{ata}, \mathbf{T}\mathsf{ask})$ 

My<sup>°</sup>ML







X Class1







# Problem: Post-training explanation argmax Q(Explanation|Model, Human, Data, Task)



One of the most popular interpretability methods for images:

# Saliency maps



Used for image classification and medical applications.

a logit 
$$\rightarrow \frac{\partial p(z)}{\partial x_{i,j}}$$
  
pixel i,j  $\rightarrow \frac{\partial x_{i,j}}{\partial x_{i,j}}$ 



 $\underset{E}{\operatorname{argmax}} Q(E|M,H,D,T)$ 

picture credit: @sayres SmoothGrad [Smilkov, Thorat, K., Viégas, Wattenberg '17] Integrated gradient [Sundararajan, Taly, Yan '17] 55 One of the most popular interpretability methods for images:

# Saliency maps



One of the most popular interpretability methods for images:

# Saliency maps



SmoothGrad [Smilkov, Thorat, K., Viégas, Wattenberg '17] Integrated gradient [Sundararajan, Taly, Yan '17] 57

$$\underset{E}{\operatorname{argmax}} Q(E|M,H,D,T)$$

Used for image classification

and medical applications.

a logit  $\rightarrow \frac{\partial p(z)}{\partial x_{i,j}}$ pixel i,j  $\rightarrow \frac{\partial x_{i,j}}{\partial x_{i,j}}$ 

Sanity check: If I change M a lot, will human perceive that E has changed a lot?

### Some confusing behaviors of saliency maps.



### Some confusing behaviors of saliency maps.





Sanity Checks for Saliency Maps Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

### Some confusing behaviors of saliency maps.





Sanity Checks for Saliency Maps Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

# Some saliency maps look similar when we randomize the network.



Sanity Checks for Saliency Maps Joint work with Adebayo, Gilmer, Goodfellow, Hardt, [NIPS 18]

### What can we learn from this?

- Potential human confirmation bias: Just because it "makes sense" to humans, doesn't mean they reflect evidence for the prediction.
- Our discovery is consistent with other findings [Nie, Zhang, Patel '18] [Ulyanov, Vedaldi, Lempitsky '18]
- Some of these methods have been shown to be useful in practice. Explaining predictions or features? More studies needed.



### What can we do better? Creating a wishlist.



### What can we do better? Creating a wishlist.







prediction: Cash machine

prediction: Sliding door Why correct? Why incorrect?

<u>https://pair-code.github.io/saliency/</u> SmoothGrad [Smilkov, Thorat, K., Viégas, Wattenberg '17]



Were there more pixels on the cash machine than on the person?

Did the 'human' concept matter? Did the 'glasses' or 'paper' matter?







Were there more pixels on the cash machine than on the person?

Did the 'human' concept matter? Did the 'glasses' or 'paper' matter?

Which concept mattered more?

Is this true for all other cash machine predictions?

Wouldn't it be great if we can **quantitatively** measure how important *any* of these **user-chosen concepts** are?

### Goal of TCAV: Testing with Concept Activation Vectors



**Quantitative** explanation: how much a concept (e.g., gender, race) was important for a prediction in a trained model.

...even if the concept was not part of the training.

### Goal of TCAV: Testing with Concept Activation Vectors



A trained machine learning model (e.g., neural network)

p(z)

**Doctor-ness**


A trained machine learning model (e.g., neural network)

p(z)

**Doctor-ness** 

Was gender concept important to this doctor image classifier?







### TCAV:

### Testing with Concept Activation Vectors



### Defining concept activation vector (CAV)



# Defining concept activation vector (CAV)

#### Inputs:



### TCAV:

### Testing with Concept Activation Vectors



### TCAV core idea: Derivative with CAV to get prediction sensitivity

#### TCAV



**Directional derivative with CAV** 

### TCAV core idea: Derivative with CAV to get prediction sensitivity

#### TCAV



$$S_{C,k,l}(\mathbf{M})$$

$$S_{C,k,l}(\mathbf{M})$$

$$S_{C,k,l}(\mathbf{M})$$

$$S_{C,k,l}(\mathbf{M})$$

$$\text{TCAV}_{Q_{C,k,l}} = \frac{|\{ x \in X_k : S_{C,k,l}(x) > 0 \}|}{|X_k|}$$

**Directional derivative with CAV** 

### TCAV:

### Testing with Concept Activation Vectors



### TCAV:

### Testing with Concept Activation Vectors



# Guarding against spurious CAV

Did my CAVs returned high sensitivity by chance?

# Guarding against spurious CAV





Learn many stripes CAVs using different sets of random images

Guarding against spurious CAV





Guarding against spurious CAV







Guarding against spurious CAV







Check the distribution of  $TCAV_{Q_{C,k,l}}$  is statistically different from random using t-test



TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it. Even if your training data wasn't tagged with the concept

Even if your input feature did not include the concept

1. Learning CAVs	2. Getting TCAV score	3. CAV validation
$f_{l}(\textcircled{)}) f_{l}(\textcircled{)}) f_{l}(f_{l}) f_{l}(f_{l})) f_{l}(f_{l}) f_{l}(f_{l})) f_{l}(f_{l}) f_{l}(f_{l}) f_{l}(f_{l})) f_{l}(f_{l}) f_{l}$	$S_{C,k,l}(\mathbb{Z})$ $S_{C,k,l}(\mathbb{Z})$ ) TCAV <sub>QC,k,l</sub>	Qualitative
$f_{l}(\blacksquare) \qquad f_{l}(\textcircled{)} \qquad f_{l}(f_{l}(f_{l})) \qquad f_{l}(f_{l}(f_{l})) \qquad f_{l}(f_{l}) \qquad f_{l}) \qquad f_{l}(f_{l}) \qquad f_{l}(f_{l$	$S_{C,k,l}(\mathbb{R})$	Quantitative

# Results

- - cab image

cab image with caption

2. Biases in Inception V3 and GoogleNet

1. Sanity check experiment

3. Domain expert confirmation from Diabetic Retinopathy



# Results



cab image

cab image with caption

- 1. Sanity check experiment
- 2. Biases from Inception V3 and GoogleNet
- 3. Domain expert confirmation from Diabetic Retinopathy



# Sanity check experiment

If we know the ground truth (important concepts), will TCAV match?



#### An image + Potentially noisy Caption



models can use either image or caption concept for classification.





concept for

classification.



Caption noise level in training set

# Sanity check experiment



#### Cool, cool. Can saliency maps do this too?

# Can saliency maps communicate the same information?



### Human subject experiment: Can saliency maps communicate the same information?



- 50 turkers are
  - asked to judge importance of image vs. caption given saliency maps.
  - asked to indicate their confidence
  - shown 3 classes (cab, zebra, cucumber) x 2 saliency maps for one model

### Human subject experiment: Can saliency maps communicate the same information?

- Random chance: 50%
- Human performance with saliency map: 52%
- Humans can't agree: more than 50% no significant consensus



Human subject experiment: Can saliency maps communicate the same information?

- Random chance: 50%
- Human performance with saliency map: 52%
- Humans can't agree: more than 50% no significant consensus
- Humans are **very** confident even when they are wrong.



# Results

1. Sanity check experiment



cab image

cab image with caption

- 2. Biases from Inception V3 and GoogleNet
- 3. Domain expert confirmation from Diabetic Retinopathy



#### TCAV in

### Two widely used image prediction models



105

#### TCAV in

### Two widely used image prediction models



### TCAV in Two widely used image prediction models



107

#### TCAV in

### Two widely used image prediction models


# Results

cab

cab image

cab image with caption

1. Sanity check experiment

- 2. Biases Inception V3 and GoogleNet
- 3. Domain expert confirmation from Diabetic Retinopathy



# Diabetic Retinopathy

- Treatable but sight-threatening conditions
- Have model to with accurate prediction of DR (85%) [Krause et al., 2017]

Concepts the ML model uses

Vs

Diagnostic Concepts human doctors use

DR level 4 Retina



## TCAV for Diabetic Retinopathy



#### TCAV for Diabetic Retinopathy



TCAV for DR level 1

**HMA** 

MA



TCAV shows the model is **inconsistent** with doctor's knowledge for classes when model is less accurate

DR level 1

Green: domain expert's label on concepts belong to the level Red: domain expert's label on concepts does not belong to the level

### TCAV for Diabetic Retinopathy



**Green:** domain expert's label on concepts belong to the level **Red:** domain expert's label on concepts does not belong to the level

#### Summary:

code: github.com/tensorflow/tcav

#### Testing with Concept Activation Vectors

Joint work with Wattenberg, Gilmer, Cai, Wexler, Viegas, Sayres

**stripes** concept (score: 0.9) was important to **zebra** class for this trained network.

Our values

TCAV provides **quantitative importance** of a concept **if and only if** your network learned about it.



Our knowledge



https://imgflip.com